



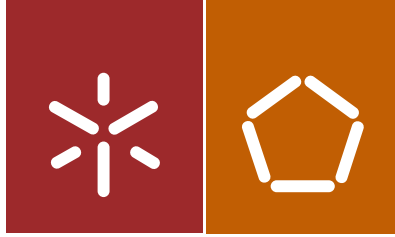
Benedita Chaves Moura

Inteligência Coletiva para Análise de  
Sentimento sobre Mensagens da  
Plataformas StockTwits

Universidade do Minho  
Escola de Engenharia







Universidade do Minho  
Escola de Engenharia

Benedita Chaves Moura

Inteligência Coletiva para Análise de  
Sentimento sobre Mensagens da  
Plataformas StockTwits

Dissertação de Mestrado  
Engenharia e Gestão de Sistemas de Informação

Trabalho efetuado sob a orientação do  
Professor Doutor Paulo Alexandre Ribeiro Cortez

Outubro de 2014

## DECLARAÇÃO

Nome: Benedita Chaves Moura

Endereço eletrónico: a54077@alunos.uminho.pt Telefone:936688339

Cartão do Cidadão: 13594548

Título da dissertação: Inteligência Coletiva para Análise de Sentimento sobre Mensagens da Plataforma StockTwits

Orientador:

Professor Doutor Paulo Alexandre Ribeiro Cortez

Ano de conclusão: 2014

Mestrado Integrado em Engenharia e Gestão de Sistemas de Informação

É AUTORIZADA A REPRODUÇÃO INTEGRAL DESTA DISSERTAÇÃO APENAS PARA EFEITOS DE INVESTIGAÇÃO, MEDIANTE DECLARAÇÃO ESCRITA DO INTERESSADO, QUE A TAL SE COMPROMETE.

Universidade do Minho, \_\_\_\_/\_\_\_\_/\_\_\_\_\_

Assinatura:

## Agradecimentos

Dado por terminado este projeto, importa agradecer aos elementos que contribuíram para o desenvolvimento do mesmo. Em primeiro agradeço ao Professor Paulo Cortez, pela orientação e apoio demonstrado ao longo de todo o projeto.

Ao professor Nelson Areal pela sua contribuição através da partilha da aplicação.

Ao Nuno Oliveira por toda a disponibilidade demonstrada para obter os dados e realização dos testes sobre os mesmos.

Aos meus pais e irmã pelo apoio incondicional e motivação demonstrada ao longo de todo o projeto.

Ao Tiago em especial pela motivação e apoio em todos os momentos.

À Marta e Nadine pelo apoio e companheirismo demonstrado.

A todos eles gostaria de agradecer pelo apoio pois sem eles este projeto não teria sido possível.



## Resumo

Hoje em dia a *Web 2.0* possibilita aos utilizadores inúmeras possibilidades de colaboração para atingir-se um determinado objetivo, realçando-se em particular a Inteligência Coletiva. Por sua vez, a Análise de Sentimentos relaciona-se com a identificação das opiniões positivas ou negativas relativamente a um texto sendo que existem já diversas formas de obter a polaridade dos sentimentos através de algoritmos e dicionários de léxicos.

Neste projeto de dissertação avaliou-se a capacidade de utilizar a Inteligência Coletiva para a Análise de Sentimento sobre mensagens de âmbito financeiro do serviço de *microblogging* *StockTwits*. Em particular, desenvolveu-se uma aplicação (Finance\$entiment) que foi disponibilizada na *Web* e que permitiu classificar o sentimento de um conjunto de mensagens *StockTwits* via uma abordagem de inteligência coletiva. Por último, a assertividade da inteligência coletiva foi comparada com a obtida com métodos automáticos, via dicionários de léxicos e algoritmos de *text mining*.





## Abstract

With the Web 2.0, there are numerous ways for users, to collaborate in order to reach a certain goal. In this work, we particularly highlight the use of Collective Intelligence. On the other hand, Sentiment Analysis relates to the identification of positive and negative opinions in a given text and there are already several ways of automatically obtaining sentiment value through text mining algorithms and lexical dictionaries.

In this dissertation project, we evaluate the ability of using Collective Intelligence to perform Sentiment Analysis of financial messages from *StockTwits* microblog. For such purpose, we developed an application (*Finance\$entiment*), and made it available on the Web so it would be used to classify the sentiment values from a set of messages from *StockTwits*, through a Collective Intelligence approach. Finally, the performance of the Collective Intelligence approach was compared with the one achieved with automated methods through the use of lexicon dictionaries and text mining algorithms.



# Índice

Agradecimentos .....	i
Resumo .....	iii
Abstract .....	v
Índice.....	vii
Índice de Figuras.....	xi
Índice de Tabelas .....	xvii
1 Introdução .....	1
1.1 Enquadramento .....	1
1.2 Objetivos.....	1
1.3 Organização.....	2
2 Análise de Sentimento.....	3
2.1 Introdução .....	3
2.2 Motivação para Análise de Sentimento.....	3
2.3 Aplicações da Análise de Sentimento.....	4
2.4 Desafios da Análise de Sentimento .....	5
2.5 Funcionamento da Análise de Sentimento.....	6
2.5.1 Classificação e Extração .....	8
2.5.2 Sumarização .....	8
2.6 Níveis de Análise de Sentimento .....	9
2.6.1 Classificação dos Sentimentos ao Nível de Documentos.....	9
2.6.2 Classificação dos Sentimentos ao Nível das Frases .....	14
2.7 Sumarização das Opiniões .....	17
2.7.1 Sumarização de Opiniões Baseadas em Aspetos.....	17
2.7.2 Sumarização Tradicional.....	17

2.7.3 Sumarização Orientada à Opinião num Documento Único .....	18
2.8 Dicionários de Léxicos .....	18
2.8.1 <i>SentiWordNet</i> .....	21
2.8.2 <i>Financial Sentiment Dictionaries (FIN)</i> .....	21
3 Inteligência Coletiva e <i>Crowdsourcing</i> .....	23
3.1 Introdução .....	23
3.2 Inteligência Coletiva.....	24
3.2.1 Definição.....	24
3.2.2 Classificação .....	25
3.2.3 Exemplos .....	26
3.2.4 Inteligência Coletiva nas Aplicações <i>Web</i> .....	27
3.2.5 Técnicas .....	29
3.2.6 Áreas de Aplicação .....	30
3.2.7 Vantagens .....	34
3.2.8 Utilizadores .....	35
3.2.9 Atribuição de <i>Tags</i> .....	42
3.2.10 Tipos de Conteúdo .....	46
3.2.11 Blogosfera e Web Crawling .....	50
3.2.12 Fatores-Chave e de Sucesso .....	52
3.3 <i>Crowdsourcing</i> .....	54
3.3.1 Definição.....	54
3.3.2 Exemplos .....	55
3.3.3 <i>Crowdsourcing</i> e Opensource .....	58
3.3.4 Desafios Principais .....	59
3.3.5 Vantagens e Desvantagens .....	60

3.3.6 Trabalho Futuro.....	61
3.4 CAPTCHA.....	62
3.4.1 Definição.....	62
3.4.2 Funcionamento .....	63
3.4.3 Aplicações.....	63
3.4.5 Caraterísticas .....	65
3.4.6 Problemas e Vulnerabilidades .....	65
3.4.7 Tipos de <i>CAPTCHAs</i> .....	66
3.4.8 <i>reCAPTCHA</i> .....	69
4 Inteligência Coletiva para Análise de Sentimento.....	71
4.1 Introdução .....	71
4.2 Planeamento do Projeto .....	71
4.3 Metodologia .....	73
4.4 Aplicação Desenvolvida .....	75
4.4.1 Arquitetura do Sistema .....	75
4.4.2 Ferramentas Utilizadas.....	76
4.5 Implementação .....	77
4.5.1 Descrição da Aplicação.....	77
4.5.2 Divulgação .....	96
4.6 Análise de Resultados .....	97
4.6.1 Análise dos Utilizadores.....	97
4.6.2 Análise das Classificações Individuais .....	111
4.6.3 Análise das Classificações Médias .....	113
4.6.4 Comparação dos Resultados com Abordagens Automáticas .....	118
4.6.5 Sumário .....	119

5 Conclusões .....	123
5.1 Síntese.....	123
5.2 Discussão .....	124
5.3 Trabalho Futuro.....	125
Referências Bibliográficas.....	127
Anexos.....	135

# Índice de Figuras

Figura 1 - Processo de extração de Sentimento de <i>Tweets</i> através da utilização de léxicos, adaptado de (Zhang et. al, 2011).....	20
Figura 2 - Classificação da Inteligência Coletiva, adaptado de (Alag, 2009).....	25
Figura 3 – Influência dos utilizadores e da Inteligência Coletiva, adaptado de (Alag, 2009) 28	
Figura 4 - Inteligência Coletiva nas aplicações <i>Web</i> (1- Permite interação dos utilizadores; 2- Aprende sobre os utilizadores em conjunto; 3 – Personaliza conteúdos tendo em conta dados das interações dos utilizadores); adaptado de (Alag, 2009).....	28
Figura 5 - Três fatores que influenciam o número mínimo de participantes, adaptado de (Lesser et. al,2012).....	40
Figura 6 - Técnicas de Motivação, adaptado de (Lesser et. al, 2012).....	41
Figura 7 - Tag Cloud, adaptado de (Alag, 2009).....	45
Figura 8 - Utilização de Tags, contexto e meta-dados para obter os resultados relevantes, adaptado de (Alag, 2009).....	46
Figura 9 – Alguns exemplos de <i>CAPTCHAs</i> baseados em Texto (Pawar & Bauskar, 2013) 67	
Figura 10 - Exemplos de <i>CAPTCHAs</i> gráficos (Asirra e ESP CAPTCHA) (Pawar & Bauskar, 2013).....	68
Figura 11 - Exemplo de CAPTCHA áudio (Pawar & Bauskar, 2013).....	68
Figura 12 - Exemplo de um <i>reCAPTCHA</i> (Vaishakh & Harish, 2011).....	70
Figura 13- Planeamento do Projeto.....	72
Figura 14 - Gráfico de Gantt do Planeamento do Projeto.....	73
Figura 15 - Arquitetura da aplicação <i>Finance\$entiment</i> .....	76
Figura 16 - Página inicial da aplicação <i>Finance\$entiment</i> .....	78
Figura 17 - Aviso de <i>email</i> inválido.....	78
Figura 18 - Política de Privacidade.....	79
Figura 19 - Informação relativa ao âmbito do projeto.....	80
Figura 20 – Primeiro exemplo de como se joga.....	80
Figura 21 – Segundo exemplo de como se joga.....	81
Figura 22 - 3º Exemplo de como se joga.....	81

Figura 23 - Página de jogo.....	82
Figura 24 - Página de jogo.....	82
Figura 25 - Página de jogo.....	83
Figura 26 - Página de ajuda.....	84
Figura 27 - Página de explicação da pontuação.....	84
Figura 28 - Aviso de nova medalha .....	85
Figura 29 - Visualização das medalhas ganhas .....	85
Figura 30 - Visualização das medalhas ganhas .....	86
Figura 31 - Classificações do utilizador .....	86
Figura 32 - <i>Ranking</i> dos utilizadores .....	87
Figura 33- Página de <i>login</i> para aceder à secção do Administrador .....	88
Figura 34 – Página onde se encontram todas as respostas dadas pelos utilizadores.....	89
Figura 35 – Visualização das classificações e classificação média de determinada frase... 89	
Figura 36 – Exemplos de respostas cujo sentimento atribuído foi “Difícil de Classificar” ... 90	
Figura 37 - Exemplos de respostas cujo sentimento atribuído foi “Neutro” .....	90
Figura 38 - Exemplos de respostas cujo sentimento atribuído foi “Negativo” .....	91
Figura 39 - Exemplos de respostas cujo sentimento atribuído foi “Positivo” .....	91
Figura 40 – Estatísticas gerais da aplicação.....	92
Figura 41 – Estatísticas das classificações individuais.....	92
Figura 42 – Estatísticas das frases tendo em conta a média das classificações .....	93
Figura 43 - Casos de Uso (UML) da aplicação <i>Finance\$entiment</i> .....	94
Figura 44 - Estrutura da Base de Dados da aplicação <i>Finance\$entiment</i> .....	95
Figura 45- Relação entre número de Utilizadores e número de classificações .....	97
Figura 46 -Utilizadores com menos de 10 respostas (Gráfico I) .....	98
Figura 47 - Utilizadores com menos de 10 respostas (Gráfico II) .....	98
Figura 48 - Número de Utilizadores por % de Respostas Certas/Erradas.....	99
Figura 49 - Utilizadores por % de Respostas Certas .....	100
Figura 50 - Utilizadores por % de Respostas Erradas .....	100
Figura 51 - Utilizadores com mais de 10 e menos de 50 respostas (Gráfico I) .....	101
Figura 52 - Utilizadores com mais de 10 e menos de 50 respostas (Gráfico II) .....	102
Figura 53 - Número de Utilizadores por % de Respostas Certas/Erradas.....	102



Figura 54 - Utilizadores por % de Respostas Certas .....	103
Figura 55 - Utilizadores por % de Respostas Erradas .....	103
Figura 56 - Utilizadores com mais de 50 e menos de 100 respostas Utilizadores com mais de 50 e menos de 100 respostas .....	104
Figura 57 - Número de Utilizadores por % de Respostas Certas/Erradas.....	105
Figura 58 - Utilizadores por % de Respostas Certas .....	105
Figura 59 - Utilizadores por % de Respostas Erradas .....	106
Figura 60 - Utilizadores com mais de 100 respostas .....	107
Figura 61 - Número de Utilizadores por % de Respostas Certas/Erradas.....	107
Figura 62 - Utilizadores por % de Respostas Certas .....	108
Figura 63 - Utilizadores por % de Respostas Erradas .....	108
Figura 64 - Número de Utilizadores por % de Respostas Certas/Erradas.....	109
Figura 65 - Utilizadores por % de Respostas Certas .....	110
Figura 66 - Utilizadores por % de Respostas Erradas .....	110
Figura 67 - Classificações atribuídas às Frases Positivas .....	112
Figura 68 - Classificações atribuídas às Frases Negativas.....	113
Figura 69 - Classificação Média das Frases Positivas .....	115
Figura 70 - Classificação Média das Frases Negativas .....	116
Figura 71 - Comparação das Classificações Individuais e da Média (Frases Positivas).....	117
Figura 72 - Comparação das Classificações Individuais e da Média (Frases Negativas) ...	117
Figura 73 - Comparação de Resultados da Média das Classificações, do algoritmo e dos léxicos FIN e SWN .....	118



## Índice de Tabelas

Tabela 1 - Vieses, mitigações e exemplos relativos à construção de soluções potenciais (Bonabeau, 2009) .....	31
Tabela 2 - Vieses, mitigações e exemplos para a avaliação de soluções potenciais (Bonabeau, 2009) .....	32
Tabela 3 - Tipos de Conteúdo, adaptado de (Alag, 2009) .....	47
Tabela 4 - Métricas e indicadores de sucesso em diferentes aplicações (Bonabeau, 2009)	53
Tabela 5 - Exemplos de Sistemas de <i>Crowdsourcing</i> (Doan et. al, 2011) .....	57
Tabela 6 - Classificações individuais .....	111
Tabela 7 - Classificações atribuídas às frases Positivas .....	111
Tabela 8 - Classificações atribuídas às Frases Negativas .....	112
Tabela 9 - Classificações Médias .....	114
Tabela 10 - Classificação Média das Frases Positivas .....	114
Tabela 11 - Classificação Média das Frases Negativas .....	115
Tabela 12 - Comparação de Resultados da Média das Classificações, do algoritmo e dos léxicos FIN e SWN .....	118



# 1 Introdução

## 1.1 Enquadramento

Atualmente a sociedade está a sofrer uma revolução cultural devido à evolução da *Internet* rumo à *Web 2.0* e devido ao facto da *Internet* estar mais próxima de todos nós e a um custo mais reduzido. Através da *Internet* as pessoas têm a possibilidade de comunicar entre si o que provoca uma crescente criação de *blogs*, fóruns e redes sociais entre outros. Desta forma as pessoas podem desfrutar da liberdade de expressão e do fácil acesso a todos os tipos de informação para expressar as suas opiniões relativas a diversos assuntos (Moreo et. al, 2012).

A *Web* permite a colaboração dos seus utilizadores para atingir determinado objetivo comum, sendo que o paradigma da chamada *Web 2.0* está fortemente vocacionado para a colaboração. Em particular, a *Web 2.0* facilita a realização da Inteligência Coletiva (Michalsky et. al, 2010). Desta forma, as organizações podem recorrer mais este método para resolver desafios complexos do mundo atual (Lesser et. al, 2012). O termo *Crowdsourcing* está também ligado à Inteligência Coletiva e fornece uma nova forma de resolver tarefas comuns e complexas, permitindo um número de soluções elevado tanto a nível quantitativo como qualitativo (Brabham, 2008). De entre as diversas formas de Inteligência Coletiva sobressai o projeto *reCAPTCHA* que permite não só a deteção de programas automáticos intrusos ao sistema (e.g., *Bot*) como reaproveita o esforço realizado pelos utilizadores humanos ao mapear digitalmente os caracteres encontrados em imagens retiradas de documentos (Schlaikjer, 2007).

## 1.2 Objetivos

O objetivo desta dissertação é realizar um estudo sobre a aplicação da Inteligência Coletiva para a Análise de Sentimento de mensagens (*tweets*) financeiras presentes no *microblog StockTwits*. Para alcançar este objetivo, desenvolveu-se um protótipo, a aplicação *Finance\$entiment*, que foi disponibilizada num servidor *Web* e divulgada a um conjunto de utilizadores humanos. Posteriormente, os resultados obtidos para a classificação do sentimento foram comparados com os obtidos via métodos automáticos, nomeadamente através do uso de dois dicionários de léxicos e um algoritmo simples de *text mining*.

### 1.3 Organização

Este documento está organizado por 5 capítulos.

O primeiro capítulo apresenta o contexto, objetivos e a organização desta dissertação.

No segundo capítulo, denominado Análise de Sentimento, são expostos conteúdos como a motivação, as aplicações, desafios e funcionamento da análise de sentimento. São descritos ainda dois dicionários de léxicos usados para analisar os resultados.

O terceiro capítulo descreve a Inteligência Coletiva, e o *Crowdsourcing*, com um particular destaque para o método *CAPTCHA*. Relativamente à Inteligência Coletiva, é realizada uma pequena definição da mesma, são descritos os tipos sob os quais esta pode surgir, são apresentados alguns exemplos e como pode ser utilizada nas aplicações *Web*. Para além disso são descritas as técnicas que pode utilizar, as áreas em que pode ser utilizada e as suas principais vantagens. De uma forma geral são descritas as formas de obter informações dos utilizadores e das suas interações, e para além disso são descritas formas de atribuir *Tags*, os tipos de conteúdo que se pode encontrar, a Blogosfera e o *Web Crawling* assim como os fatores-chave para aplicar a inteligência coletiva com sucesso. Relativamente ao subcapítulo *Crowdsourcing*, é realizada uma definição do mesmo, são dados alguns exemplos e é explicada a diferença entre *Crowdsourcing* e *Open Source*. Mais ainda, são descritos os seus desafios, vantagens e desvantagens e trabalho futuro desta área. Finalmente na secção *CAPTCHA*, é dada uma definição da mesma e do seu funcionamento, são descritas as suas possíveis aplicações, características, problemas e vulnerabilidades, assim como os tipos de *CAPTCHAs* existentes e o exemplo mais conhecido, o *reCAPTCHA*.

O quarto capítulo, intitulado de Inteligência Coletiva para análise de Sentimento, descreve a implementação da aplicação realizada, desde a escolha das ferramentas necessárias, arquitetura, desenvolvimento, funcionalidades, método de divulgação, obtenção dos dados e análise de resultados.

O quinto capítulo refere-se às conclusões resultantes da análise cuidadosa dos resultados obtidos. É realizada também uma síntese do projeto e é feita uma pequena análise de trabalho que possa vir a ser desenvolvido.

## 2 Análise de Sentimento

### 2.1 Introdução

Neste momento estamos a sofrer uma revolução cultural devido à evolução da *Internet* (e.g. *Web 2.0* e globalização da *Internet*). Atualmente as pessoas têm inúmeras possibilidades de comunicação assistida por tecnologias de informação, o que provocou a crescente criação de blogs, fóruns e redes sociais entre outros. Desta forma as pessoas podem desfrutar de uma maior liberdade de expressão bem como de um fácil acesso a todos os tipos de informação e assim expressar as suas opiniões relativas a diversos assuntos (Moreo et. al, 2012).

A **Análise de Sentimento**<sup>1</sup> é também designada por *Opinion Mining* e pode ser definida como o campo de estudo que analisa os sentimentos, opiniões, avaliações e atitudes que as pessoas apresentam, relativas a entidades tais como produtos, serviços, organizações, indivíduos, eventos, temas, entre outros (Liu, 2012). O foco da Análise de Sentimentos é identificar a forma como os sentimentos são expressos em textos. Estas expressões devem indicar opiniões positivas ou negativas relativamente a algo. Análise de Sentimento compreende as fases de identificação do sentimento expresso nos textos, a identificação da sua polaridade e a relação que estas compreendem com o assunto (Nasukawa & Yi, 2003). Em suma, a Análise de Sentimento preocupa-se com a análise de textos baseados no sentido, ou seja texto que possui opiniões ou emoções (Xu, 2012).

### 2.2 Motivação para Análise de Sentimento

Desde sempre que “aquilo que os outros pensam” é importante para o processo de tomada de decisão nas empresas. Neste momento a *Internet* tornou possível a obtenção da opinião e experiências de pessoas que não estão diretamente relacionadas com o desenvolvimento de

---

<sup>1</sup> Tradução adotada para o termo *Sentiment Analysis*.

determinado produto ou serviço. As pessoas estão cada vez mais a partilhar as suas opiniões com estranhos através da *Internet*. Segundo o estudo de Pang & Lee (2008): 81% dos utilizadores da *Internet* já fizeram pesquisas sobre determinado produto pelo menos uma vez; de 73% a 87% dos leitores de revisões *on-line* de restaurantes, hotéis, entre outros serviços, dizem que as revisões têm uma grande influência na sua compra; alguns consumidores estão dispostos a pagar de 20% a 99% mais por um item ou serviço de cinco estrelas do que por um que teve a cotação de quatro estrelas; para além disso 32% já forneceu uma classificação sobre um produto, serviço ou pessoa através dos sistemas de classificação *on-line* e 30% já comentou ou fez alguma revisão em relação a algum produto ou serviço.

### 2.3 Aplicações da Análise de Sentimento

Hoje em dia, é vital para as organizações saber quais as opiniões dos consumidores em relação aos seus produtos ou serviços. No passado, estas realizavam estudos, votações e discussões de grupo sempre que necessitavam de saber a opinião pública ou do consumidor. A obtenção da opinião pública tornou-se também por si mesma em um negócio de *marketing*, incluindo empresas de campanhas políticas e relações públicas, entre outras. Contudo com o crescimento das redes sociais, fóruns de discussão, *blogs*, *micro-blogs* (e.g. *Twitter*), as organizações estão a adotar cada vez mais a *Internet* para obter informação relevante para a tomada de decisão. Torna-se agora mais fácil para as organizações obter informação, sem necessitar de realizar estudos, votações e discussões de grupo. As aplicações de Análise de Sentimento têm proliferado em diversos domínios tais como saúde, serviços das finanças, eventos sociais, eleições políticas e produtos e serviços para o consumidor (Liu, 2012). Segundo Pang & Lee (2008), algumas das aplicações bem-sucedidas de Análise de Sentimento são por exemplo *websites* relacionados com avaliações, aplicações como uma subcomponente tecnológica e aplicações em negócios e inteligência governamental.

Segundo Osimo & Mureddu (2011), Análise de Sentimento ou *Opinion Mining* pode ter diversas aplicações como: *softwares* de mapeamento de *software* úteis na organização de afirmações políticas, explicando as conexões entre elas; aplicações de aconselhamento sobre o voto, que ajudam os eleitores a entenderem qual o partido político que possui políticas mais parecidas com as suas. Um exemplo deste tipo de *software* é o SmartVote.ch; aplicações que



identificam comentários relevantes e que lhes atribuem conotação positiva ou negativa, ou seja sentimentos. Estas ferramentas de *Opinion Mining* são bastante úteis nesta área, política, pois ajuda a obter *feedback* dos cidadãos (Osimo & Mureddu, 2011).

## 2.4 Desafios da Análise de Sentimento

A tarefa de Análise de Sentimento torna-se difícil no que toca a extração de informação relevante, devido ao constante crescimento de diversos sítios *Web* que possuem um grande número opiniões. Esta dificuldade leva à necessidade da utilização de Sistemas de Análise de Sentimento automatizado (Liu, 2012).

Existem alguns fatores que tornam o *Opinion Mining* difícil, tais como a polaridade dos sentimentos na classificação de textos, a possibilidade de expressar o mesmo sentimento de formas diferentes e a difícil identificação da entidade a quem se refere a expressão (Pang & Lee, 2008). Segundo Liu (2012), podemos classificar as opiniões como sendo regulares ou comparativas e explícitas ou implícitas. Opiniões regulares podem ter opiniões diretas e indiretas. Uma opinião direta refere-se a diretamente a um sujeito ou a um aspeto deste, contrariamente à indireta que é expressa, como o próprio nome indica, indiretamente sobre uma entidade ou aspeto desta baseando-se nos efeitos que esta provoca noutras entidades. Opinião comparativa refere-se às opiniões que relacionam duas ou mais entidades através das suas semelhanças ou diferenças. Opiniões explícitas podem ser regulares ou comparativas assim como as implícitas, sendo a primeira subjetiva e a segunda objetiva.

No caso de serviços de *microblogging*, tal como o *Twitter*, uma das principais questões a ter em conta é a classificação dos *tweets* entre os que mencionam atividades e naqueles que não mencionam. Apesar dos modelos anteriores funcionarem bem, estes necessitam de conseguir lidar com um número limitado de *tweets* e é necessário agregar um resultado global (sumário) para que este seja reconhecido como uma atividade (Weerkamp & Rijke, 2012). Segundo Liu (2012) as opiniões são subjetivas ao contrário das informações fatuais, e a opinião de uma pessoa normalmente não é suficiente para a definição de uma ação. Assim torna-se necessária a criação de um sumário para as diversas opiniões desejadas. O reconhecimento das palavras que representam uma indicação temporal é reduzido, sendo necessária uma forma automática de

reconhecimento destes termos. A segmentação dos *tweets* é necessária para por exemplo dividir a frase em atividades temporais (“*esta manha tenho escola, mas à noite vou ao cinema!*”). Os métodos mais simples identificam a “*escola*” como estando relacionada com “*noite*”, assim torna-se importante definir o âmbito das referências ao futuro. Uma das tarefas mais difícil é a avaliação da previsão de atividades, pois levanta várias questões (Weerkamp & Rijke, 2012): a atividade extraída é realmente uma atividade? É adequada para determinado período? é uma atividade popular? O sistema extrai todas as atividades populares para determinado período?

## 2.5 Funcionamento da Análise de Sentimento

O *Opinion Mining* pode ser definido como uma subcategoria de linguagem computacional que se foca na extração da opinião das pessoas da *Internet*. A Análise de Sentimento determina a subjetividade, polaridade e força desta, sobre determinado texto, ou seja verifica qual a opinião do autor. *Opinion Mining* analisa textos e extrai (Osimo & Mureddu, 2011):

- qual a parte que expressa opinião;
- quem deu a opinião; e
- o que se está a comentar.

Uma ferramenta de *Opinion Mining* deveria no seu ideal processar um conjunto de resultados relativos a determinado item, gerar uma lista de atributos do produto e posteriormente agregar opiniões relativas a cada atributo (Pang & Lee, 2008).

Segundo Liu (2012) e Aggarwal e Zhai (2012), uma opinião regular é um quintuplo que se exhibe na seguinte forma:

(*ei*, *aij*, *sijkl*, *hk*, *tl*),

onde *ei* é o nome da entidade, *aij* é um aspeto dessa entidade, *sijkl* é o sentimento relativo ao aspeto *aij*, da entidade *ei*. *hk* é o detentor da opinião *sijkl* e por último *tl* é o tempo em que esta opinião é expressa. Estes cinco componentes são essenciais pois a falta de qualquer um deles pode provocar problemas. Esta definição proporciona uma ferramenta que permite transformar textos não estruturados em dados estruturados que podem ser guardados em bases de dados (Liu, 2012).

Segundo Aggarwal e Zhai (2012), o objetivo do *Opinion Mining* é descobrir todos os quintuplos de opiniões de uma dada coleção de documentos  $D$ . Para a concretização deste objetivo são executadas seis tarefas:

Tarefa 1 – Extração da entidade e agrupamento;

Tarefa 2 – Extração dos aspetos e agrupamento;

Tarefa 3 – Extração do detentor da opinião e do tempo;

Tarefa 4 – Classificação do sentimento do aspeto; e

Tarefa 5 – Geração do quintuplo da opinião.

Na Tarefa 1 são extraídas todas as expressões de entidades de  $D$  e agrupadas com as que sejam sinónimas em *clusters* de entidades. Na Tarefa 2 são recolhidas todas as expressões relativas aos aspetos das entidades e guardadas também em *clusters*. Na Tarefa 3 são também recolhidas informações relativas ao detentor da opinião e ao tempo em que a opinião foi gerada. Na Tarefa 4 é determinado se cada opinião é positiva, negativa ou neutra, e por fim geram-se todos os quintuplos que sumarizam as opiniões encontradas na coleção de documentos  $D$ .

Para uma melhor compreensão das tarefas da Análise de Sentimentos é utilizado um exemplo semelhante utilizado por Aggarwal & Zhai (2012) e Liu (2012):

**“Exemplo:** Publicado por: João

Data: Setembro, 15, 2013

Comprei uma camara Samsung e o meu amigo comprou ontem uma câmara Canon. Na última semana, ambos usamos bastante as câmaras. As fotos da minha Samy não são grande coisa, e a duração da bateria também é pouca. O meu amigo estava muito contente com a sua câmara e adora a qualidade das suas imagens. Eu quero uma câmara que possa tirar boas fotos. Amanhã vou devolvê-la.”

Para começar deve realizar-se a Tarefa 1 e extrair as entidades “Samsung”, “Samy” e “Canon” e agrupar “Samsung” e “Samy” pois representam a mesma entidade. Em segundo lugar a Tarefa 2 deve extrair as expressões de aspetos ou seja, “imagens”, “fotos” e “duração da bateria” e agrupar “imagens” e “fotos”, pois são sinónimos. Seguidamente, a Tarefa 3 reside na procura dos detentores das opiniões, que são o João e o amigo do João, e o tempo, que é Set-15-2013. Na Tarefa 4 deve verificar-se qual a orientação da opinião, que é negativa para a

qualidade das fotos e da duração da bateria da camara Samsung, contudo esta é positiva no que diz respeito à câmara Canon de forma geral e relativamente às qualidades das fotografias que tira. Por último, é necessário realizar a Tarefa 5, que é relativa à produção dos quintuplos:

(Samsung, qualidade\_fotos, negativa, João, Set-15-2013)

(Samsung, duração\_bateria, negativa, João, Set-15-2013)

(Canon, GENERAL, positiva, amigo\_João, Set-15-2013)

(Canon, qualidade\_fotos, positiva, amigo\_João, Set-15-2013)

Para realizar a abstração da classificação de sentimentos é necessário ter em conta a definição da opinião e a sua sumarização. Estas são descritas de uma forma mais completa nas próximas subsecções.

### 2.5.1 Classificação e Extração

A **classificação** é fundamental em diversas aplicações atuais de *Opinion Mining*. A classificação engloba tarefas como: ordenação de um conjunto de textos, por exemplo segundo o seu grau de positivismo; fornecer uma única etiqueta à coleção inteira; e também realiza uma categorização das relações existentes entre duas entidades baseando-se em evidências textuais, como por exemplo “a entidade A aprova as ações da entidade B?” (Pang & Lee, 2008).

Determinadas aplicações, como as de sumarização e de resposta a perguntas, requerem informação de uma ou mais unidades textuais. A extração de informação foi concebida exatamente para resolver estas questões. Relativamente à extração de informação pode-se considerar que esta possui duas tarefas principais, que são a identificação das características do produto e a extração de opiniões relativas a estes (Pang & Lee, 2008).

### 2.5.2 Sumarização

A **sumarização** representa as tarefas de agregação e representação da informação relativa aos sentimentos encontrados num documento ou numa coleção. Esta tarefa é interessante pois um utilizador pode desejar ver os diversos pontos de vista encontrados nos documentos. Outra aplicação pode ser na determinação automática dos sentimentos de mercado, ou a inclinação

maioritária de um conjunto de investidores em relação às observações individuais destes (Pang & Lee, 2008). Estes sumários incluem as frases mais importantes do *input* que podem ser relativas a um documento ou a um *cluster* de documentos relacionados (Aggarwal & Zhai, 2012).

## 2.6 Níveis de Análise de Sentimento

### 2.6.1 Classificação dos Sentimentos ao Nível de Documentos

A classificação de sentimentos ao nível de documentos é realizada para detetar se o sentimento de um documento inteiro é positivo ou negativo. Contudo esta classificação assume que a opinião expressa é apenas relativa a uma entidade, logo não pode ser aplicada em documentos que avaliem ou comparem diversas entidades (Liu, 2012).

A maioria das técnicas de classificação dos sentimentos ao nível de documentos utiliza métodos de classificação de aprendizagem supervisionada, apesar de também existirem métodos de aprendizagem não supervisionada (Liu, 2012).

#### 2.6.1.1 Classificação de Sentimentos com Utilização da Aprendizagem Supervisionada

Para realizar o treino e teste desta aprendizagem são normalmente utilizadas *reviews* de produtos (e.g. avaliação numa escala de estrelas de 1 a 5, sendo valores entre 1 e 2 considerados negativos e os restantes positivos). Técnicas de classificação tais como *Support Vector Machines* (SVM) e *Naïve Bayes* utilizam a aprendizagem supervisionada para mapear um conjunto de entradas numa saída desejada (valor do sentimento). Para a definição das entradas e saídas são analisadas algumas características importantes, tais como:

**Termos e a sua frequência** – São utilizados “*unigrams*” para registar certas palavras e a quantidade de vezes que surgem;

**POS (*Part of Speech*)** – Por exemplo os adjetivos são considerados bastante importantes para indicar as opiniões;

**Palavras e frases de sentimentos** – Existem palavras e frases que conseguem expressar mais facilmente as opiniões. Os verbos e adjetivos são os mais habituais;

**Regras de opiniões-** Existem também determinadas expressões que tem significado sentimental;

**Deslocadores de sentimentos-** São utilizados para mudar a opinião de positiva para negativa ou vice-versa. Um exemplo destes deslocadores de sentimento é a negação;

**Dependências sintáticas-** Também podem ser utilizadas palavras que tenham dependências sintáticas de outras (Liu, 2012).

### 2.6.1.2 Classificação de Sentimentos com Utilização da Aprendizagem Não Supervisionada

A técnica fornecida pelo autor Turney (2002), é um exemplo de aprendizagem não supervisionada, pois segue padrões sintáticos. Esta segue três passos essenciais, começando pela utilização de *tags* POS e extração de duas palavras consecutivas se estas coincidirem com um dado padrão. O segundo passo é a estimação da orientação do sentimento através da medida *Pointwise Mutual Information* (PMI):

$$PMI(term_1, term_2) = \log_2 \left( \frac{\Pr(term_1 \wedge term_2)}{\Pr(term_1) \Pr(term_2)} \right)$$

$$SO(phrase) = PMI(phrase, \text{"excellent"}) - PMI(phrase, \text{"poor"})$$

$$SO(phrase) = \log_2 \left( \frac{hits(phrase \text{ NEAR } \text{"excellent"})hits(\text{"poor"})}{hits(phrase \text{ NEAR } \text{"poor"})hits(\text{"excellent"})} \right)$$

O terceiro e último passo resume-se à computação da média da Orientação do Sentimento (SO) de todas as frases. O resultado será um sentimento positivo ou negativo dependendo do resultado da SO (Turney, 2002).

Outra abordagem que utiliza aprendizagem não supervisionada é um método que utiliza um dicionário de palavras com sentimentos (léxico) e uma associação com a sua orientação e força (Liu, 2012).

### 2.6.1.3 Previsão da Classificação dos Sentimentos

Mais do que classificar mensagens cujo sentimento já é conhecido, pretende-se prever qual o sentimento de uma nova mensagem de texto. Para isto são usadas técnicas com capacidades de generalização via uma aprendizagem supervisionada, como o *SVM* e classificadores *Bayesianos*, entre outros (Liu, 2012).

### 2.6.1.4 Classificação de Sentimentos com Cruzamento de Domínios

Um classificador que seja treinado utilizando somente um determinado domínio tenderá a obter um desempenho fraco quando for aplicado a mensagens de textos que pertencem a outro domínio. Por exemplo, é comum que uma mesma palavra em contextos diferentes possa ter significados diferentes. Assim, é relevante, conseguir ter métodos que sejam capazes de classificar sentimentos em múltiplos domínios. Segundo o Liu (2012), diversas técnicas foram utilizadas nesta classificação com cruzamento de domínios, tais como *SVM*, *Structural Correspondence learning* (SCL) e *Spectral Feature Alignment* (SFA), entre outras.

### 2.6.1.5 Classificação de Sentimentos com Cruzamento de Línguas

Este tipo de classificação pretende encontrar os sentimentos em documentos que apresentem diversas línguas. Existem duas razões maioritárias para a existência deste tipo de classificação, que são a necessidade que os investigadores de diversos países apresentam de construir sistemas de Análise de Sentimento na sua própria língua. Para além desta, as empresas desejam saber mais sobre as opiniões dos seus consumidores (e.g. não só aqueles que as expressam em inglês mas também aqueles que são de outros países). Segundo o autor esta classificação pode ser realizada com abordagens com base em vocabulário, aprendizagem máquina, métodos de “*co-training*” que utilizam *SVM*, *Structural Correspondence learning* (SCL) e, método *Supervised Latent Dirichlet Allocation* (SLDA), entre outros (Liu, 2012).

### 2.6.1.6 Detecção de Polaridade e Grau de Positivismo

A classificação de polaridade de sentimentos ou simplesmente chamada de polaridade de sentimentos é uma tarefa de classificação binária que rotula um documento como sendo em

geral positivo ou negativo. Segundo o autor esta tarefa tem sido designada, por vezes, na literatura, por classificação sentimental.

É necessário entender o contexto, pois nem sempre a avaliação do texto em “positivo” ou “negativo” é uma tarefa simples. Em determinados casos, como por exemplo em textos orientados à área política, a questão da orientação das visões liberais ou conservadoras tem vindo a ser explorada. Para além disso *input* de um classificador de sentimentos não precisa de ser necessariamente opinante.

Quando um autor expressa explicitamente a sua opinião sobre determinado assunto (e.g., “Este computador é muito bom”), algumas informações mencionadas no texto (e.g., “bateria de longa duração”) são geralmente utilizadas para ajudar na determinação do sentimento global do texto. Contudo, determinar se uma parte da informação é “positiva” ou “negativa” não é exatamente o mesmo que classificar essa informação em classes, sendo que mesmo a distinção entre informação subjetiva e objetiva pode ser subjetiva, por exemplo dizer “a bateria dura horas” é diferente de “a bateria apenas dura horas”.

As **Categorias relacionadas** são uma forma alternativa de sumarização de revisões. Expressões de “prós” e “contras” podem diferir de expressões de opinião positiva ou negativa. Por exemplo podemos considerar a opinião “considero este computador incrível” e a razão “apenas custa 400 euros”. A identificação de prós e contras pode ajudar na construção de mais sumários orientados aos sentimentos e podem ser ajudar a decidir se a avaliação do autor foi útil ou não.

O caso mais comum de **Inferência de classificação** é a determinação da avaliação do autor relativamente a uma escala de vários valores, como por exemplo revisões de zero a cinco estrelas. E pode ser visto como um problema de categorização com várias classes. Saber qual o grau de positividade pode fornecer informação mais detalhada em relação à revisão feita. Para além disso, cada classe pode possuir o seu próprio vocabulário.

Devido à oposição entre a natureza da polaridade das classes dá-se espaço para a exploração da **deteção de concordância**. Por vezes surge a necessidade de decidir se dados dois textos, estes devem receber o mesmo rótulo de sentimento ou não, tendo em conta a relação



existente entre os elementos. O seu resultado é utilizado para melhorar a rotulação das opiniões portadas por diferentes partes (Pang & Lee, 2008).

#### 2.6.1.7 Determinação de Sentimentos Coletivos

Existem diversos métodos que utilizam a determinação de métodos coletivos, no caso do *Twitter* são utilizados os parâmetros como a contagem dos *Tweets*, o número de *followers*, tempo da publicação, entre outros (Xu, 2012).

#### 2.6.1.8 Articulação da Análise Temática de Sentimentos

Para simplificar a classificação dos sentimentos por vezes assume-se que o documento todo se refere ao sujeito em que se tem interesse. Isto acontece porque por vezes quem realiza a classificação assume que a recolha de documentos sobre um único tópico se tenha realizado anteriormente. Contudo é possível que haja interações entre tópicos e opiniões que faz com que seja melhor considerá-los simultaneamente em vez de isolados. Para além disso, pode surgir um documento que contenha material sobre diversos temas, de interesse para o utilizador, tornando assim necessária a identificação dos tópicos e posterior associação das opiniões ao correspondente. Alguns exemplos de documentos que utilizam esta análise são estudos comparativos de produtos relacionados, e textos que abrangem vários aspetos ou atributos de determinado tema (Pang & Lee, 2008).

#### 2.6.1.9 Perspetivas e Pontos de Vista

As perspetivas e pontos de vista são casos mais complexos de Análise de Sentimento, como é o caso já visto anteriormente relativo a textos orientados à política, onde se determinam orientações políticas. Estes podem ser classificados por exemplo de liberal, conservador, ou outro. Mais ainda, pode ser-lhes também atribuída uma escala relativa às classificações que obtiveram. Uma outra forma de representar estes pontos de vista pode ser pela comparação de duas ideologias, ou perspetivas. Normalmente as classes não correspondem a opiniões isoladas mas a coleções de atitudes e crenças.

A motivação para este tipo de abordagem da Análise de Sentimento foi a possibilidade de responder a perguntas de várias perspectivas, em vez de questões baseadas em factos (Pang & Lee, 2008).

#### 2.6.1.10 Outros Textos com Informação Não Factual

Segundo Pang & Lee (2008), podemos considerar as seis emoções “universais” que são: medo, raiva, desgosto, felicidade, tristeza e a surpresa. Uma possível aplicação desta classificação é uma aplicação de interação entre o utilizador e o computador, em que se o sistema determina qual o estado emocional do seu utilizador, sendo que se este se encontrar de alguma forma aborrecido, o interface pode alterar de modo automático para outro modo de interação. Outra área que utiliza este método é a de determinação do género dos textos (Pang & Lee, 2008).

### 2.6.2 Classificação dos Sentimentos ao Nível das Frases

Neste tipo de análise é verificado se determinada frase possui uma opinião positiva, negativa ou neutra, com vista a uma classificação subjetiva. Um documento possui diversas opiniões, enquanto uma frase, geralmente, tem apenas uma. Segundo Liu (2012), esta deteção de subjetividade torna-se mais difícil do que classificar a polaridade. Determinar a força da opinião é diferente da inferência de classificação, pois classificar determinado texto como neutro não é dizer que este é objetivo, pois pode mostrar falta de opinião (Pang & Lee, 2008).

Este tipo de classificação pode ser resolvido seguindo dois passos distintos de classificação. O primeiro passo resulta na determinação da existência ou não de opinião na frase e denomina-se de classificação de subjetividade. O segundo passo permite determinar de esta opinião encontrada no primeiro é positiva ou negativa.

#### 2.6.2.1 Classificação Subjetiva

Neste passo é determinado se a frase apresenta uma opinião subjetiva ou objetiva, ou seja, se esta apresenta opinião ou sentimentos (opinião subjetiva) ou não (opinião objetiva). Uma

opinião objetiva expressa informação factual ao contrário da subjetiva que pode demonstrar crenças, opiniões, emoções, especulações, avaliações, entre outros (Wiebet et. al, 1999).

Os adjetivos são bons indicadores de subjetividade. Pode ser utilizada a propriedade semântica “*gradability*”, pois esta permite aos adjetivos apresentar diversos graus de força (Liu, 2012). Segundo o autor, podem ser utilizados diversos padrões, já estudados por diversos investigadores, para gerar métodos baseados em regras, e a informação resultante destes é utilizada para treinar classificadores como *Naive Bayes*. No seu artigo, Barbosa & Feng (2010), descrevem a utilização do algoritmo SVM, devido ao seu melhor desempenho na utilização de classificadores como *Unigrams* e *TwitterSA*.

Por vezes uma frase pode apresentar partes subjetivas e partes objetivas, logo é útil a identificação da força da subjetividade. Segundo Liu (2012):

**S** (Subjetiva e avaliativa) - Sentimento positivo ou negativo;

**OO** – Opinião positiva ou negativa implícita numa frase objetiva;

**O** – Frase objetiva sem opinião;

**SN** – Subjetiva sem avaliação, ou seja, não possui nem sentimento positivo nem negativo;

#### 2.6.2.2 Classificação de Sentimentos de Frases

Aqui assume-se que uma frase possui apenas um sentimento e um único possuidor da mesma. Para descobrir a orientação do sentimento, ou seja, se este é positivo ou negativo, alguns investigadores sugerem a utilização de algoritmos baseados em léxicos. Para descobrir esta orientação são atribuídos valores a todas as palavras com sentimentos (e.g. palavra positiva equivale a +1 valor, palavra negativa equivale a -1 valor) e no final realiza-se uma soma destes valores (Liu, 2012).

Na investigação realizada por (Davidov, et al., 2010), relativa à classificação de sentimentos na rede social *Twitter*, utilizam, para além das características usuais, *hashtags*, *smiles*, pontuações e os seus padrões de frequência. A utilização destas características provou-se bastante benéfica na obtenção de resultados corretos.

### 2.6.2.3 Frases Condicionais

Relativamente a frases condicionais pouco se investigou, contudo o problema existe e pode influenciar muitas vezes o sistema de forma a prejudicar os resultados. Uma frase condicional pode possuir por vezes, palavras que expressam um sentimento positivo e no entanto ter no seu contexto um sentido negativo ou neutro, ou vice-versa. Na frase “Se alguém fizer um carro fiável, eu irei comprá-lo” são utilizadas palavras que expressam um sentimento positivo, contudo, não se pode afirmar que esta possua uma opinião positiva. Para resolver esta situação alguns investigadores propõem a utilização de características linguísticas como a existência e localização das palavras com sentimentos, *tags POS* também das palavras que apresentam sentimentos, padrões de tensão, conectores condicionais, entre outros. Para além das frases condicionais, existem também as interrogações, que a grande parte das vezes apresentam um sentimento neutro (Liu, 2012).

### 2.6.2.4 Frases Sarcásticas

Relativamente à Análise de Sentimentos em frases sarcásticas, o problema surge pois o autor pretende dizer o oposto do significado das palavras que transmite. Assim dificulta-se a tarefa de detetar se a frase é positiva ou negativa. Apesar de ser difícil de lidar com este tipo de frases, segundo o autor, estas são pouco frequentes em avaliações de serviços ou produtos, no entanto surgem muitas vezes em comentários políticos e discussões *on-line* (Liu, 2012).

Os autores (González-Ibáñez et. al, 2011) realizaram um estudo sobre esta problemática e utilizaram para sua resolução, um classificador *SVM*, regressão logística e como características utilizaram *unigrams* e outras baseadas em dicionários.

### 2.6.2.5 Subjetividade entre Linguagens

Relativamente a este tópico o significado é o mesmo do enunciado relativamente à classificação de sentimento realizada ao nível de documentos. Atualmente existem três estratégias para resolver esta problemática (Liu, 2012):

- traduzir frases de teste da língua de destino para a língua de origem e classificá-las utilizando um classificador fonte de linguagens;

- traduzir o corpus de treino da linguagem de origem para a língua de destino e construir um classificador com base no corpus na língua de destino;
- traduzir um sentimento ou um léxico de subjetividade na língua fonte para a língua de destino e construir um classificador com base no léxico na língua de destino.

## 2.7 Sumarização das Opiniões

O sumário é necessário quando são estudadas as opiniões de muitas pessoas, ou seja, quando a opinião de uma única pessoa não é suficiente. Os sumários realizados com base em quintuplos  $(e_i, a_{ij}, o_{ijkl}, h_k, t_l)$ , são designados por sumários baseados em aspetos. O *output* da sumarização pode ser um documento estruturado ou não estruturado (Liu, 2012).

### 2.7.1 Sumarização de Opiniões Baseadas em Aspetos

A sumarização de opiniões baseada nos aspetos tem duas características essenciais, ou seja esta capta a essência das opiniões (os alvos da opinião e os sentimentos relativos a estes) e apresenta os dados de forma quantitativa. De forma quantitativa pois expõe os números ou percentagens de pessoas que possuem sentimentos positivos e negativos. Esta sumarização pode ser apresentada sob várias formas, como por exemplo em gráficos (Liu, 2012).

### 2.7.2 Sumarização Tradicional

Um sumário de opiniões pode ser diferente dos sumários tradicionais de um documento ou de vários documentos. Os sumários de opiniões baseados nos aspetos focam-se nas entidades, aspetos, sentimentos relacionados, e quantidades dos mesmos. Da sumarização tradicional de um documento só, resulta um texto pequeno que possui as frases mais importantes. Na sumarização tradicional de vários documentos, realiza-se uma filtragem de informação repetida e tendo em conta os algoritmos de sumarização são definidas as frases importantes (Liu, 2012).

### 2.7.3 Sumarização Orientada à Opinião num Documento Único

Existem diversas abordagens que criam sumários de sentimentos baseados na extração de frases ou unidades de texto similares. O autor (Pang & Lee, 2008) sugere que através da localização do fluxo de sentimentos do documento, pode criar-se sumários de sentimentos escolhendo as frases encontradas nos extremos do fluxo. Esta abordagem torna-se interessante pois ao considerar o fluxo do documento, esta técnica considera o documento de forma global. Outra forma é a extração das frases subjetivas, que podem ser usados como sumários. Para além destes métodos podemos ainda encontrar outros métodos que trabalham diretamente com o *output* de sistemas de extração de informação orientados à opinião.

Um aspeto interessante da extração de informação relativa aos sentimentos de um documento único é que por vezes, o *output* baseado em grafos parece mais apropriado ou útil do que o *output* baseado em textos. Os sumários baseados em grafos são apropriados para situações em que a informação mais importante é o conjunto de entidades descritas e as opiniões que estas possuem umas sobre as outras.

Os elementos gráficos podem ser utilizados para representar apenas uma parte, como um sumário dos sentimentos encontrados num documento (Pang & Lee, 2008).

## 2.8 Dicionários de Léxicos

Atualmente, existem diversas propostas de métodos e recursos léxicos para analisar os sentimentos de uma frase. Os autores Bravo-Marquez, et. al, (2013), sugerem a utilização de um método que combine indicadores de força, emoção e polaridade. Os autores Alborn et. al, (2012), sugerem a realização de quatro tarefas: deteção de subjetividade, classificação da polaridade, classificação da intensidade e identificação da emoção. Palanisamy et. al, (2013), utilizam uma abordagem baseada em léxicos que consiste na deteção de palavras e frases positivas, negativas e de paragem. Devido ao conteúdo dos *tweets* possuir, muitas vezes, *hashtags*, *emoticons* e palavras alteradas, os autores sugerem a realização de um pré-processamento. As abordagens baseadas em léxicos baseiam-se na suposição de que o sentimento geral é o resultado do somatório das orientações sentimentais de cada palavra ou frase.

Esta abordagem usa dicionários de palavras com opiniões para ser possível determinar a orientação do sentimento expresso. Estes dicionários são designados léxicos de opinião e a abordagem é baseada em léxicos. Contudo os *emoticons*, abreviaturas, expressões, entre outros, não se encontram nestes dicionários. Isto é uma desvantagem desta abordagem, pois estas palavras ou expressões possuem muitas vezes sentimentos que não são considerados. Estas expressões e palavras podem ser adicionadas aos léxicos, contudo estas estão em constante alteração, e surgem cada vez mais e de formas diferentes, seguindo modas na *Internet*. Em alternativa a esta abordagem baseada em léxicos pode utilizar-se os métodos baseados em aprendizagem máquina, para detetar os sentimentos (Zhang et. al, 2011).

Segundo Zhang et. al. (2011), e como se pode verificar na Figura 1, o processo de classificação de Sentimento de *Tweets* utilizando léxicos, apresenta diversas fases. A fase inicial representa a obtenção dos *Tweets* para posterior análise dos mesmos. Seguidamente é realizado um pré-processamento para limpar os dados. Como é sabido, o *Twitter* apresenta uma linguagem com convenções próprias tais como: RT (utilizado para *Retweets*); # (o chamado *hashtag* utilizado para marcar determinado *Tweet* com determinado tópico ou categoria); @*username* (representa uma resposta a determinado utilizador cujo nome é “username”); *Emoticons* ou expressões como “lowve”, “:-)”, “lmao”, “wknd” “”; e *Links* externos. Estas palavras devem ser removidas, no caso dos *Retweets*, *links* externos e nomes de utilizadores, ou melhoradas como é o caso das abreviaturas. Seguidamente é realizado o método baseado em léxicos de sentimentos e por fim são extraídos os Sentimentos destes e Classificados.

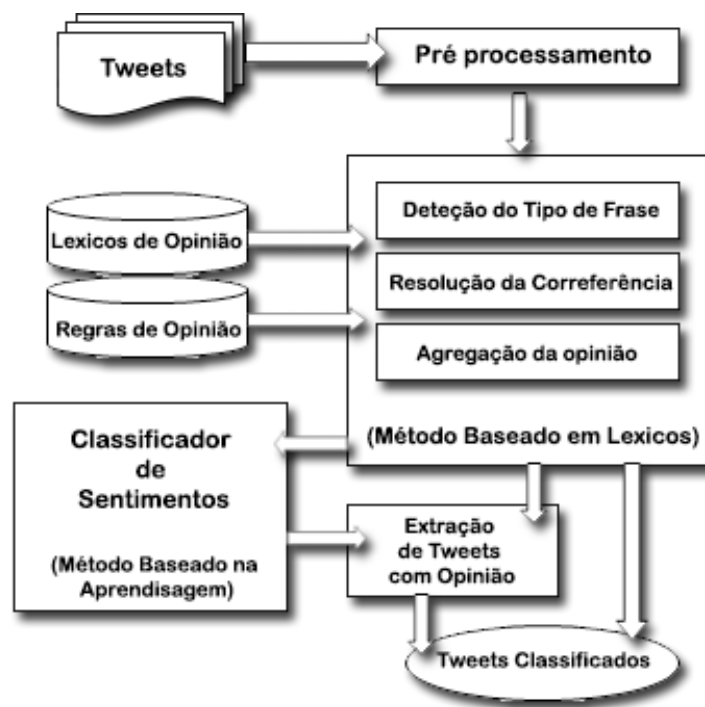


Figura 1 - Processo de extração de Sentimento de *Tweets* através da utilização de léxicos, adaptado de (Zhang et. al, 2011)

Segundo Feldman (2013), um Léxico de Sentimentos pode ser adquirido de três formas diferentes. Uma delas é a abordagem manual, na qual o Léxico é codificado manualmente, outra abordagem é baseada em dicionários, em que um conjunto de palavras é expandido utilizando recursos tais como *WordNet*. Por fim, outra abordagem é baseada no corpus, que expande o conjunto de palavras utilizando um conjunto elevado de documentos de um único domínio. De acordo com o autor a abordagem manual não é muito viável pois cada domínio possui o seu próprio Léxico o que leva a uma sobrecarga de trabalho. A abordagem baseada em dicionários inicia-se com um conjunto pequeno de palavras com sentimento que são então expandidas através da utilização de sinónimos e antónimos da *WordNet*. Esta abordagem possui a desvantagem de não ser independente do domínio em que o texto se encontra, ou seja não captura as peculiaridades específicas de cada domínio. Isto pode ser melhorado com a abordagem baseada no corpus, que produz um Léxico de sentimentos específico de determinado domínio.

Segundo Herbrich e Graepel (2010), a abordagem baseada em dicionários é das técnicas mais simples e baseia-se em *bootstrapping* utilizando um conjunto de palavras com Sentimento e um dicionário *on-line* como o caso do *WordNet*. Segundo estes autores, é necessário recolher



manualmente um conjunto de palavras com opinião com orientações conhecidas e posteriormente aumentar este conjunto através da utilização de um dicionário para encontrar sinónimos e antónimos destas. E finalmente estas palavras são adicionadas à lista dos termos pertencentes ao léxico.

Existem diversos léxicos que podem ser utilizados para obter sentimentos das frases, tais como *SentiWordNet*, *Q-WordNet*, *WordNet-Affect*, *Web Search*, *WordNet*, *OpinionFinder*, *Sentistrength*, *Sentiment140*, *SentiSence*, *MPQA* (Liu, 2012).

### 2.8.1 *SentiWordNet*

Este Léxico é utilizado para obtenção de sentimentos nas frases e pode ser obtido em “<http://sentiwordnet.isti.cnr.it/>” (Feldman, 2013). *SentiWordNet* 3.0 é a versão atual do Léxico e é uma versão melhorada do *SentiWordNet* 1.0. Este está disponível para fins de investigação e é usado em projetos a nível mundial. É o resultado da anotação automática dos “*synsets*” do *WordNet* de acordo com as noções de positividade, negatividade e neutralidade (Baccianella et. al, 2010). “*Synsets*” são grupos de palavras ou expressões semanticamente equivalentes em determinado contexto (Oliveira et. al, 2014). Este Léxico é uma extensão do *WordNet* e também uma base de dados lexical de palavras inglesas que são agrupadas por sinónimos. Estes são anotados no *SentiWordNet* num intervalo de -1 (extremamente negativo) a 1 (extremamente positivo) de acordo com o sentimento expresso, no qual 0 equivale a um sentimento neutro. Estes valores são obtidos a partir de algoritmos com “semi-supervisão” (Bravo-Marquez et. al, 2013).

### 2.8.2 *Financial Sentiment Dictionaries (FIN)*

O objetivo deste léxico é classificar palavras comuns em textos financeiros. Segundo Wang et. al, (2013), três quartos das palavras financeiras encontradas em relatórios do ano 1994 até 2008, consideradas negativas pelo dicionário Psicossociológico de Harvard, não possuem essa conotação negativa em contextos financeiros.

Este léxico possui seis listas de palavras aplicadas a contextos financeiros. Uma das listas possui palavras negativas, outra possui palavras positivas. Outra das listas possui palavras

incertas, que demonstram imprecisão. A lista de palavras litigiosas apresenta palavras que têm tendência a surgir em contextos legais. As listas de palavras fortes e de palavras fracas apresentam palavras níveis de confiança, fortes ou fracos respectivamente (Wang et. al, 2013).

## 3 Inteligência Coletiva e *Crowdsourcing*

### 3.1 Introdução

Hoje em dia a *Web* é fortemente colaborativa, ou seja possibilita aos utilizadores individuais, colaborar com outros indivíduos para atingir um determinado objetivo. Esta é uma característica marcante na *Web 2.0*, pois permite a realização da Inteligência Coletiva (Michalsky et. al, 2010). Graças à *Internet* e às Tecnologias de Informação, é cada vez mais fácil aceder à informação desejada. Contudo numa organização é sempre necessário tomar decisões, e é nessa fase que a Inteligência Coletiva se apresenta vantajosa (Bonabeau, 2009). Sob as circunstâncias certas, um grupo de pessoas medianas, consegue obter melhores resultados do que, cada elemento do grupo obteria individualmente. Isto demonstra que a vantagem da Inteligência Coletiva não está no consenso nem nos compromissos, mas sim na competitividade e nas opiniões diversificadas e independentes que esta permite (Leimeister, 2010). A Inteligência Coletiva tem sido mais divulgada e os indivíduos estão a aderir cada vez mais a esta ficando cada vez mais à vontade no que diz respeito à partilha de ideias em espaços virtuais. Desta forma as organizações podem utilizar mais este método para resolver desafios complexos do mundo atual. Assim a Inteligência Coletiva pode desempenhar um papel importante nas empresas, permitindo a introdução de ideias novas e diversificadas, resolvendo problemas antigos, desagregando e distribuindo o trabalho de formas mais inovadoras e diferentes do habitual, e sobretudo permitindo tomar decisões, com mais informação, sobre o futuro (Lesser et. al, 2012).

Importa ainda referir um outro termo que está intimamente ligado à Inteligência Coletiva: o *Crowdsourcing* – Trata-se de uma área relativamente recente que emergiu nos últimos anos e que denota um modelo distribuído de resolução de problemas *on-line* (Brabham, 2008). Sendo esta uma área emergente, naturalmente, este conceito surge com diferentes nomes, dos quais, “*peer production*”, “*user-powered systems*”, “*user-generated content*”, “*collaborative systems*”, “*community systems*”, “*social systems*”, “*social search*”, “*social media*”, “*collective intelligence*”, “*wikinomics*”, “*crowd wisdom*”, “*smart mobs*”, “*mass collaboration*”, “*human computation*” são alguns exemplos (Doan et. al, 2011). As aplicações que utilizam *Crowdsourcing* fornecem uma nova forma de resolver problemas que podem ser generalizados e aplicados a diversas indústrias para resolver tarefas tanto comuns como complexas. Esta é uma

forma de atrair pessoas capazes de fornecer ideias e soluções com grande qualidade e quantidade (Brabham, 2008).

Existem diversos tipos de participação na inteligência coletiva. Uma dessas formas é o *reCAPTCHA*, um tipo de *CAPTCHA* que permite não só a detecção de programas automáticos intrusos ao sistema como reaproveita o esforço realizado pelos utilizadores humanos ao mapear os caracteres encontrados nas imagens (Schlaikjer, 2007).

## 3.2 Inteligência Coletiva

### 3.2.1 Definição

A Inteligência Coletiva é usada de diversas formas à décadas, e tem-se tornado cada vez mais popular e importante devido às novas tecnologias. Apesar desta expressão poder parecer relacionar-se com grupos de consciencialização e fenómenos sobrenaturais, esta é na verdade utilizada para definir um grupo de pessoas preocupadas em obter determinada solução. Logicamente a Inteligência Coletiva era possível antes de existir a *Internet*, contudo esta veio facilitar em grande escala a realização desta (Segaran, 2007).

**Inteligência Coletiva** é definida por Michalsky et al. (2010) como sendo um grupo de indivíduos que colabora entre si ou compete por um determinado objetivo. É o conhecimento e experiência de determinados indivíduos, que se podem encontrar dentro ou fora das barreiras formais de uma empresa (Lesser et. al, 2012). Segundo Glenn (2013) Inteligência Coletiva é uma propriedade emergente de sinergias entre três elementos: os dados, informação ou conhecimento, o *software* ou *hardware* e os especialistas. Segundo este, um sistema de Inteligência Coletiva eficiente, deveria combinar estes três elementos numa plataforma interoperável, para que os processos e produtos pudessem ser visualizados e modificados por outros utilizadores. Os seus utilizadores deveriam ter permissões para comentar qualquer informação, *software* ou modelo computacional desse mesmo sistema (Glenn, 2013). Segundo Alag (2009), a Inteligência Coletiva significa a utilização de grupos com sabedoria para realizar a sua aplicação, ou seja usar a informação fornecida por uma *crowd* para melhorar determinada aplicação de outra pessoa. Os participantes da *crowd* interagem entre si através da *Web* e

expressão as suas opiniões e influenciam os outros participantes. Isto gera um círculo de influências que cresce rapidamente e pode moldar as opiniões dos seus participantes.

Para compreender melhor o termo de Inteligência Coletiva, pode analisar-se a sua decomposição. **Inteligência** refere-se à capacidade de aprender, compreender e adaptar a um ambiente usando o seu próprio conhecimento (Leimeister, 2010). Por sua vez, **Coletiva** descreve um grupo de indivíduos, que não têm necessariamente as mesmas atitudes, nem os mesmos pontos de vista.

### 3.2.2 Classificação

A Inteligência Coletiva pode ser dividida em três tipos: Inteligência Coletiva consciente, Inteligência Coletiva inconsciente e Inteligência Coletiva derivada. Esta classificação é realizada de acordo com a forma como a Inteligência Coletiva é gerada. A Inteligência Coletiva inconsciente permite ao utilizador, não se aperceber de que está a melhorar o próprio sistema (Michalsky et. al, 2010).

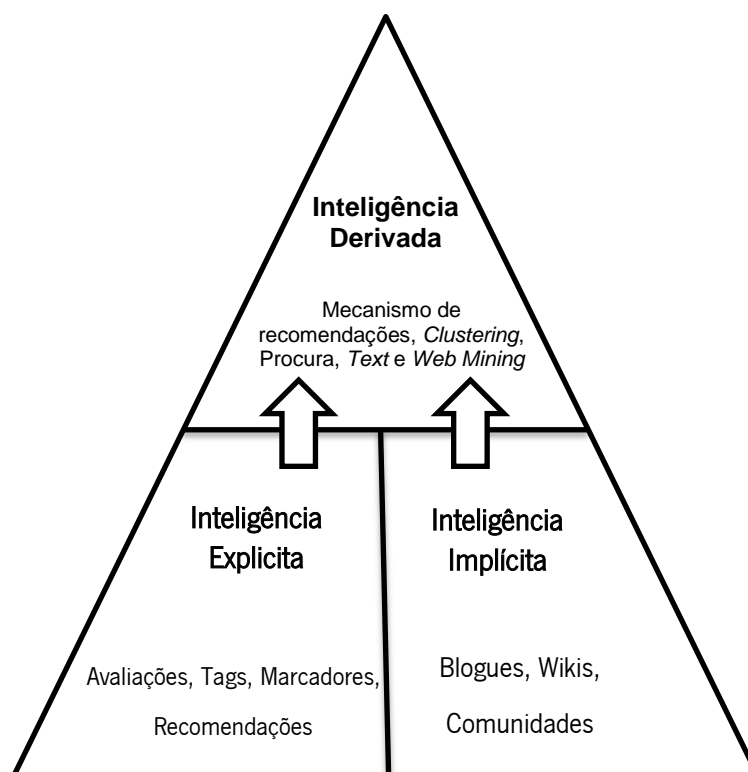


Figura 2 - Classificação da Inteligência Coletiva, adaptado de (Alag, 2009)

Os dados podem apresentar-se em dois formatos, ou seja, podem ser estruturados ou não estruturados. Os dados estruturados, têm uma forma bem definida, o que facilita o seu armazenamento e por sua vez as consultas aos mesmos. O tipo de dados influencia a forma como a Inteligência Coletiva deve ser utilizada. Segundo Alag (2009), como mostra a Figura 2, a inteligência pode ser explícita, implícita ou derivada. A inteligência explícita, é utilizada em recomendações e avaliações, *tags*, votações, entre outros. A inteligência inconsciente pode ser expressa através de mensagens, blogs, avaliações, entre outros, que permitam ao utilizador exprimir a sua opinião. Enquanto a inteligência derivada lida com informação obtida a partir dos utilizadores, por exemplo através de *Text* e *Data Mining*, análise preditiva, *clustering*, procura inteligente e mecanismos de recomendações (Alag, 2009).

### 3.2.3 Exemplos

Embora já existissem métodos de Inteligência Coletiva antes da *Internet*, esta permitiu a obtenção de informação de milhares ou mesmo milhões de pessoas, o que levou ao surgimento de novas possibilidades (Segaran, 2007).

A *Wikipédia* é um dos exemplos mais conhecidos de Inteligência Coletiva, pois é uma enciclopédia *on-line* criada somente pelas contribuições dos seus utilizadores. Esta possui um número reduzido de administradores e as páginas podem ser criadas e editadas por qualquer pessoa. O seu *software* apenas procura as alterações e a versão mais recente para apresentação das mesmas. Ou seja, este é um exemplo de Inteligência Coletiva, pois cada artigo é criado e modificado por um vasto grupo de pessoas e como resultado surge a maior enciclopédia existente. Apesar da manipulação de certos utilizadores maliciosos, esta enciclopédia é geralmente precisa na maior parte dos assuntos (Segaran, 2007). Esta apresenta já mais de dois milhões de artigos só na versão inglesa. Apesar de ser a versão mais popular do paradigma das *wikis*, esta não é o único exemplo. Existem diversas *wikis*, utilizadas por exemplo em empresas, como a *socialtext.com*, nas escolas, na análise de inteligência, como a *Intellipedia*, nas ciências, como a *Scholarpedia* e a *OpenWetwave*, entre outras. O *software* que a *Wikipedia* utiliza, *MediaWiki*, é também utilizado por outras *wikis* como é o caso da do jogo de computador *World of Warcraft* com mais de um milhão de utilizadores. As páginas da *wikipedia* podem ser visualizadas por qualquer pessoa, e esta possui assuntos em diversas línguas, estruturas fortes

tais como as subcategorias, *links*, avaliações de qualidade, e *templates* que podem ser extraídos pelo utilizador (Kittur & Kraut, 2010).

O Google é outro exemplo de utilização de Inteligência Coletiva. Este é o motor de busca mais popular e foi o primeiro a classificar páginas *Web* baseando-se na quantidade de páginas que se ligam a esta. Este método de utilização de Inteligência Coletiva é um pouco diferente da usada na wikipedia, pois recolhe a informação que milhares de pessoas disseram sobre determinada página *Web* e utiliza essa informação para posicionar os resultados de uma procura (Segaran, 2007).

A Amazon possui um serviço designado por *Mechanical Turk* que possibilita um mercado para trabalhos que necessitem Inteligência Coletiva. Os empregadores publicam determinadas tarefas e os trabalhadores realizam-nas por uma recompensa monetária (tipicamente de valor reduzido a nível individual). Esta agrupa um vasto número de trabalhadores assim como tarefas bastante diversificadas (Kittur et. al, 2011).

Convém referir que existem opiniões contrárias ao que foi exposto anteriormente. Por exemplo, segundo o autor Glenn (2013), a wikipedia, Google, *Crowdsourcing*, mercados de previsão, entre outros, podem não ser considerados exemplos puros de Inteligência Coletiva no que diz respeito à sua definição desta. Tal afirmação baseia-se numa definição de Inteligência Coletiva que inclui *feedback* sistemático e de forma contínua entre os elementos que o compõem. Segundo o autor, os exemplos anteriormente mencionados não produzem uma inteligência emergente e contínua, mas antes uma imagem no tempo.

#### 3.2.4 Inteligência Coletiva nas Aplicações *Web*

Neste trabalho assume-se que a Inteligência Coletiva dos utilizadores é a inteligência recolhida de um conjunto de interações e contribuições realizada por estes. Pode também ser a utilização da inteligência para filtrar o que interessa a determinado utilizador, numa aplicação. Este filtro pode utilizar as preferências ou as interações do utilizador para lhe poder fornecer informação mais relevante (Alag, 2009).

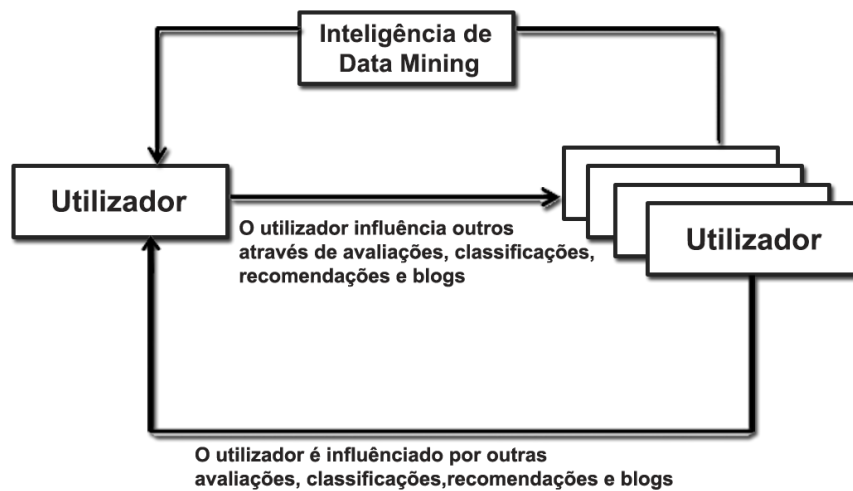


Figura 3 – Influência dos utilizadores e da Inteligência Coletiva, adaptado de (Alag, 2009)

O filtro pode agir de duas formas, ou seja, através de uma simples influência que a informação coletiva tem no indivíduo, como por exemplo através de uma avaliação de determinado produto (como mostra na Figura 3), ou através da construção de modelos de recomendação de conteúdo para os indivíduos (Alag, 2009).

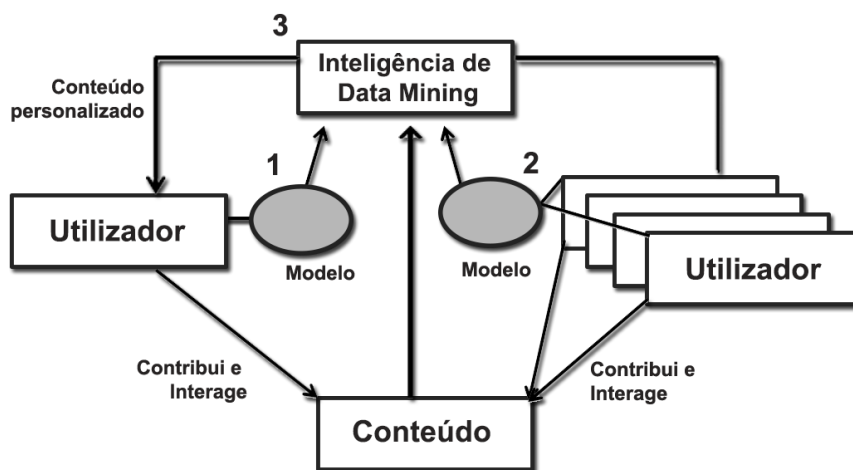


Figura 4 - Inteligência Coletiva nas aplicações Web (1- Permite interação dos utilizadores; 2- Aprende sobre os utilizadores em conjunto; 3 – Personaliza conteúdos tendo em conta dados das interações dos utilizadores); adaptado de (Alag, 2009)

Para poder personalizar as aplicações Web os sistemas devem proporcionar modelos diferentes construídos para o efeito. Na Figura 4 são apresentados os elementos da Inteligência Coletiva nas aplicações Web. Esta deve permitir aos utilizadores a interação com a aplicação e



entre si, fazendo com que estes possam aprender com as interações e contribuições dos outros. Para além disso, a aplicação deve utilizar modelos que permitam agregar a informação recolhida relativa aos utilizadores e as suas contribuições. A aplicação deve também influenciar os modelos de forma a recomendar a informação mais relevante para um dado utilizador (Alag, 2009).

### 3.2.5 Técnicas

A Inteligência Coletiva pode manifestar-se numa aplicação de diferentes formas. Seguidamente são mostrados alguns exemplos de formas de aplicar a Inteligência Coletiva numa aplicação (Alag, 2009):

**Agregar Informação através de listas** - Criação de uma lista de itens gerados e agregados pelos utilizadores. Por exemplo, via uma lista de produtos mais comprados, lista de itens recentes ou lista de produtos mais procurados.

**Classificações, Avaliações e Recomendações** - Informação coletiva que os utilizadores produzem e que vai influenciar outros.

**Conteúdo gerado pelo utilizador** - Os utilizadores podem extrair a inteligência de informação gerada e contribuições de outros utilizadores. Este conteúdo influencia muitas vezes outros utilizadores. Por exemplo: *Wikis, Blogs* e *Fóruns*.

**Votar, guardar, atribuir tags, bookmark** - A Inteligência Coletiva pode ser utilizada para emergir o conteúdo relevante, aprender mais relativamente aos utilizadores e assim poder conectá-los.

**Navegação através da *clouds de tags*** - A classificação dinâmica de conteúdo pode usar técnicas como *machine-generated, professionally-generated ou user-generated*.

**Análises de conteúdos para construir perfis de utilizadores** - Através da análise de conteúdos associados ao utilizador, são extraídas palavras-chave que são utilizadas para construir perfis destes.

**Modelos preditivos e *Clustering*** - Agrupando utilizadores e itens e construindo modelos preditivos.

**Mecanismos de recomendação** - A obtenção de inteligência dos indivíduos a partir da interação dos utilizadores e da análise de conteúdo permitem a recomendação de conteúdos relacionados com aqueles que o utilizador interagiu.

**Procura** - Permite apresentar resultados de procuras, mais pertinentes, utilizando perfis dos utilizadores.

**Incorporar conteúdos externos** - Fornece informações relevantes provenientes da blogosfera e de *websites* externos.

### 3.2.6 Áreas de Aplicação

A divulgação de tecnologias fáceis de usar, que permitem aos utilizadores interagir e construir aplicações *Web* sem necessidade de possuir conhecimentos de programação, levou à geração de muito conteúdo por parte dos utilizadores. Os utilizadores podem contribuir diretamente e de forma mais crítica através da *Web*, o que contribui para uma maior força coletiva. Este comportamento demonstra a capacidade que as massas têm de atingir determinado objetivo, através participação e colaboração na *Web*. O desafio é tentar perceber como libertar o conhecimento e experiência dos empregados, consumidores ou parceiros para obter aproveitamento da sua Inteligência Coletiva. De seguida são expostas algumas abordagens e áreas de aplicação segundo Alag (2009):

#### **Apoio na Tomada de Decisão**

A Inteligência Coletiva pode ser utilizada na tomada de decisão e pode ser resumida em duas tarefas principais: a construção de potenciais soluções e a sua avaliação. Estas podem ser influenciadas negativamente por diversos vieses humanos, e mitigadas através do uso de três abordagens: divulgação, aglomeração adicional e auto-organização (Alag, 2009). Segundo Bonabeau (2009), existem diversos vieses que podem influenciar a formação de soluções potenciais assim como a sua avaliação. Estes podem ser mitigados de diversas formas como se pode verificar nas Tabelas 1 e 2. Na Tabela 1 estão descritos os vieses e mitigações para o processo de construção de soluções potenciais e na Tabela 2 estão também os vieses e

mitigações, mas neste caso para a avaliação das soluções potenciais. Em ambos os casos são dados alguns exemplos de organizações que utilizam estas mesmas mitigações.

Tabela 1 - Vieses, mitigações e exemplos relativos à construção de soluções potenciais (Bonabeau, 2009)

Vieses	Mitigações	Exemplos
<b>Viés do próprio serviço, procura confirmar suposições;</b>	Divulgação para obter uma grande diversidade de suposições;	Google, Affinova, InnoCentive, Threadless, Bell Canada's I.D.ah!,
<b>Interferências sociais provocadas pela influência dos outros;</b>	Aglomerção aditiva para obter participantes individuais;	ManyEyes, Swivel, Marketocracy, Goldcorp, Delicious, Digg, Procter + Gamble's Connect and Develop, Salesforce.com's Idea Exchange, Dell's IdeaStorm, Cajun Navy, Netflix's contest, blogs, wikis, Delphi method, lead-user tool kits, software open-source, motores de recomendação, fóruns de suporte
<b>Viés da disponibilidade, em o utilizador pensa que o problema pode ser resolvido com uma solução fácil;</b>	Divulgação para obter uma grande diversidade de soluções fáceis;	
<b>Viés da confidencialidade própria, em que se acredita ter encontrado uma solução prematuramente;</b>	Divulgação para obter uma diversidade de soluções;	
<b>Viés da âncora, em que é explorada a vizinhança de uma determinada âncora;</b>	Divulgação para obter uma diversidade de âncoras;	
<b>Crença na perseverança, em que o participante continua a acreditar, apesar das provas em contrário;</b>	Divulgação para obter uma diversidade de crenças;	
<b>Simulação, ou seja, quando só acredita numa solução quando a vê.</b>	Divulgação e auto-organização para obter diversidade de estímulos.	

Tabela 2 - Vieses, mitigações e exemplos para a avaliação de soluções potenciais (Bonabeau, 2009)

Vieses	Mitigações	Exemplos
Viés da linearidade, em que se procura uma relação de causa-efeito;	Auto-organização para obter iterações não-lineares;	Digg, HSX, Zagat, “American Idol”, Affinova, Threadless, Intrade, Google, StumbleUpon, Bell Canada’s I.D.ah!, Delicious, Mechanical Turk, Marketocracy, Salesforce.com’s Idea Exchange, software open-source, Delphi method, Mercados de informação e previsão.
Local versus Global, em que se confundem os efeitos locais com os globais;	Auto-organização para obter iterações não-lineares;	
Viés estatístico, em que se evitam as análises estatísticas;	Agregação aditiva para utilizar a legislação de largos números;	
Padrão da obsessão, em que se vêem padrões onde não existem;	Agregação aditiva e divulgação para obter diversidade de detetores de padrões;	
Conceção, ou seja, o participante é influenciado pela apresentação de soluções;	Agregação aditiva para obter diversidade de influências;	
Desconto hiperbólico, dominado pelo efeito de curto prazo;	Agregação aditiva para obter diversidade de escalas de tempo;	
Viés da doação, em que existe aversão ao risco e à perda;	Agregação aditiva para obter diversidade de perfis de risco.	

Para tomar decisões precisas é necessário uma dose elevada de processamento de informação e uma posterior avaliação das soluções potenciais. Ao longo dos anos estas tarefas deixaram de ser executadas unicamente por equipas e grupos de foco das empresas e passaram também a integrar a Inteligência Coletiva para apoiar estes processos. Podem ser distinguidas três formas de mitigação, ou seja, a divulgação, a auto-organização e a agregação aditiva. A divulgação permite aumentar o número de participantes envolvidos no processo de forma a identificar ideias diferentes e em maior quantidade. A agregação aditiva permite combinar e condensar informação obtida de muitos contribuidores. E a auto-organização permite interações para acrescentar valor às contribuições (Leimeister, 2010).

### **Inovação Aberta**

Outra área de aplicação da Inteligência Coletiva é a Inovação aberta, que se proporciona com a abertura das empresas e dos seus processos de inovação. As empresas podem envolver a Inteligência Coletiva e o potencial da inovação dos utilizadores da *Internet* em diversas fases do desenvolvimento do produto. Por exemplo, no caso da empresa LEGO, os consumidores podem contribuir para o desenvolvimento dos seus produtos através do desenho dos modelos. A Lego Digital Designer fornece kits para os utilizadores poderem fazer as suas contribuições. Para além desta, as empresas SAP, BMW e IBM proporcionam também ao consumidor uma possibilidade de contribuir para a conceção inovadora dos produtos e serviços. É muito importante para a sobrevivência das empresas, encontrar abordagens para capturar o conhecimento dos consumidores e apresentar ideias. Esta forma de proporcionar uma contribuição do consumidor no processo de desenvolvimento do produto pode ser uma estratégia interessante e com sucesso (Leimeister, 2010).

### **Crowdsourcing**

*Crowdsourcing* é baseado no conceito de *outsourcing* pois surgiu no contexto do *outsourcing* de determinadas atividades corporativas, realizadas por um grupo de pessoas independentes, a “*crowd*”. Esse grupo de indivíduos independentes realizam tarefas tais como resolver questões de investigação ou reconhecimento de padrões. Para além disso, esse grupo

de indivíduos realiza melhor essas tarefas ou de forma mais econômica do que as máquinas ou os mesmo os especialistas. Um bom exemplo de utilização de *Crowdsourcing* são os mercados de previsão, pois utilizam as opiniões e expectativas das grandes massas para prever probabilidades de ocorrências de eventos no futuro. Outro exemplo de aplicação do *Crowdsourcing* é o *Mechanical Turk* da *Amazon*, que permite às empresas expor as suas tarefas, e a um conjunto de pessoas, a possibilidade de as resolver por um preço mínimo (Leimeister, 2010).

### Colaboração Social

Na área de colaboração social a Inteligência Coletiva possui um grande potencial, pois a criação de valor é feita a partir de pequenas contribuições de um conjunto de indivíduos. Um exemplo muito popular de contribuição da Inteligência Coletiva nesta área é a *Wikipedia*. Esta é a maior enciclopédia de língua Inglesa, com mais de quatro milhões de artigos escritos nesta, apesar de possuir muitos temas traduzidos noutras línguas. Qualquer pessoa pode ter acesso aos artigos escritos nesta e para além disso qualquer pessoa pode criar artigos e editá-los. Apesar disso, esta é considerada por alguns autores de qualidade equivalente à Enciclopédia *Britannica*. Esta abordagem foi copiada por diversas empresas e indivíduos que criaram as suas próprias *Wikis*. Existem outras abordagens com o mesmo objetivo, tais como as plataformas de Partilha Social, que permitem aos seus utilizadores guardar, gerir e partilhar conteúdos. O cruzamento de referências e as categorias são garantidos a partir da utilização de *tags* (Leimeister, 2010).

#### 3.2.7 Vantagens

A utilização de Inteligência Coletiva pode trazer diversas vantagens para uma organização. Segundo Alag, 2009, os seus impactos podem ser:

**Taxas de retenção mais elevadas** - Quanto mais utilizadores interagirem com a aplicação, maiores são as probabilidades de que estes repitam as suas visitas.

**Melhores oportunidades de compra para o utilizador** – Quanto mais interações, maior será o número de páginas visitadas pelo utilizador, o que aumenta as oportunidades de compra ou de comunicação com o utilizador.

**Maior probabilidade de que um utilizador termine uma transação e encontre informação de interesse** – Quanto mais relevante contextualmente for a informação que o utilizador encontra, mais provável será que este encontre a informação que necessita para completar a operação ou encontrar informação de interesse. Isto traduz-se em mais cliques e taxas de conversão para publicidade.

**Impulsionar o motor de *ranking* de pesquisas** – Quanto mais contribuições os utilizadores fazem, mais conteúdo se torna disponível e este é indexado por motores de procura. Ou seja isto pode tornar a aplicação mais fácil para que outros a encontrem.

Para além destas vantagens, a utilização da Inteligência Coletiva pode apoiar a tomada de decisões complexas e de políticas. É um grande desafio, para não dizer que é quase impossível para quem toma decisões, reunir e compreender toda a informação necessária para tomar a decisão ótima ou suficientemente robusta. Isto deve-se à complexidade e às constantes mudanças, juntamente com os aumentos das capacidades tecnológicas. Devido a estas constantes mudanças, caso não se tomem as decisões mais acertadas, os impactos podem ser maiores no futuro do que eram anteriormente. São necessários sistemas que nos ajudem a perceber a situação de uma forma global, contudo é necessário que não haja sobrecarga de informação. Assim é necessário reunir apenas a informação relevante e aí surge a contribuição da Inteligência Coletiva. Esta permite de diferentes formas gerir conteúdos, organizar experiências, encontrar comentários e modificações nos documentos e ajudar na atribuição de prioridades (Glenn, 2013).

### 3.2.8 Utilizadores

#### 3.2.8.1 Utilizadores e Itens

De uma forma geral as aplicações podem ser geralmente compostas por utilizadores e itens. Os utilizadores interagem com os itens, que por sua vez podem possuir meta-dados. Estes compram, contribuem, recomendam, vêem, colocam *tags*, avaliam, guardam e marcam itens. Artigos, vídeos, fotografias, publicações em *blogs*, perguntas e respostas publicadas em fóruns de discussões, ou mesmo produtos ou serviços à venda na aplicação podem ser considerados itens, pois são entidades de interesse. Para além destes podemos também considerar os

utilizadores como itens, se estivermos perante uma aplicação que seja uma rede social. Cada item pode ter meta-dados associados, que por sua vez se podem apresentar sob a forma de palavras-chave geradas profissionalmente, ou extraídas através de algoritmos após análise do texto, *tags* gerados pelos utilizadores, avaliações, *rankings* de popularidade, ou qualquer outro elemento que possa fornecer informação que interligue os itens. Os meta-dados associados a um item podem ser definidos como um conjunto de atributos que permitem caracterizar ou qualificar um item. Existem três fontes principais para gerar meta-dados de um item: baseada nos atributos, baseada no conteúdo ou baseada nas ações dos utilizadores (Alag, 2009).

Relativamente à geração de meta-dados baseada nos atributos, pode dizer-se que é realizada a partir da visualização dos atributos dos utilizadores ou dos itens. Estes atributos podem ser muito variados dependendo do domínio em que a aplicação se encontra. Por exemplo, se considerarmos um utilizador como um item, no caso de uma rede social, os seus atributos podem ir desde a sua idade, à sua profissão passando por diversas informações pessoais. Para itens que não sejam utilizadores, os atributos podem ser por exemplo o preço, fabricante, disponibilidade geográfica, data de fabrico, caso se trate de um produto ou serviço (Alag, 2009).

Os meta-dados podem também ser gerados a partir de conteúdos, ou seja, através da análise do conteúdo de determinado documento. Esta abordagem extrai os meta-dados a partir de texto não estruturado, através de *text mining* e retorno de informação. Dados como o título, subtítulo, palavras-chave, número de vezes que determinadas palavras surgem ao longo do documento, entre outra informação, podem ser convertidos em meta-dados para um dado item. Por último, os meta-dados podem ser gerados a partir da análise das interações dos utilizadores. Os utilizadores realizam interações que podem ser bastante explícitas em relação aos seus interesses ou preferências. Exemplos destas são a compra de determinado produto, avaliação de um item, etc. Contudo nem todas as interações são fáceis de compreender. Os meta-dados podem ser vistos como um vetor de atributos associado a cada item ou utilizador, que depois de obtidos podem ser comparados e medidos (Alag, 2009).



### 3.2.8.2 Interações dos Utilizadores

As técnicas utilizadas para obter a informação desejada a partir das interações dos utilizadores, são segundo Alag (2009):

**Histórico de transações** – Lista de itens comprados pelo utilizador, lista de itens favoritos, ou itens que se encontram atualmente no carrinho de compras;

**Conteúdo visitado** – Conteúdo procurado ou lido pelo utilizador, publicidade em que este *clicou*;

**Caminho seguido** – Como é que o utilizador conseguiu chegar a determinado conteúdo, qual o objetivo deste;

**Seleções no perfil** – Escolhas que os utilizadores realizam quando selecionam os perfis, por exemplo o aeroporto escolhido por defeito quando este visita uma aplicação de viagens;

**Feedback de votações e perguntas** – Quando um utilizador responde a votações ou questões *on-line*;

**Avaliações** – Avaliação de determinado conteúdo;

**Atribuir Tags** – Associar *tags* a determinados itens; e

**Votar, guardar, *bookmarking*** – Expressar interesse em determinado item.

Para extrair a inteligência da interação do utilizador em determinada aplicação é necessário ir além do conteúdo que este visitou ou visualizou, é necessário quantificar e qualificar a sua interação. Para saber se um utilizador gostou ou não de determinado artigo, é necessário ter uma forma de quantificar o quanto este gostou do artigo em comparação com outros. Para uma melhor compreensão destas interações segue-se uma distinção das formas de interações dos utilizadores segundo Alag (2009):

**Votação e avaliação** – Para obter uma avaliação do item, a aplicação pode pedir diretamente ao utilizador o seu *feedback* em relação a determinado item, assim a informação obtida é quantificável e pode ser usada diretamente sem necessidade de análise. Geralmente a avaliação realizada pelos utilizadores é positiva pois estes tem tendencia a avaliar somente os itens consumidos ou com que interagiram e normalmete só consomem ou interagem com itens de que gostam. A escala mais utilizada para avaliações é de zero a cinco, contudo podem

permitir também a utilização de votações e nesse caso a escala será de um, caso o voto seja positivo ou menos um caso seja negativo. Permitir aos utilizadores atribuir as suas votações, é outra forma de obter mais informação destes.

**Enviar um *Email* ou Reencaminhar um *link*** – É muito frequente encontrar nos *websites* uma forma de partilhar conteúdos ou mandar *email* com determinada página a outros. A partilha de conteúdo realizada por um utilizador a outro pode ser considerada como um voto positivo atribuído pelo utilizador.

**Guardar e *bookmarking*** – Existem serviços *on-line* de *bookmarking*, que permitem aos seus utilizadores, guardar e recuperar *URLs*, também conhecidos como *bookmarks*. Para além de guardar a sua lista de interesses, os utilizadores podem também ver outros *URLs* interessantes que foram guardados por outros utilizadores. Quando um utilizador guarda um *URL*, está a expressar interesse pelo material contido neste. Algumas aplicações permitem ainda a criação de pastas, que são coleções de *URLs* que podem ou não estar relacionados, ou seja, pode possuir meta-dados associados. Assim uma pasta é também um item, pois pode ser partilhado, marcado e avaliado.

**Compra de Itens** – A avaliação de um item pode ser também concluída a partir da resolução de uma compra, ou seja, se um item é comprado demonstra um voto de confiança, o que pode ser considerado positivo. Caso o item seja devolvido, demonstra um voto negativo. Para além disso utilizadores que comprem itens semelhantes podem receber uma lista de recomendações construída com base no que os outros compraram. Ou seja, se utilizador tem gostos em comum com os de outro, as suas compras do primeiro serão utilizadas como sugestão para as compras do segundo.

**Fluxo de cliques** – Quando o utilizador é apresentado com um conjunto de itens, é muito provável que este vá *clicar* em pelo menos um, baseando-se no título ou descrição deste. Contudo após visualizar rapidamente o conteúdo pode aperceber-se que aquela informação não lhe interessa ou não lhe é útil. Uma forma simples de quantificar a relevância de um artigo, utilizada por exemplo pela *Google News* para personalizar o seu *website*, é gravar um voto positivo por cada vez que um item é *clicado*. Isto inclui as vezes que o utilizador saiu da página porque o conteúdo não lhe interessou, para resolver esta situação pode verificar-se o tempo que este ficou na página. Esta informação pode não ser infalível, pois o utilizador pode deixar a

pagina aberta sem estar a olhar para ela, contudo de forma geral esta informação deve ser considerada útil.

**Avaliações** – Os gostos e opiniões são geralmente expressos através de avaliações e recomendações que vão ter um impacto posterior noutros utilizadores, especialmente quando quem expressa a opinião é imparcial, tem uma opinião similar à de quem lê ou é uma pessoa com influência. Desta forma adicionamos mais uma entidade, ou seja o revisor, e a sua relação com o item, a revisão. As revisões encontram-se sob a forma de texto não estruturado o que implica a utilização de mecanismos de procura.

### 3.2.8.3 Identificar e motivar os utilizadores

Existem dois pontos muito importantes aquando da aplicação das técnicas de Inteligência Coletiva, ou seja, identificar os indivíduos ou grupos cuja inteligência ou experiência pode fazer parte da “*crowd*” e fornecer as motivações necessárias à sua participação (Lesser et. al,2012).

Todos os métodos de Inteligência Coletiva necessitam de um número mínimo de participantes ativos para gerar perceções com o valor necessário. Segundo Lesser, et. al (2012), são necessárias pelo menos vinte ou trinta pessoas na comunidade para que haja as interações necessárias para manter a motivação do grupo. Para gerar ideias e idealização de eventos, são necessárias centenas ou milhares de indivíduos. Dependendo da área o número mínimo necessário de indivíduos vai variando. O número apropriado de indivíduos é baseado em três fatores importantes como se pode ver na Figura 5.

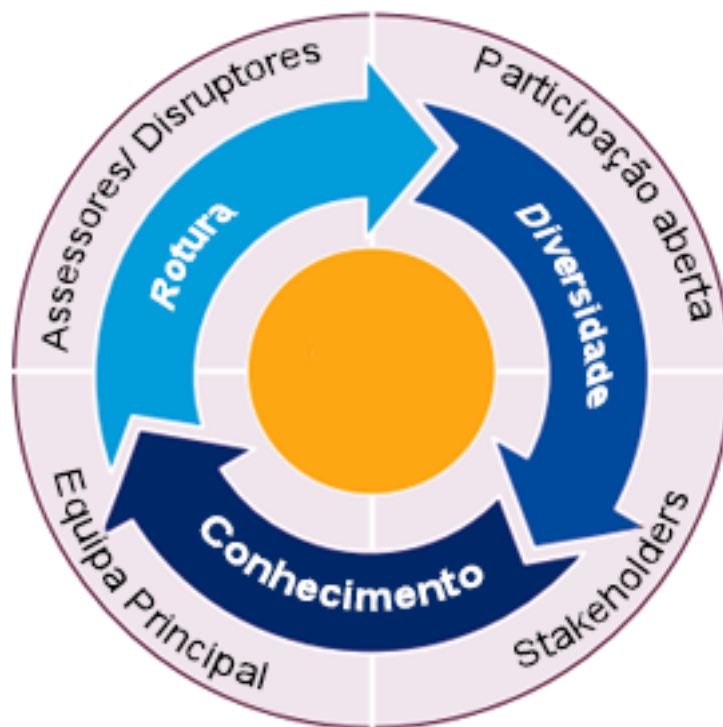


Figura 5 - Três fatores que influenciam o número mínimo de participantes, adaptado de (Lesser et. al,2012)

Esses fatores são então o conhecimento, a diversidade e a rotura. Um dos fatores é o conhecimento pois pode ser necessário um conhecimento contextual e uma familiaridade com o tema a explorar para se poder criar uma opinião informada e uma perspectiva sobre o mesmo. Sem esta perspectiva sobre o assunto a participação dos intervenientes pode ser baseada somente em especulações. Outro fator importante é a diversidade, pois caso não haja participantes suficientes pode não existir variedade de perspectivas o que faz com que o assunto possa ser abordado de um único ponto de vista. Relativamente à rotura, pode ser vantajoso ter elementos que consigam contribuir com perspectivas de rotura, que provoquem um pensamento inovador. Estes indivíduos são importantes pois caracterizam-se por questionar suposições com pensamentos independentes e geralmente são pessoas com visão de futuro. Apesar de uma grande diversidade de ideias, estes indivíduos podem apresentar mais vantagens pois estes são capazes de criticar o pensamento de grupo e conduzir a soluções mais profundas e completas (Lesser et. al, 2012).

Depois de identificados os potenciais participantes, devem ser utilizadas formas de encorajar os mesmos a partilhar a sua experiência ou conhecimento. Esta fase requer uma

articulação clara sobre a importância do contributo que o utilizador realiza e é desejável incorporar motivadores extrínsecos como dinheiro ou medindo a performance e intrínsecos como a satisfação, lealdade e divertimento pessoal (Lesser et. al, 2012). Na Figura 6 podem ver-se as técnicas de motivação dos utilizadores assim como a sua fonte, natureza e complexidade.

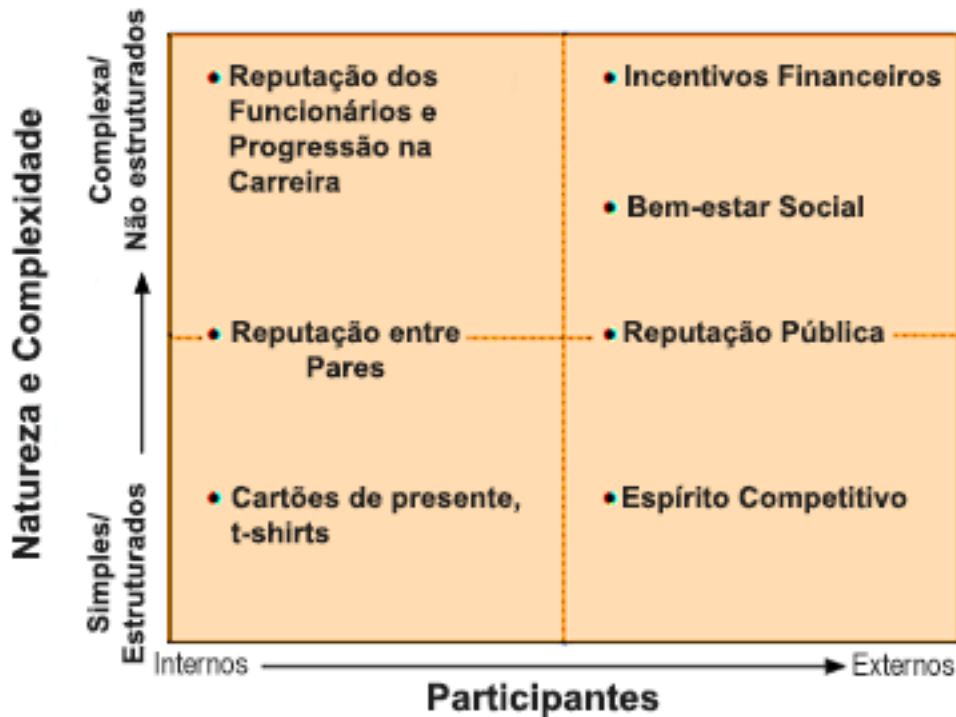


Figura 6 - Técnicas de Motivação, adaptado de (Lesser et. al, 2012)

Segundo Alag (2009), muitos sítios fornecem incentivos monetários para motivar o seu público. O valor dos incentivos monetários vai desde uns cêntimos por transação até milhões de dólares, como no caso da *Netflix Challenge*, que propõe aos participantes o desenvolvimento de um algoritmo de recomendação de filmes por um milhão de dólares (Lesser et. al, 2012).

Também os incentivos para um espírito competitivo, a reputação a pares, os objetivos comuns ou trabalho para ajudar outros podem ser considerados como motivadores para contribuição para a Inteligência Coletiva (Lesser et. al, 2012). Os avaliadores apreciam o reconhecimento e gostam de contribuir para poder obtê-lo. Muitas aplicações apresentam uma lista dos melhores avaliadores, ou seja, uma espécie de *ranking*. Por vezes as aplicações até

realçam alguns dos seus melhores avaliadores para se gerar um incentivo à partilha e contribuição (Alag, 2009). Segundo o autor, Lesser et al. (2012) os jogos têm um grande potencial para serem usados na inteligência Coletiva, pois estes conseguem motivar os seus participantes através da atribuição de prémios, simulando desafios e vitórias épicas.

Motivar participantes internos, tais como empregados, pode ser mais desafiante, pois os participantes externos são facilmente motivados através de cartões, *t-shirts*, e outras pequenas ofertas (Lesser et. al, 2012).

Uma das razões que move mais os utilizadores a participar e partilhar a sua experiência é a necessidade de ser descobertos por outros (Alag, 2009). Os participantes são muitas vezes motivados pela visibilidade que as suas contribuições lhes provocam ou até mesmo a procura por oportunidades de carreira. Por vezes somente a ideia de que a sua opinião está a ser analisada por pessoas responsáveis por tomar decisões já motiva determinados contribuidores. As empresas dão valor aos empregados com conhecimento mas também àqueles que o partilham com outros e os indivíduos cada vez mais têm noção dessa realidade (Lesser et. al, 2012).

Para além destes, a confiança é um fator chave para motivar os participantes e pode ser demonstrada de várias formas: a sugestão de ideias é respeitada e o esforço valorizado; concordar com a partilha de propriedade intelectual é também valorizada; *feedback* dos participantes é reconhecido (Lesser et. al, 2012). Segundo Alag (2009), os utilizadores têm mais tendência a fazer contribuições em sítios com grandes audiências, pois apresentam um nível de confiança elevado.

### 3.2.9 Atribuição de *Tags*

De uma forma geral as aplicações permitem aos seus utilizadores uma navegação pelo conteúdo através de categorias ou menus, o que se pode revelar para o utilizador um pouco entediante tentar chegar ao conteúdo de interesse. Cada categoria possui subcategorias que o utilizador vai seleccionando até se aproximar do conteúdo desejado. Isto pode ser entediante para o utilizador pois pode ser necessário navegar através de vários subtópicos para obter a informação que deseja. Para além disso, existem sítios que mostram o quão absurdo seria ter

que categorizar todos os seus elementos, como é o caso do *Flickr*. Com os seus milhões de fotografias seria demasiado custoso ter que categorizar cada uma manualmente. Aquilo que é exposto nesta secção está descrito na referência (Alag, 2009).

Uma alternativa a esta é a navegação através de um conjunto de *links* ou *hyperlinks* dinâmicos, construídos com base em texto similar ao que o utilizador se mostrou interessado. Através deste método o utilizador atribui rótulos ou *tags* aos itens, faz *bookmarking* destes, partilha-os e visualiza-os e evita a necessidade de colocar profissionais a categorizar os itens da aplicação. Em suma, a inteligência pode ser encontrada sob a forma de itens relacionados com outros que se encontram rotulados ou com *tags*, sob a forma de conexões com outros utilizadores que possuem itens rotulados similares, ou mesmo mostrando *tags* alternativos que foram associados a um item de interesse e que permitem mostrar itens relacionados.

Atribuir *tags* é o processo de adicionar texto, palavras ou pequenas frases aos itens. Estes *tags* podem ser associados a itens tais como os utilizadores, fotografias, artigos, produtos, publicações nos blogs, *podcasts*, vídeos, entre outros. O conjunto de termos ou tags usados na aplicação definem o vocabulário desta. E tendo em conta o seu contexto e a quem estes aparecem, estes termos podem ser usados como *links* de navegação dinâmica, como por exemplo as *Tag Clouds*.

Assim sendo, segundo o autor Alag, 2009, a utilização de *tags* numa aplicação permite: a construção de um vetor de termos para associar aos utilizadores e itens; a construção de *links* de navegação dinâmicos como a *Tag Cloud* e os textos com hiperligações; o uso de meta-dados para personalizar e relacionar os utilizadores; a construção de um vocabulário para a aplicação; e a marcação de itens para poderem ser partilhados com outros utilizadores.

Os meta-dados podem ser baseados no conteúdo ou na colaboração, sendo o primeiro obtido através da análise do conteúdo do item e o segundo através das ações realizadas pelo utilizador. Os meta-dados dos utilizadores ou itens obtidos através das *tags*, são simples vetores de termos com os seus pesos associados.

Segundo Alag, 2009, as *tags* podem ser divididas em três categorias tendo em conta quem as gerou. Estas podem ser geradas profissionalmente, geradas pelo utilizador ou por máquinas. As *tags* geradas profissionalmente, são realizadas por especialistas daquele domínio e tem

características como: mostrar os conceitos relacionados com o texto; obter o valor semântico, o que pode levar à utilização de palavras não encontradas no texto; podem fornecer uma visão mais global, fora da área de interesse; podem apresentar sinónimos; possuir frases com mais que um termo; podem controlar o vocabulário para possuir apenas um conjunto possível de termos.

Quando existe uma grande quantidade de conteúdo a ser continuamente gerado é melhor realizar uma categorização por parte dos utilizadores. A criação de *tags* por parte dos utilizadores é uma boa forma de utilização da Inteligência Coletiva pois o poder coletivo destes é aproveitado. Geralmente as *tags* geradas por utilizadores possuem as seguintes características: utilização de termos familiares ao utilizador; estes realçam os conceitos relacionados com o texto; encontram também o valor semântico associado, utilizando palavras que não se encontram no texto; definem frases com vários termos; fornecem informação com valor sobre o utilizador e o item; necessitam de ser verificados devido à diversidade de termos ou formas que estes possuem para se referir a uma coisa, como por exemplo os plurais.

As *tags* geradas por máquinas são realizadas a partir de algoritmos automatizados, que realizam uma análise do texto e encontram termos e frases. Estas possuem geralmente as seguintes características: utilização de termos do texto, à exceção dos sinónimos; utilização de termos com uma palavra; possibilidade de gerar *tags* com ruído, pois é possível que esta possua diferentes significados em diferentes contextos;

A razão para que os utilizadores coloquem *tags* é para poder organizar os itens. Para além disso também o fazem para partilhar informação, encontrar itens relacionados que outros utilizadores tenham rotulado e quando querem ser encontrados por outros. Os utilizadores conseguem aceder a informação similar à que estão a ver e que lhes interessa e para além disso são muitas vezes confrontados com uma diversidade de termos semelhantes ou relacionados com esse conteúdo, que lhes permite aumentar a amplitude da procura.

As *Tag Clouds* e os *hyperlinks* dentro dos conteúdos são exemplos de navegação dinâmica. Existem aplicações que utilizam os *hyperlinks* nas cidades, números de telefone e palavras-chave, como por exemplo o *Gmail* e *Yahoo!*. Uma *Tag Cloud* é uma lista de *tags*, geralmente ordenada alfabeticamente, que apresenta cada *tag* com tamanho de fonte de acordo com a frequência de utilização. Ou seja, quanto mais usada for uma *tag* maior será o tamanho da sua



fonte. Por vezes, as *tags* podem aparecer também com cores diferentes. A Figura 7 pode-se ver como se forma uma *Tag Cloud*, através da combinação de vetores de termos. Para se influenciar corretamente as *tags* é necessário tratar das palavras com plural, detetar frases com mais que um termo e lidar com *tags* sinónimas. Quando um utilizador *clica* numa *tag* da *Tag Cloud*, obtemos um contexto para essa *tag*. Este contexto pode ser aproveitado e introduzido no mecanismo de recomendações, onde se realizam procuras de conteúdo relacionadas com o utilizador ou *tag* de interesse.

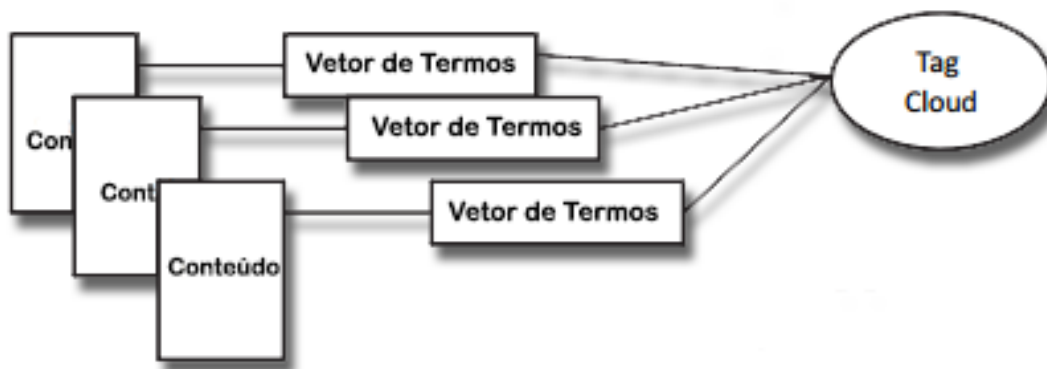


Figura 7 - Tag Cloud, adaptado de (Alag, 2009)

Tal como o sistema de *ranking* do Google verifica o número de *links* para determinada página como métrica para quantificar a importância desta, as *tags* trabalham com uma métrica similar. Ou seja, caso um artigo esteja a ser rotulado por muitos utilizadores com a mesma *tag*, é bastante provável que este seja relevante para determinado tópico e que seja interessante para outros utilizadores. Na Figura 8 mostra-se, a combinação do contexto das *tags* com os metadados do utilizador para realizar uma *query* no motor de procura que devolve os resultados relevantes para a situação.

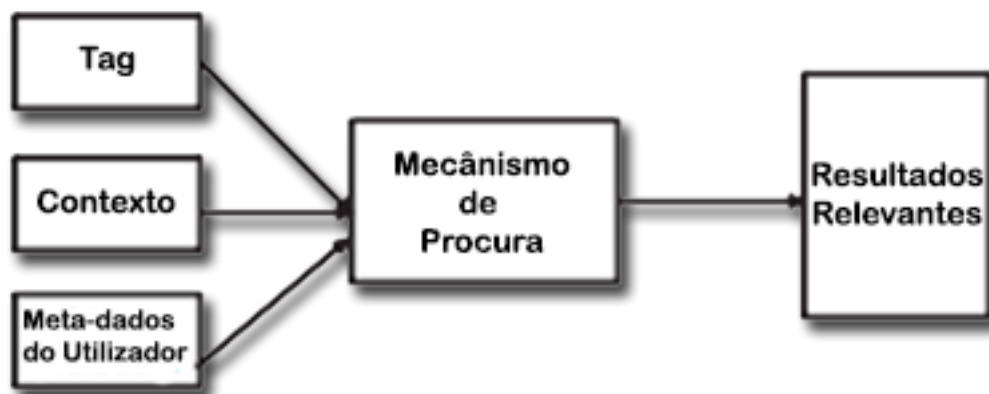


Figura 8 - Utilização de Tags, contexto e meta-dados para obter os resultados relevantes, adaptado de (Alag, 2009)

*Folksonomia* é o processo de classificação que permite aos utilizadores recuperar informação usando termos com que estão familiarizados. Esta faculta aos utilizadores uma forma de encontrar outros com interesses semelhantes aos seus e ver o conteúdo pelo qual se interessam. As *folksonomias* são criadas a partir das *tags* geradas pelos utilizadores (Alag, 2009). Segundo Michalsky et. al, (2010) esta é a sua grande vantagem, pois é utilizada a linguagem dos próprios utilizadores e proporciona procuras mais precisas por parte dos outros utilizadores. Segundo Herzog et. al, (2007) *folksonomia* é apenas um conjunto de termos, que do ponto de vista matemático pode ser um grafo triplice com três arestas, o utilizador, a *tag* e o item.

### 3.2.10 Tipos de Conteúdo

O conteúdo pode ser desenvolvido por profissionais, pelos próprios utilizadores, como se viu anteriormente, ou através de *sites* externos. Aquilo que é exposto nesta secção está descrito na referência (Alag, 2009) Estes podem ser: artigos, produtos, termos de classificação, *blogs*, *wikis*, grupos e fóruns de mensagens, fotografias e vídeos, votações, termos de procura, páginas de perfil, ferramentas e folhas de cálculo, registos em chats, classificações, publicidade e listas como se pode ver na Tabela 3.

Tabela 3 - Tipos de Conteúdo, adaptado de (Alag, 2009)

Tipos de Conteúdo	Descrição	Fontes
Artigos	Texto sobre determinado tópico, com título, corpo e possíveis subtítulos	Gerado profissionalmente, pelo utilizador, <i>feed</i> de notícias, agregado de outros sites
Produtos	Item a ser vendido. Possui, geralmente, título, descrição, palavras-chave, avaliações, classificações, preço, produtor e disponibilidade geográfica	Gerado pelo site e pelo utilizador
Termos de classificação	Termos com palavras-chaves ou <i>tags</i> associadas, criadas para navegação por parte do utilizador	Gerado profissionalmente, por mecanismos, e pelo utilizador
Blogs	Diários pessoais <i>on-line</i> , onde o utilizador partilha informação sobre temas que considera interessantes para outros lerem e comentarem	Gestão do site, empregados da empresa ou gerado pelo utilizador
Wikis	Ferramenta de colaboração <i>on-line</i> onde os utilizadores podem facilmente criar, editar ou apagar páginas	Gerado quase sempre pelo utilizador
Grupos e fóruns de mensagens	Locais onde se colocam perguntas, ou se respondem a estas. As respostas podem ser classificadas como uteis ou não	Gerado quase sempre o utilizador, contudo podem ser obtidas respostas de utilizadores especialistas que trabalham no site
Fotografias e vídeos	Formas de comunicação sob a forma de fotografias e vídeos	Gerados profissionalmente e pelo utilizador

Votações	O utilizador é confrontado com questões às quais tem que responder com opções já definidas	Gerados profissionalmente ou pelo utilizador
Termos de procura	Questões procuradas pelo utilizador, similar à classificação dinâmica	Gerados pelo utilizador
Páginas de perfil	Páginas de perfil dos utilizadores geralmente criadas pelos próprios, onde colocam informação sobre eles mesmos	Gerados pelo utilizador
Ferramentas e folhas de cálculo	Ferramentas e folhas de cálculo disponibilizadas pelos próprios sites	Gerados profissionalmente
Registos em chats	Transcrições de conversas nos chats <i>on-line</i>	Especialistas comunicam com utilizadores e os utilizadores comunicam entre si
Classificações	Avaliações de itens	Gerados profissionalmente ou pelo utilizador
Anúncios	Publicidade com título e corpo da mensagem, com possíveis palavras-chave associadas	Gerados profissionalmente ou pelo utilizador
Listas	Lista de itens combinados entre si	Gerados profissionalmente ou pelo utilizador

De entre estas formas de conteúdo, os blogs, as *wikis* e os fóruns de mensagens estão geralmente associados à Inteligência Coletiva, pois são gerados pelos utilizadores através das suas contribuições e opiniões.

## Blogs

*Blog* é a abreviatura de *weblog* corresponde a simples diários *on-line*, onde as pessoas escrevem sobre temas que desejam partilhar com outros. Estes por sua vez podem ler e comentar os tópicos abordados no *blog*. A popularidade de um *blog* pode ser vista a partir do número de vezes que este é referenciado nas publicações de outras pessoas e pelo número de blogs que possui na blogosfera. Estes surgem em três contextos diferentes: em *websites* corporativos, dentro de aplicações e outros. Os *blogs* nas empresas são utilizados para conectar os seus consumidores, acionistas, funcionários e outros. Este pode ser um bom mediador para a colaboração dentro da empresa, sobretudo quando esta se encontra geograficamente dispersa. Os *blogs* surgem também num contexto das aplicações, para permitir aos utilizadores obter e fornecer informação relevante sobre estas. Os outros *blogs* contidos na blogosfera podem ter um impacto nas marcas. Enquanto uma publicação favorável pode aumentar a publicitação de determinada marca, um comentário negativo pode levar á destruição da imagem desta.

## Wikis

As *wikis* potenciam uma boa colaboração *on-line*, pois permitem aos seus utilizadores a criação edição e eliminação de páginas. Podendo ter revisões das suas modificações, permite o retrocesso ao estado anterior. Estas podem ser utilizadas para colaboração *on-line*, na mobilização para a contribuição dos utilizadores, no impulsionamento de *ranking* de mecanismos de procura e como repositório de conhecimento. Quando são utilizadas para colaboração *on-line* entre grupos, estes podem partilhar os seus pensamentos e informação de interesse, pois qualquer elemento com a acesso à *wiki* poderá ver o conteúdo. Quando são utilizadas na mobilização de utilizadores para contribuição, estas podem ser úteis por exemplo no desenvolvimento de páginas com as perguntas mais frequentes, guia de instalação, guia de ajuda, entre outros. Uma *wiki* pode ajudar a impulsionar a visibilidade de uma aplicação. E pode servir como repositório de conhecimento sendo acessível a todos e proporcionando uma forma de procura e retorno de informação fácil.

## Fóruns de Mensagens

Os fóruns de mensagens consistem em locais virtuais onde se podem colocar questões e outros podem responder as estas, permitindo atribuir avaliações de utilidade a cada uma. Estes

fóruns estão, geralmente, associados a grupos de indivíduos, que partilham interesses. Estes fóruns são vantajosos na medida em que juntam pessoas com interesses e gostos semelhantes para resolver questões que podem facilmente ser respondidas por outros elementos.

### 3.2.11 Blogosfera e Web Crawling

É importante encontrar a informação que outros partilham sobre o nosso produto ou aplicação pois esta tem um grande impacto na marca. Existem empresas que se dedicam à procura de *blogs*, fornecendo as *APIs* necessárias para realizar as *queries* relevantes. A informação partilhada por milhões de pessoas é uma excelente forma de Inteligência Coletiva que pode fornecer informação muito interessante. De forma similar, a blogosfera fornece um conjunto de conteúdos agregados e que são expostos quando se tornam relevantes. Para encontrar este conteúdo relevante, são necessárias duas etapas: agregar ou encontrar o conteúdo e determinar se este é relevante ou não. Aquilo que é exposto nesta secção está descrito na referência (Alag, 2009)

Para poder manter os artigos e publicações sobre determinados temas em diversos sítios na *Internet*, atualizados é necessário realizar *updates* frequentes. A maior parte destes publica os seus conteúdos no formato RSS que permite partilhar conteúdos na *Web* em formato XML, e encontrar outros sites com *updates* que usem um formato similar. Fazer pesquisas na blogosfera implica quatro passos: criar a *query* de procura, enviá-la a um fornecedor de procura de blogs num formato em que este a possa entender, analisar a resposta e converte-la num formato standard.

Contudo dada a quantidade de informação disponível na *Internet*, como é possível encontrar apenas a informação que nos interessa? *Web Crawling* permite obter essa informação. O *Web crawling* é um programa ou um processo automatizado de visita de páginas *Web* com o objetivo de retornar conteúdos, extrair *URLs* para outros *links* interessantes, e visitá-los caso seja possível. Este conteúdo pode encontrar-se em diversas formas como texto, imagens ou vídeos. Os motores de busca, como são exemplos o Yahoo! e o Google, fazem '*crawling*' da *Web* para indexar novos conteúdos disponíveis. Os *Web crawlers* são geralmente conhecidos como *Web spiders*, *bots* ou indexadores automatizados.

O objetivo principal do *Web crawling* é então a recolha de dados de sítios externos e estes são geralmente utilizados para: agregação de conteúdos e indexação de conteúdos externos; procura de informação específica; acionar eventos; deteção de *links* corrompidos; procura de violações dos direitos de autor. Os passos essenciais do funcionamento de *Web crawling* são: fazer *seeding* do *Web crawler*, com um conjunto de *URLs* para visitar; verificar critérios de paragem, para que o *Web crawler* possa parar de retornar *URLs* quando este deteta determinada situação; obter o próximo URL a visitar, que possua permissões de visualização por parte dos *crawlers*; retornar o conteúdo dos *URLs* visitados sem duplicar informação; verificar se o conteúdo é relevante, podendo ser um ponto opcional, contudo permite verificar se os conteúdos encontrados se encaixam no modelo de informação relevante procurada; extração dos *URLs*, dos conteúdos que foram considerados relevantes e devem ser colocados na lista de *URLs* a verificar; injetar atrasos, em situações onde pode ser conveniente, como a utilização de várias *threads*, ou quando o processo é demasiado rápido.

Os *Web crawlers* são geralmente combinados com uma biblioteca de procura, utilizada na indexação e procura de conteúdo. Devido à constante modificação, criação e eliminação das páginas *Web*, os *crawlers* são realizados periodicamente para manter os conteúdos atualizados. Um *crawler* pode deparar-se com alguns desafios ao longo da sua utilização, sendo um deles a '*spider trap*'. Esta é criada para guardar os sites contra *crawlers spam*, que utilizam métodos para enganar os *crawlers* e substituir os motores de *ranking* de procuras pelos seus.

O *crawling* focado geralmente utiliza uma lista de *URLs* para visitar e atribui uma medida de relevância ao conteúdo. *Crawling* focado baseia-se no princípio de que quanto mais relevante um site é para um tópico de interesse, maior a probabilidade de que as páginas com ligação a este contenham informação relevante. Portanto é interessante verificar também estas páginas. Contudo fazer *crawling* a toda a *Web* implica custos significativos, que podem ir desde *software* e *hardware*, até à necessidade de uma rede de *Internet* de acesso rápido, passando pela necessidade de possuir dispositivos de armazenamento e administração destas infraestruturas.

### 3.2.12 Fatores-Chave e de Sucesso

Para obter um sistema de Inteligência Coletiva com sucesso é necessário identificar os fatores essenciais para que isso seja possível. Segundo Bonabeau (2009) e Leimster (2010), estes fatores são:

**Controlo** – Ao aplicar a Inteligência Coletiva numa empresa, a estrutura desta pode mudar drasticamente, o que pode levar a uma perda de controlo. A estrutura que era hierarquicamente fechada passa a ser aberta, com processos a serem realizados por *outsourcing*. Uma das maiores questões relacionados com o controlo é o incluir ou não de pessoas estranhas ao processo, pois a existe informação sobre a organização que será exposta para qualquer um ver. Esta informação pode não ser vista pelas pessoas com as melhores das intenções, o que pode trazer consequências para a empresa. Contudo partilhar informação com uma grande diversidade de pessoas e com experiências diferentes para partilhar pode produzir resultados mais valiosos quando executados corretamente (Bonabeau, 2009).

**Diversidade versus conhecimento aprofundado** – Para cada tarefa é necessário manter um balanço entre a diversidade e o conhecimento aprofundado do coletivo. A diversidade a mais pode levar a uma infinidade de respostas e ideias, contudo estas soluções podem ser viáveis e interessantes. As empresas devem decidir quais as pessoas que querem ver envolvidas nos projetos, tendo em conta a capacidade destes para compreender os problemas e contribuir coletivamente para resolve-los. Existem aplicações que exigem a participação de um grande número de indivíduos para garantir a qualidade do resultado, contudo é necessário que estes possuam o conhecimento necessário para poder contribuir com sugestões úteis. Por outro lado, existem decisões que necessitam de ser tomadas com base em muita experiência que pode ser encontrada em poucas pessoas. Em situações deste género podem ser utilizadas ferramentas como o método Delphi para poderem ser obtidas soluções com qualidade.

**Compromisso** – Para incentivar os indivíduos a participar, incentivos monetários por si só não são suficientes. Incentivos como o altruísmo, realização pessoal e identificação com o grupo podem ser incentivos importantes. Quando uma empresa procura soluções for a da empresa é necessário determinar se e como vai assumir a propriedade sobre os resultados obtidos.



**Fiscalização** – Quanto mais utilizadores envolvidos num projeto, maior a probabilidade de aparecerem comportamentos maliciosos. Para resolver casos destes, podem ser atribuídos castigos.

**Propriedade Intelectual** - Se há criação de soluções e ideias, é necessário que a empresa possa adquirir a propriedade intelectual destas. Quando os indivíduos se envolvem nestes projetos de Inteligência Coletiva devem reconhecer a perda da posse sobre esta.

Segundo Lesser et. al (2012), para o sucesso na área da Inteligência Coletiva, são necessários esforços como: incorporar fontes de resistência como desafios operacionais, conflitos nos contratos existentes, perda de controlo percebida e troca de papéis e responsabilidades; integrar a Inteligência Coletiva no ambiente de trabalho cultural e tecnológico; agir sobre as descobertas, comunicando o valor e resultados para o indivíduo e organização.

Para verificar a eficiência da Inteligência Coletiva em diferentes tipos de aplicações é necessário ter em conta várias métricas e indicadores-chave como podemos verificar na Tabela 4 (Bonabeau, 2009).

Tabela 4 - Métricas e indicadores de sucesso em diferentes aplicações (Bonabeau, 2009)

Tipo de Aplicação	Exemplos	Métrica-chave	Indicador-chave
Investigação, Desenvolvimento e Inovação	InnoCentive; TopCoder; Netflix;	Qualidade das soluções; Consistência do <i>output</i> ;	Acesso ao talento; Diversidade de participantes; Compromisso dos participantes ao longo do tempo;
Estudos de Mercado	Mechanical Turk; Affinnova;	Capacidade de descoberta; Descobrir respostas verdadeiras;	Tamanho exemplo e se é representativo do mercado; Compromisso dos participantes;
Previsões	Mercados Informativos como <i>Intrade</i> e <i>NewsFutures</i> ;	Precisão das previsões;	Capacidade de localizar quantidades reais; Compromisso dos participantes ao longo do tempo;

Serviços Customizados	Comunidades de utilizadores;	<p>Porcentagem de problemas resolvidos;</p> <p>Tom das conversações;</p> <p>Descoberta antecipada dos problemas;</p>	<p>Capacidade de resposta a problemas não resolvidos;</p> <p>Compromisso dos participantes;</p>
	Wikis;	<p>Qualidade, precisão e frequência das contribuições;</p> <p>Utilização do <i>output</i> em situações reais</p>	<p>Acesso a conteúdo somente de leitura;</p> <p>Compromisso e atividade dos participantes;</p>
	NIST; Peer-to-Patent	<p>Número, qualidade e âmbito de questões não esperadas e não abrangidas;</p>	<p>Progresso dos testes;</p> <p>Compromisso dos participantes;</p>
	Cajun Navy	<p>Acesso a informação difícil de obter;</p> <p>Minimização dos danos causados pela crise;</p>	<p>Formas de comunicação utilizadas;</p> <p>Compromisso dos participantes;</p>

### 3.3 Crowdsourcing

#### 3.3.1 Definição

Na última década o *Crowdsourcing* surgiu como um paradigma de recolha de informação e uma forma de resolver problemas na *Web* (Doan et. al, 2011). Este descreve um modelo de

negócio baseado na *Web* que proporciona soluções criativas através de uma rede de indivíduos que colocam as suas propostas (Brabham, 2008).

O *Crowdsourcing* baseia-se na exposição de uma “função”, por parte de uma empresa ou instituição, que anteriormente seria realizada por trabalhadores desta, a uma rede indefinida de pessoas na forma de convite aberto. Esta pode ser realizada de forma colaborativa ou individualmente, o único pré-requisito é a utilização de um concurso público/aberto a uma comunidade que possibilitam uma vasta rede de trabalhadores potenciais. Estamos a falar em *Crowdsourcing* apenas em situações em que uma empresa publica um problema *on-line*, recebe um grande número de soluções para este, de diversos indivíduos, as ideias vencedoras são recompensadas e a empresa produz a ideia para seu próprio benefício (Brabham, 2008).

Um sistema de *Crowdsourcing* é um sistema que abrange um grupo de pessoas que são utilizadas para resolver um problema definido pelos donos desse mesmo sistema (Doan et. al, 2011).

Segundo Braham, et al. 2008, sob certas circunstâncias, um grupo de pessoas demonstra uma inteligência extraordinária e é frequente que a inteligência do grupo seja maior do que a das pessoas mais inteligentes deste (Brabham, 2008).

### 3.3.2 Exemplos

Alguns exemplos bem conhecidos de *Crowdsourcing* são a Wikipedia, *Linux*, *Yahoo!*, *Answers*, *Youtube* (Doan et. al, 2011), *Flickr*, *ehow.com*, *Quora*, *Swivel*, *Demand Media*, *ESP* e *Mechanical Turk-based applications* (Doan et. al,2011). Para além destes, existem muitos outros exemplos menos conhecidos, que também utilizam este modelo tais como *Threadless*, *iStockphoto*, *InnoCentive* e *Goldcorp Challenge* (Brabham, 2008).

Threadless.com é uma empresa de *t-shirts* baseada na *Web* que realiza o processo de design através de *Crowdsourcing*, ou seja através da realização de uma competição *on-line*. Qualquer pessoa se pode juntar a esta comunidade, e qualquer elemento desta pode votar nos designs ou mesmo submeter as suas próprias ideias. Ou seja esta comunidade é composta por profissionais e amadores. Os designs têm que seguir um *template* e uma série de diretrizes para

poder realizar o seu design. Por fim os designs com melhores resultados são escolhidos pelos membros da empresa *Threadless* para serem impressos e postos à venda (Brabham, 2008).

*iStockphoto* é também uma empresa com negócio baseado na *Web*, que vende fotografias, animações e vídeos. Para poder pertencer a esta comunidade é necessário apresentar provas da sua identidade, para serem confirmadas por uma equipa da *iStockphoto*. Só assim será possível submeter as suas fotos, que são guardadas com palavras-chave. Os clientes podem depois comprar estas imagens para usar por exemplo em *websites*, em apresentações de negócios, entre outros. Os fotógrafos são recompensados com vinte por cento do preço de venda do item, por cada vez que se realiza *download* do mesmo. Contudo estes podem evoluir a sua relação com o *iStockphoto* e ganhar contratos que podem chegar a quarenta por cento do valor de venda por *download*. Assim como na *Threadless*, esta comunidade é composta por profissionais mas também amadores nesta área (Brabham, 2008). Segundo Howe (2006), nesse ano existiam já vinte e dois mil contribuidores neste *website* e os preços por imagem variavam entre \$1 e \$5 por cada imagem básica (Howe, 2006).

*InnoCentive* é uma organização que permite aos cientistas participantes obter reconhecimento profissional e financeiro através da resolução de problemas e a empresas tais como a *Boeing*, *DuPont* e *Proctor and Gamble*, receber soluções inovadoras de uma comunidade científica global. As categorias dos desafios vão passar pelas ciências da vida, química e ciências aplicadas. As soluções são submetidas pela comunidade, revistas pela empresa que propôs o desafio, que se mantém no anonimato durante o mesmo. Por fim se uma solução possuir os requisitos técnicos desejados para o desafio, a empresa que lançou o desafio, recompensa o indivíduo que submeteu a proposta (Brabham, 2008).

As aplicações de *Crowdsourcing* são desenvolvidas em plataformas públicas ou privadas. Nos últimos dez anos foram desenvolvidas diversas plataformas por várias empresas, tais como *Mechanical Turk*, *Turkit*, *Mob4hire*, *uTest*, *Freelancer*, *eLance*, *oDesk*, *Guru*, *Topcoder*, *Trada*, *99design*, *Innocentive*, *CloudCrowd* e *CloudFlower* (Doan et. al, 2011). Estas plataformas servem para mais fácil e rapidamente se poder desenvolver um sistema de *Crowdsourcing* em diversos domínios (Doan et. al, 2011).

Na Tabela 5 são dados alguns exemplos de sistemas que utilizam *Crowdsourcing* segundo Doan et al. (2011).

Tabela 5 - Exemplos de Sistemas de *Crowdsourcing* (Doan et. al, 2011)

Natureza	Necessita recrutar utilizadores?	Função dos utilizadores	Exemplos	Problemas alvo
Explícita	SIM	Avaliar Avaliar, Votar, Colocar <i>Tags</i>	Amazon Del.ici.ous.com Google Co-op	Avaliar uma coleção de itens
		Partilhar Itens, Conhecimento Textual, Conhecimento estruturado	Napster, Youtube, Flickr, CPAN, programmableweb.com Mailing Lists, Yahoo! Answers, QUIQ, ehow.com, Quora Swivel, Many Eyes, Google Fusion Tables, Google Base, bmr.wisc.edu, galaxyzoo, Plazza, Orchestra	Construir uma compilação de itens que pode ser partilhada entre utilizadores
		Networking	• LinkedIn, MySpace, facebook	Construir redes sociais
		Construir Artefactos Software Bases de Conhecimento textual Bases de Conhecimento Estruturado Sistemas Outros	Linux, Apache, Hadoop Wikipedia, openmind, Intellipedia, ecolicomunity Wikipedia Infoboxes/DBpedia, IWP, Google Fusion Tables, YAGO-NAGA, Cimple/DBLife WikiaSearch, mahalo, Freebase, Eureka Digg.com, SecondLife	Construir artefactos físicos
		Execução de Tarefas	Encontrar extraterrestres, eleições, encontrar pessoas, criação de conteúdo	Possivelmente qualquer problema
		Implícita	SIM	Jogar jogos com um objetivo Apostar em mercados de previsão Usar contas privadas Resolver <i>CAPTCHAS</i>

		Comprar, vender, leiloar, jogar jogos <i>multiplayer</i>		
	NÃO	Procura de palavras-chave Compra de produtos Procurar <i>Websites</i>	Google, Microsoft, Yahoo Amazon Yahoo	Correção ortográfica, previsão epidémica Recomendação de produtos Reorganizar o <i>Website</i> para melhorar o acesso

A Tabela 5 apresenta uma divisão em termos de natureza e em termos de necessidade de recrutar utilizadores. Ou seja, se os utilizadores colaboram explicitamente ou implicitamente e se é mesmo necessário recrutar utilizadores ou não. De uma forma simplificada são expostas as funções que os utilizadores realizam ao participar em *Crowdsourcing*, assim como uma descrição resumida do problema que estão a tentar resolver e são dados alguns exemplos.

### 3.3.3 Crowdsourcing e Opensource

O *Opensource* envolve a permissão de acesso aos elementos essenciais de um produto a qualquer pessoa, para um fim de melhoria colaborativa de um produto já existente. Envolve também uma transparência e distribuição grátis do produto ao longo das várias etapas do seu desenvolvimento aberto. Ou seja permite a disponibilização de um produto para que as pessoas possam contribuir para melhorá-lo. Alguns exemplos de sucesso deste modelo são o sistema operativo Linux e o *browser Mozilla Firefox*. Contudo este modelo é mais apropriado para desenvolvimento de *software* do que para outras aplicações, pois não se preocupa com interesses pessoais e exigências materiais de produção. É nesta perspetiva que o *Crowdsourcing* supera as limitações de *opensource*, pois fornece uma forma clara de compensação dos seus contribuidores e um modelo híbrido que conjuga elementos transparentes e democratizados de *opensource* num modelo de negócio vantajoso, tudo proporcionado através da *Web* (Brabham, 2008).

A produção em *opensource* defende o contrário desta noção de *Crowdsourcing* pois torna disponível código para toda a gente. Existem produtos *opensource* de qualidade superior à de outros que são realizados pela forma habitual. Apesar de nobre, esta visão é um pouco ingénua

pois qualquer produto material não se produz por si só e não é livre de custos e riscos (Brabham, 2008).

### 3.3.4 Desafios Principais

Os sistemas de *Crowdsourcing* apresentam quatro desafios chave: como recrutar contribuidores, que contribuições conseguem fazer, avaliar os utilizadores e as suas contribuições (Doan et. al, 2011).

Recrutar e manter contribuidores é um dos desafios mais importantes do *Crowdsourcing*. O recrutamento pode ser realizado por uma entidade autoridade para realizar esta tarefa, como por exemplo um gestor. Para além disso os utilizadores são recompensados monetariamente, o que torna estas propostas mais atrativas. Uma terceira solução é a requisição de voluntários. Para além destas existe outra solução que passa pelo pagamento dos próprios utilizadores pelo serviço. Por outras palavras, um utilizador paga pela utilização de um Sistema A para contribuir para outros Sistema B (Doan et. al, 2011).

Depois de escolhida a estratégia de recrutamento é necessário pensar numa forma de encorajar e manter os utilizadores. Isto pode ser feito a partir de uma forma de gratificação instantânea, de uma experiência agradável, do estabelecimento de renome, confiança ou reputação, ou mesmo através do estabelecimento de competições (Doan et. al, 2011).

Outro desafio está relacionado com as contribuições que os utilizadores podem fazer. A contribuição dos utilizadores, em alguns sistemas, é bastante limitada, ou seja estes podem avaliar, atribuir *tags*, partilhar, entre outros. Contudo em sistemas de *Crowdsourcing* mais complexos, as contribuições podem ser muito variadas. O desafio é criar uma extensão de possíveis contribuições. Para isso é necessário ter em conta o quão exigente cognitivamente os contribuidores são, pois quanto mais cognitivos, mais relutante será a sua participação. Ou seja utilizadores que estão mais acima no *ranking* têm mais motivação para participar nas contribuições mais complicadas. Para além disso é preciso ter em conta qual o impacto que uma contribuição deve ter. Este pode ser medido considerando o quanto essa contribuição pode afetar o sistema de *Crowdsourcing*. Deve ser também considerado o facto de alguns contribuidores poderem ser máquinas, e nesse caso as tarefas devem ser divididas para que os

utilizadores humanos recebam as tarefas mais fáceis para estes, mas difíceis para as máquinas e vice-versa (Doan et. al, 2011).

Outro desafio é a combinação das contribuições dos utilizadores. Apesar de ainda haver muitos sistemas de *Crowdsourcing* que ainda não combinam as contribuições dos utilizadores, esta tarefa pode ser vantajosa. Um caso bastante conhecido de combinação de contribuições é a *Wikipédia*, pois permite aos utilizadores fundir as suas edições, mas há também outros casos em que esta combinação é feita automaticamente. Independentemente da forma como a combinação das contribuições é realizada, o grande desafio é decidir o que fazer quando dois utilizadores (ou um número par de indivíduos) diferem numa determinada opinião. Normalmente o que uma solução automática faz é medir o peso das contribuições a partir da pontuação atribuída ao utilizador. Numa solução manual, os utilizadores conversam sobre o assunto entre si. As soluções automáticas são mais eficientes, contudo não funcionam para todos os tipos de contribuições. Em situações que envolvam contribuições mais conflituosas, que necessitem de mais discussão, as soluções manuais são mais adequadas. As discussões não são possíveis quando envolvem contribuidores que são máquinas, pois estes não poderiam explicar a um outro utilizador o seu ponto de vista ou a sua explicação (Doan et. al, 2011).

Por fim o último desafio é a avaliação das contribuições realizadas pelos utilizadores. Os sistemas de *Crowdsourcing* devem ter uma forma de conseguir gerir utilizadores maliciosos, através de técnicas de bloqueio, deteção e detenção. Primeiramente é necessário bloquear utilizadores maliciosos através da limitação de utilizadores em determinadas contribuições. Em segundo lugar estes devem ser detetados através de diversas técnicas automáticas e manuais. E em terceiro lugar os utilizadores maliciosos devem ser punidos. A forma mais comum de punição destes é bani-los do sistema, ou então a “humilhação pública” de serem expostos aos outros utilizadores como sendo utilizadores maliciosos (Doan et. al, 2011).

### 3.3.5 Vantagens e Desvantagens

Apesar de juntar os benefícios da filosofia de *opensource* com os benefícios gerais de qualquer negócio, o *crowdsourcing* pode afetar negativamente a força de trabalho. Ou seja as soluções vencedoras valem muito mais do que o valor que os contribuidores recebem. Por exemplo na *Threadless* os designers ganham muito menos do que um profissional ganharia se



recorressem ao seu trabalho através de outsourcing. Um outro exemplo é o *iStockphoto* em que os fotógrafos amadores não se importam de colocar os seus trabalhos por preços mínimos, o que torna o trabalho dos fotógrafos profissionais obsoleto (Brabham, 2008).

Para além disso este tipo de trabalho não está acessível a todos pois ainda existem muitas pessoas que ainda não possuem ligação à *Internet*, ou não possuem uma ligação de alta velocidade limitando-as da participação em certos projetos. Ou seja isto significa que não se pode garantir uma diversidade de opiniões na *crowd*. Segundo o autor, a diversidade deve ser dividida em partes mais pequenas como diversidade de identidades, de capacidades e de investimento político. A divisão por diversidade de identidades é importante pois cada indivíduo tem características específicas e próprias. Estas diferenças podem assim proporcionar visões diferentes relativamente a determinado problema e assim obter soluções superiores (Brabham, 2008).

Uma vantagem do *Crowdsourcing* é a capacidade de oferecer aos seus colaboradores uma possibilidade de empreendedorismo. São satisfeitas as necessidades de adquirir novas capacidades, e de resolver problemas. Esta vontade de participar nestes projetos parte do desejo de adquirir novas capacidades para uma procura de melhor emprego ou mesmo como forma de trabalho empreendedor como independente. Contudo esta possibilidade de sobressair alguns indivíduos que fornecem as melhores soluções e torná-los profissionais mais capazes, não é a grande vantagem do *Crowdsourcing*, mas sim as mentes jovens capazes de produzir ideias inovadoras e o produto resultante destas ideias adotado pelas organizações. Estas são as mais beneficiadas por este sistema (Brabham, 2008).

### 3.3.6 Trabalho Futuro

Existe ainda muito trabalho a ser feito nesta área, assim como investigação para entender como é que a *crowd* se sente em relação ao seu papel como trabalhador de outras empresas, como é que a *crowd* resiste às formas de manipulação exercidas pelas empresas. Para além disso é também necessário ter em atenção o tipo de indivíduos que ainda não participam nestes projetos, as novas barreiras que se apresentam à sua participação. Investigar que tipos de projetos é que são bem-sucedidos e quais os que falham, mais aspetos da produção, análise legal e ética destes, assim como *standards* das melhores práticas do *crowdsourcing* (Brabham, 2008).

## 3.4 CAPTCHA

O constante crescimento dos serviços baseados na *Web* revolucionaram a forma como as pessoas comunicam e partilham informação. Esta situação requer um reforço nas medidas de segurança para evitar que programas maliciosos automatizados diminuam a qualidade dos serviços fornecidos (Sharma et. al, 2013). Devido ao aparecimento deste *software* e *scripts* automatizados, que corrompem as aplicações *Web*, é necessário diferenciar o utilizador humano comum de um computador (Payal et. al, 2012). Estes podem afetar diversos serviços como o *email*, blogs, votações, entre outros. *CAPTCHA* significa *Completely Automated Public Turing test to tell Computers and Humans Apart* e surge neste contexto, para fazer esta diferenciação entre o utilizador humano e os processos automáticos (Sutherland, 2012). Este surge por volta do ano 1997 em AltaVista, como forma de prevenir a submissão automática de URLs, que provocava uma distorção no seu mecanismo de *ranking* de URLs (Abrich et. al,2011).

### 3.4.1 Definição

*CAPTCHAS* são problemas de resolução simples para os humanos, que os programas computacionais não conseguem resolver. Assim, este torna possível a distinção entre humanos e programas automáticos (Sutherland, 2012). Pode ser considerado como um teste para deteção de programas que tentam aceder ao serviço, sem autorização ou de forma abusiva. Qualquer utilizador que consiga fornecer a solução correta é assumido como sendo humano e é garantido acesso ao serviço, caso contrário, é barrado (Sharma et. al, 2013). Este sistema é bastante útil na prevenção de ataques como *denial-of-service*, de força bruta, comentários não desejáveis em *blogs* e anúncios *spam* realizados por outros computadores (Saxena, 2013). Um programa malicioso pode também criar milhares de contas de *emails* grátis por minuto e assim enviar milhões de *emails* com “lixo”. Este mecanismo foi adotado por vários *websites* comerciais como a *Yahoo!*, a *Google* e *Baidu* (GAO et. al, 2014).

### 3.4.2 Funcionamento

O *CAPTCHA* é utilizado como um simples *puzzle*, como por exemplo a identificação de determinadas letras a partir de uma imagem distorcida. Este processo tem uma resolução complicada por parte dos programas automáticos, o que lhes impede o acesso. A maioria dos *CAPTCHAs* possui códigos aleatórios sob a forma de imagens, letras e números que se podem sobrepor. Esta sobreposição aumenta ainda mais a dificuldade para o acesso de programas automáticos (Pawar & Bauskar, 2013). Já existem *softwares* que conseguem passar a barreira dos *CAPTCHAs* baseados em texto, contudo os *CAPTCHAs* baseados em imagens foram introduzidos de forma a aumentar a dificuldade desta barreira. Estes testes, apesar de serem criados facilmente por máquinas automatizadas, são de elevada dificuldade de resolução para estas (Payal et. al, 2012). Um bom *CAPTCHA* deve ser de fácil resolução para humanos e de impossível resolução para uma máquina. Este teste deve ser de fácil utilização e robusto o suficiente para resistir a ataques de programas (Saxena, 2013). Para conseguir passar o teste com sucesso é necessário possuir três capacidades: reconhecimento invariável, segmentação e análise. Reconhecimento invariável, para conseguir reconhecer os formatos dos caracteres, segmentação e análise para separá-los e posteriormente traduzi-los nos símbolos corretos (Shi et. al, 2013). Diversos *websites* adotaram o princípio de *Connecting characters together* (CCT) para resistir ao método de segmentação. Desta forma os caracteres são conectados uns aos outros para que não se consigam extrair da imagem individualmente (GAO et. al, 2014).

### 3.4.3 Aplicações

Os *CAPTCHAs* têm então uma função de segurança e podem ser encontrados em diversas aplicações com funções de:

**Prevenir comentários *spam* em *blogs*** – A ideia é evitar a submissão de comentários falsos com o propósito de aumentar posições nos mecanismos de *ranking* (Payal et. al, 2012). Isto é designado por comentários *spam* e o *CAPTCHA* permite impedir esses programas automáticos de partilhar os seus comentários premeditados. A utilização de uma conta de utilizador no serviço *Web* seria muito custoso, numa situação em que o utilizador apenas deseja publicar um comentário (Shanker et. al, 2013);

**Proteger os serviços de *email* gratuitos** - Existem diversas empresas, como a *Yahoo!*, *Google*, *Rediff* e *Microsoft*, que fornecem serviços de *email* gratuitos e a maior parte deles sofrem do mesmo tipo de ataque: as *bots*. Estas são capazes de criar milhares de contas *email* a cada minuto, o que pode ser evitado com a utilização de um *CAPTCHA* (Saxena, 2013).

**Proteger os endereços *email*** – Existem *spammers* que percorrem a *Web* à procura de endereços *email*. O *CAPTCHA* fornece um mecanismo para ocultar o endereço *email* destes. A ideia é pedir ao utilizador que responda a um *CAPTCHA* antes de mostrar o *email* (Shanker et. al, 2013)

**Proteger os registos em *websites*** – Todos os serviços gratuitos devem ser protegidos com *CAPTCHA* para evitar o abuso por parte de programas automatizados como as *bots* (Saxena, 2013).

**Proteger os jogos *on-line*** – Também os jogos *on-line* devem ser protegidos para evitar que robots joguem (Saxena, 2013).

**Proteger de *worms* e *spam*** – O *CAPTCHA* fornece também proteção contra *spam* e *worms* em geral (Shanker et. al, 2013).

**Evitar ataques de dicionário** – A ideia é evitar que um computador faça várias iterações para encontrar a *password* certa. É também recomendável a utilização de *CAPTCHAs* em sistemas baseados em *passwords* para evitar que outros computadores consigam entrar em contas que não as suas (Saxena, 2013). Assim ao fim de um certo número de tentativas falhadas é apresentado o *CAPTCHA* e caso seja um programa automático o acesso será barrado (Shanker et. al, 2013).

***Bots* de mecanismos de procura** – Os donos de alguns *websites* não querem os seus sites indexados em mecanismos de procura. Não desejam isso para evitar que outros encontrem estes *websites* facilmente. São então utilizadas umas *tags* HTML para evitar que as *bots* leiam as páginas (Saxena, 2013). Contudo isto não garante que as páginas não sejam lidas pois esta *tag* apenas serve para dizer “sem *bots*, por favor”. Estas são geralmente respeitadas, contudo a única forma de garantir que não serão usados em mecanismos de procura é a através da utilização de *CAPTCHAs* (Shanker et. al, 2013).

**Votações *on-line*** – Para além da verificação do endereço IP, pode ser também utilizado o *CAPTCHA* para se verificar se um utilizador já realizou o seu voto (Saxena, 2013).

**Prevenir acesso não autorizado** – O *CAPTCHA* pode prevenir que um hacker tente descobrir uma *password* utilizando métodos de força bruta ou outros métodos de descobrir *passwords* (Shanker et. al, 2013).

### 3.4.5 Caraterísticas

Algumas caraterísticas desejáveis para o sucesso do *CAPTCHA* (Saxena, 2013) (Pawar & Bauskar, 2013):

- O teste deve ser gerado automaticamente por uma máquina, contudo esta não pode ser capaz de o resolver;
- Deve ser fácil e de rápida resolução para o humano (tempo de resposta deve ser até trinta segundos);
- O teste deve passar, com alta confiança, todos os utilizadores humanos que tentem realiza-lo. Ou rejeitar um número mínimo, pois a função do *CAPTCHA* é apenas impedir o acesso de máquinas e programas automáticos;
- Deve ser desafiante para máquinas, possuindo problemas complicados e difíceis de resolver para estas, e rejeitando assim o acesso a grande parte dos utilizadores máquina;
- A base de dados da máquina que gera os *CAPTCHAs* deve ser suficientemente para impedir ataques;
- O teste deve ser universal, ou seja, deve ser independente da língua do utilizador, do conhecimento em termos de educação e área geográfica.

### 3.4.6 Problemas e Vulnerabilidades

Os *CAPTCHAs* podem apresentar alguns problemas para pessoas com problemas de visão, ou com deficiência (Pawar & Bauskar, 2013). Já foram desenvolvidos diversos tipos de *CAPTCHAs* para contornar este problema, como os *CAPTCHAs* audíveis. Contudo estes podem não ser tão fáceis de decifrar pois devem possuir também algum tipo de ruído para impedir a

compreensão deste por parte dos programas automáticos. Isto pode levar a que o *CAPTCHA* se torne também difícil de decifrar pelo humano (Schlaikjer, 2007).

Um *CAPTCHA* tem uma média de tempo de resposta por volta dos dez segundos, o que realizado diariamente e várias vezes ao dia, pode ser custoso em termos de tempo. Para contornar este problema, Luis von Ahn criou um projeto designado *reCAPTCHA*. O seu objetivo é utilizar os *CAPTCHAs* já resolvidos para que este esforço não seja desperdiçado (Sutherland, 2012).

### 3.4.7 Tipos de *CAPTCHAs*

O *CAPTCHA* pode ser dividido em dois, tendo em conta o método que utiliza: *CAPTCHA* baseado em *Optical Character Recognition (OCR)* e *CAPTCHA* não baseado em *OCR*. Os *CAPTCHAs* que utilizam o método *OCR* mostram imagens com palavras distorcidas e com efeitos sobre estas. Estes efeitos aumentam a dificuldade no reconhecimento das palavras. Os métodos que não utilizam *OCR*, são mais simples e mais fáceis de decifrar e podem incluir *CAPTCHAs* lógicos, com imagens, áudios e vídeos (Shanker et. al, 2013).

A maior parte das técnicas existentes são derivadas da ideia original de geração de texto sob a forma de imagens (Abrich et. al, 2011). Estas técnicas podem ser:

***CAPTCHA* baseado em Texto** – O utilizador deve introduzir o texto que aparece distorcido (Abrich et. al, 2011). Este é o *CAPTCHA* mais comuns e foi inicialmente introduzido por *AltaVista*. Utilizava a distorção de texto o que provocava uma redução na precisão do *OCR*. Contudo este possuía uma desvantagem, pois era vulnerável à segmentação, que permite o reconhecimento e isolamento de cada caracter. Atualmente existem abordagens que utilizam os textos sob a forma de imagens (Pawar & Bauskar, 2013). A abordagem mais simples desta técnica é a utilização de questões que apenas um utilizador humano pode responder, como por exemplo “Qual a última letra da palavra PORTUGAL?”, “Quanto é dez menos dois?” ou “Se amanhã é domingo, que dia é hoje?” (Sharma et. al, 2013). Contudo este possui algumas desvantagens pois estas perguntas encontram-se guardadas numa base de dados juntamente com as suas soluções e até já existem *websites* que respondem a questões, resolvem problemas matemáticos, fornecem factos instantaneamente, fazem cálculos, conversações, e fornecem

dados e estatísticas quantitativas em tempo-real, o que fornece uma forma fácil de obter respostas ao mecanismo *CAPTCHA* (Vaishakh & Harish, 2011). Um tipo de *CAPTCHA* baseado em texto interessante é o *reCAPTCHA* pois melhora o processo de digitalização de livros ao apresentar palavras difíceis de decifrar, para que os utilizadores as decifrem (Saxena, 2013). Neste método, os utilizadores deparam-se com duas palavras, uma cujos caracteres são conhecidos e outra desconhecida. Caso o utilizador acerte na palavra conhecida, a palavra desconhecida é guardada para ajudar a identificar os caracteres da imagem (Pawar & Bauskar, 2013). Este serviço é gratuito e é possível pois os programas *OCR* alertam quando uma palavra não é reconhecida corretamente. Isto aumenta a viabilidade deste método pois as máquinas já falharam o teste de reconhecimento da palavra pelo menos uma vez (Sharma et. al, 2013). Na figura 9 podem verificar-se alguns exemplos de *CAPTCHAs* baseados em texto.

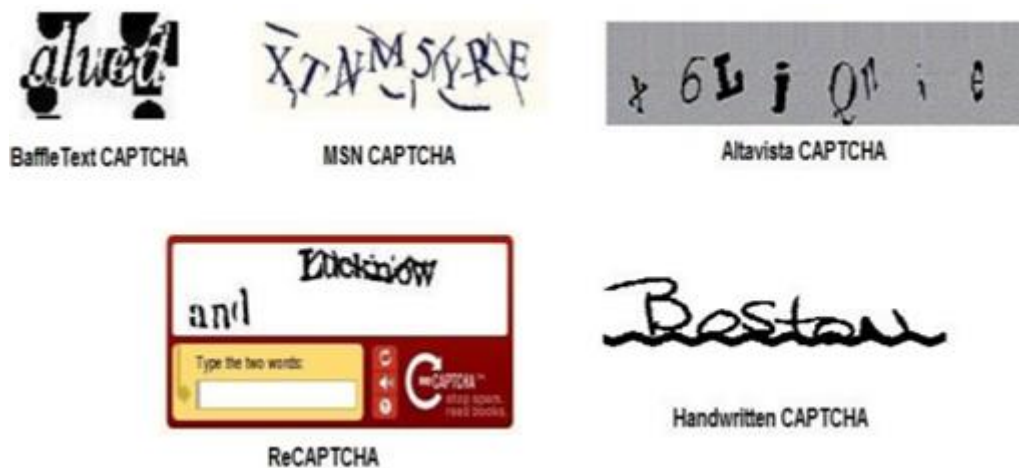


Figura 9 – Alguns exemplos de *CAPTCHAs* baseados em Texto (Pawar & Bauskar, 2013)

***CAPTCHA* Gráfico** – *CAPTCHAs* gráficos são desafios que envolvem imagens ou objetos que têm semelhanças e que os utilizadores devem adivinhar. O programa gera uma espécie de *puzzle* que posteriormente não será capaz de resolver (Sharma et. al, 2013). Existem diversos exemplos de aplicações com este tipo de *CAPTCHA*, por exemplo o *CAPTCHA ESP Pix*, em que são fornecidas quatro imagens e o utilizador tem que escolher de uma lista de setenta e duas opções, a palavra que se relaciona com as imagens. Outro exemplo desta abordagem é o *Animal species image recognition for restricting access (Asirra)*, em que o utilizador deve selecionar de todas as imagens que lhe são apresentadas, aquelas que possuem ou gatos ou cães, dependendo do que lhes é indicado (Saxena, 2013). Na figura 10 pode ver-se um exemplo de

um *Asirra CAPTCHA* e um *ESP CAPTCHA*. Esta técnica tem mais probabilidades de falhar pois é apresentado, ao utilizador, um número finito de respostas possíveis (Pawar & Bauskar, 2013).



Figura 10 - Exemplos de *CAPTCHAs* gráficos (Asirra e ESP CAPTCHA) (Pawar & Bauskar, 2013)

**CAPTCHA de Áudio** – O programa seleciona uma palavra ou uma sequência de números, constrói um clip de som com estes caracteres e distorce-o adicionando barulho, para aumentar a dificuldade para o decifrar. Depois, este fornece o clip de som ao utilizador e pede-lhe que adicione os caracteres corretos (Sharma et. al, 2013). Na Figura 11 podem ver-se exemplos de *CAPTCHA* áudio.



Figura 11 - Exemplo de CAPTCHA áudio (Pawar & Bauskar, 2013)



### 3.4.8 *reCAPTCHA*

A ideia chave por detrás do projeto *reCAPTCHA* é permitir realizar determinada tarefa, ao utilizador, apenas após resolver um determinado *puzzle*, denominado por *CAPTCHA* que permite provar que o utilizador é humano (Doan et. al, 2011). Este novo tipo de *CAPTCHA*, é uma mais-valia no sentido em que, realiza tarefas cujas respostas são guardadas para reutilização futura. Neste caso o *reCAPTCHA* ajuda a definir palavras de livros ou jornais digitalizados, que não os programas não conseguem extrair corretamente (Schlaikjer, 2007). O que torna o *reCAPTCHA* único é a sua divisão em dois problemas separados. Ou seja, possui um problema similar ao dos *CAPTCHAS* tradicionais, em que é necessário identificar uma sequência de caracteres alterados para resistir a possíveis ataques, cuja resposta é conhecida. Para além deste possui também uma sequência de caracteres que foram digitalizados de um livro ou jornal e o *software OCR* falhou em identificar (Sutherland, 2012). O *software OCR* não consegue atingir a perfeição e ser preciso em todas as transcrições que realiza, especialmente em livros velhos cujas palavras podem estar um pouco apagadas ou distorcidas. Nestes casos o *software OCR* fornece a sua resposta mais aproximada mas com um valor de confiança baixo (McMillen & Veloso, 2008). As palavras com valor de confiança baixo ou que falharam a sua identificação por pelo menos dois programas *OCR*, são enviadas para o *reCAPTCHA*, e também caso estas não façam parte do dicionário. A primeira palavra é considerada como sendo o controlo. Caso o utilizador acerte na palavra de controlo, a segunda palavra é assumida como correta e guardada pelo sistema. Assim que uma dada palavra é resolvida um número de vezes suficiente para criar um intervalo de confiança razoável, o *reCAPTCHA* devolve o valor ao processo de digitalização. Para além da sua dupla função o *reCAPTCHA* é um dos *CAPTCHAs* mais seguros e mais fáceis de utilizar. Apesar do utilizador perder tempo a responder aos *reCAPTCHAs*, este pode não ser completamente desperdiçado, pois está a contribuir para a resolução de palavra difíceis de identificar para um programa computacional (Sutherland, 2012). O *reCAPTCHA* é um serviço gratuito e é aplicado atualmente a diversos *websites*, tais como o *Facebook*, *Microsoft* e *Twitter* (GAO et. al, 2014). Na Figura 12 pode ver-se um exemplo de um *reCAPTCHA*.



Figura 12 - Exemplo de um *reCAPTCHA* (Vaishakh & Harish, 2011)

Geralmente o *reCAPTCHA* possui características como: seis a oito caracteres em cada desafio; não possui arcos de ruído nas imagens; não utiliza números, utiliza apenas onze caracteres em maiúsculas e vinte e cinco minúsculas; os caracteres podem estar conectados, mas não se sobrepõem (GAO et. al, 2014). Para além disso o *reCAPTCHA* adota diversas técnicas defensivas como: a utilização de caracteres contínuos; comprimento dos textos aleatórios; transformação linear; transformação em onda (Sano et. al, 2014).

## 4 Inteligência Coletiva para Análise de Sentimento

### 4.1 Introdução

Para a aplicação da Inteligência Coletiva para a Análise de Sentimento de mensagens do serviço *StockTwits*, desenvolveu-se uma aplicação que permitisse aos seus utilizadores classificar pequenas frases. Esta secção explica o desenvolvimento, aplicação e resultados obtidos dessa aplicação designada de "*Finance\$entiment*". Resumidamente a aplicação permite aos seus utilizadores classificar uma frase em "Positiva", "Negativa", "Neutra" ou "Difícil de Classificar". As respostas de cada utilizador são armazenadas e assim que uma frase possua três classificações, é realizada uma agregação das respostas obtidas. Esta agregação será efetuada via uma moda, ou seja, a classificação final atribuída a uma mensagem será obtida pelo valor mais comum (maioria) das três classificações. Tal resulta numa avaliação final de "Positiva", "Negativa", "Neutra" ou "Difícil de Classificar" e que é atribuída. Para além destas categorias, foi ainda definido o valor de "Indefinida", quando nenhum dos valores anteriores se encontre em maioria. Por uma questão de simplificação de texto, este método de agregação será designado de "média das avaliações" ao longo deste capítulo.

Os dados utilizados para estudo foram obtidos a partir da plataforma *StockTwits* disponível em <http://stocktwits.com>. Daqui foram extraídas as frases para classificação, assim como o sentimento atribuído pelo seu autor. Este sentimento representa-se como "bearish" ou "bullish", sendo o primeiro considerado negativo e o segundo positivo.

### 4.2 Planeamento do Projeto

Este projeto seguiu o planeamento apresentado nas Figuras 13 e 14. Resumidamente foram redefinidos os objetivos principais do projeto e realizada uma revisão de literatura sobre Análise de Sentimento e Inteligência Coletiva. Posteriormente foi desenvolvida a aplicação "*Finance\$entiment*" e testada através da sua disponibilização. De seguida foram analisados os dados obtidos a partir desta e comparados com os resultados obtidos através dos Léxicos e Algoritmo de *Text Mining*. Por fim, a aplicação e os resultados obtidos são documentados e a escrita do documento é finalizada e o mesmo é entregue.

	Task Name	Duration	Start	Finish	Predecess
0	<b>Desenvolvimento Da Dissertação</b>	<b>218 days</b>	<b>Wed 01/01/14</b>	<b>Fri 31/10/14</b>	
1	Redefinição dos Objetivos da Dissertação	6 days	Wed 01/01/14	Wed 08/01/14	
2	<b>Revisão de Literatura</b>	<b>46 days</b>	<b>Thu 09/01/14</b>	<b>Thu 13/03/14</b>	<b>1</b>
3	<b>Procura de Documentação Relevante</b>	<b>14 days</b>	<b>Thu 09/01/14</b>	<b>Tue 28/01/14</b>	
4	Análise de Sentimento	7 days	Thu 09/01/14	Fri 17/01/14	
5	Inteligência Coléitiva	7 days	Mon 20/01/14	Tue 28/01/14	4
6	Estruturação do Documento	1 day	Wed 29/01/14	Wed 29/01/14	5
7	<b>Escrita do Documento</b>	<b>30 days</b>	<b>Fri 31/01/14</b>	<b>Thu 13/03/14</b>	<b>6</b>
8	Análise de Sentimento	15 days	Fri 31/01/14	Thu 20/02/14	6
9	Inteligência Coléitiva	15 days	Fri 21/02/14	Thu 13/03/14	8
10	<b>Desenvolvimento da Aplicação</b>	<b>90 days</b>	<b>Fri 14/03/14</b>	<b>Thu 17/07/14</b>	<b>2</b>
11	Análise dos Requisitos	3 days	Fri 14/03/14	Tue 18/03/14	
12	Especificação das Deadlines	1 day	Wed 19/03/14	Wed 19/03/14	11
13	Design da Aplicação	4 days	Thu 20/03/14	Tue 25/03/14	12
14	Instalação das Ferramentas	2 days	Wed 26/03/14	Thu 27/03/14	13
15	Criação da Base de Dados	4 days	Fri 28/03/14	Wed 02/04/14	14
16	Desenvolvimento da Aplicação	74 days	Thu 03/04/14	Tue 15/07/14	15
17	Testes da Aplicação	2 days	Wed 16/07/14	Thu 17/07/14	16
18	Integração da Aplicação com o Facebook	15 days	Thu 03/04/14	Wed 23/04/14	15
19	Verificação da Aplicação	4 days	Thu 24/04/14	Tue 29/04/14	18
20	<b>Disponibilização da Aplicação</b>	<b>60 days</b>	<b>Fri 18/07/14</b>	<b>Thu 09/10/14</b>	<b>10</b>
21	Preparação dos Dados a ser Testados na Aplicação	2 days	Fri 18/07/14	Mon 21/07/14	
22	Disponibilização da Aplicação	57 days	Tue 22/07/14	Wed 08/10/14	21
23	Recolha dos dados	1 day?	Thu 09/10/14	Thu 09/10/14	22
24	<b>Análise dos Dados</b>	<b>7 days</b>	<b>Fri 10/10/14</b>	<b>Mon 20/10/14</b>	<b>20</b>
25	Estruturação dos dados recolhidos	1 day	Fri 10/10/14	Fri 10/10/14	
26	Análise dos Resultados Obtidos	3 days	Mon 13/10/14	Wed 15/10/14	25
27	Análise dos Resultados Obtidos Através dos Léxicos FIN e SWN	1 day	Thu 16/10/14	Thu 16/10/14	26
28	Análise dos Resultados Obtidos Através do Algoritmo de Text Mining	1 day	Fri 17/10/14	Fri 17/10/14	27
29	Comparação dos Resultados obtidos com os do Algoritmo e Léxicos	1 day	Mon 20/10/14	Mon 20/10/14	28
30	Documentação da Aplicação Desenvolvida	2 days	Tue 21/10/14	Wed 22/10/14	24
31	Documentação da Análise de Resultados	2 days	Thu 23/10/14	Fri 24/10/14	30
32	Finalização da Escrita do Documento	4 days	Mon 27/10/14	Thu 30/10/14	31
33	Entrega da Dissertação	1 day	Fri 31/10/14	Fri 31/10/14	32

Figura 13- Planeamento do Projeto



Figura 14 - Gráfico de Gantt do Planeamento do Projeto

### 4.3 Metodologia

Metodologia é a orientação de uma dada pesquisa, que segue um conjunto de normas e utiliza diversas técnicas com o objetivo de desenvolver o processo de verificação empírica (Pardal & Correia, 1995). Estes consideram técnica como sendo um instrumento de trabalho utilizado para a realização de uma pesquisa. Segundo Hegenberg (1976), método é o caminho a seguir para chegar a determinado resultado. O método seguido apresenta sete etapas (Quivy & Campenhoudt, 2005):

Etapa 1: Criação da pergunta – Uma investigação implica a procura de algo, que pode passar por hesitações, desvios ou incertezas. Contudo é necessário escolher um fio condutor claro para que o trabalho seja estruturado de forma coerente. A pergunta deve exprimir o mais exatamente possível, o que se pretende saber, elucidar e compreender melhor.

“Qual a assertividade humana (Inteligência Coletiva) na Análise de Sentimento de pequenas mensagens de âmbito financeiro? ”

Etapa 2: Exploração: Leituras – Após a definição da pergunta é necessário perceber como proceder para garantir a qualidade da informação obtida. Esta comporta operações como a leitura. Esta deve possuir uma ligação à pergunta, apresentar uma dimensão razoável, e alguma diversidade de abordagens. Devem ser consultados artigos de revistas e livros, cujas referências se podem encontrar em índices e páginas de procura como o *Google scholar*. Os resumos devem destacar as ideias principais de modo a fazer surgir o pensamento do autor.

Etapa 3: A Problemática – É a abordagem escolhida para resolver o problema formulado na pergunta inicial. Esta não depende do acaso ou inspiração pessoal, mas deve implicar debates e correntes de pensamento em evolução. Deve realizar-se um balanço das problemáticas possíveis e escolher a mais relevante.

Etapa 4: Construção do Modelo de Análise – Esta etapa constitui a transformação das perspetivas e ideias numa linguagem e forma de trabalho sistemático de análise e recolha de dados.

Etapa 5: Observação – Esta engloba as operações através das quais o modelo de análise é submetido a teste e confrontado com dados observáveis. Deve responder a perguntas como:

**O que se deve observar?** Se as pessoas conseguem acertar no sentimento expresso na frase.

**Em quem?** Em indivíduos que possuam conhecimentos de gestão ou financeiros.

**Como?** Através da utilização da aplicação desenvolvida.

Etapa 6: Análise da Informação- Esta etapa explora a verificação empírica. A sua função é revelar factos esperados ou não e interpretá-los. Para além disso esta deve propor aperfeiçoamentos do modelo de análise ou pistas de investigação para o futuro. Ao longo desta etapa são retirados os ensinamentos necessários para prosseguir com a criação de conclusões.

Para uma **análise quantitativa**, neste projeto são comparadas as classificações de frases obtidas com o sentimento apresentado nestas. Para facilitar este trabalho são utilizadas pequenas frases com um sentimento associado pelo seu autor, o que permite uma comparação

de classificações. Desta análise resulta uma ou mais percentagens relativas ao asserto das classificações.

Para além disso, para realizar uma **análise qualitativa**, são comparados estes mesmos resultados com os que podem ser obtidos através da utilização de léxicos e um algoritmo. Através da utilização dos mesmos dados de teste. Através desta análise pode verificar-se o desempenho dos utilizadores comparando a métodos automatizados.

Etapa 7: Conclusões- As conclusões devem apresentar uma retrospectiva das grandes linhas de procedimento, uma descrição dos contributos para o conhecimento originados pelo trabalho, considerações finais e trabalho futuro.

## 4.4 Aplicação Desenvolvida

### 4.4.1 Arquitetura do Sistema

O sistema utiliza uma arquitetura cliente-servidor, como pode ser visto na Figura 15. Quando um utilizador realiza um pedido através do cliente (*Browser*), este pedido vai ser enviado ao servidor da aplicação, sendo aí processado. O servidor em causa é um servidor disponibilizado pela Universidade do Minho (saxofone.dsi.uminho.pt). Este processa os pedidos dos clientes, realiza acessos ao servidor da Base de Dados quando necessário, como por exemplo na classificação de uma frase por parte de um utilizador. O servidor da aplicação retorna os devidos ficheiros e valores pedidos pelo cliente.

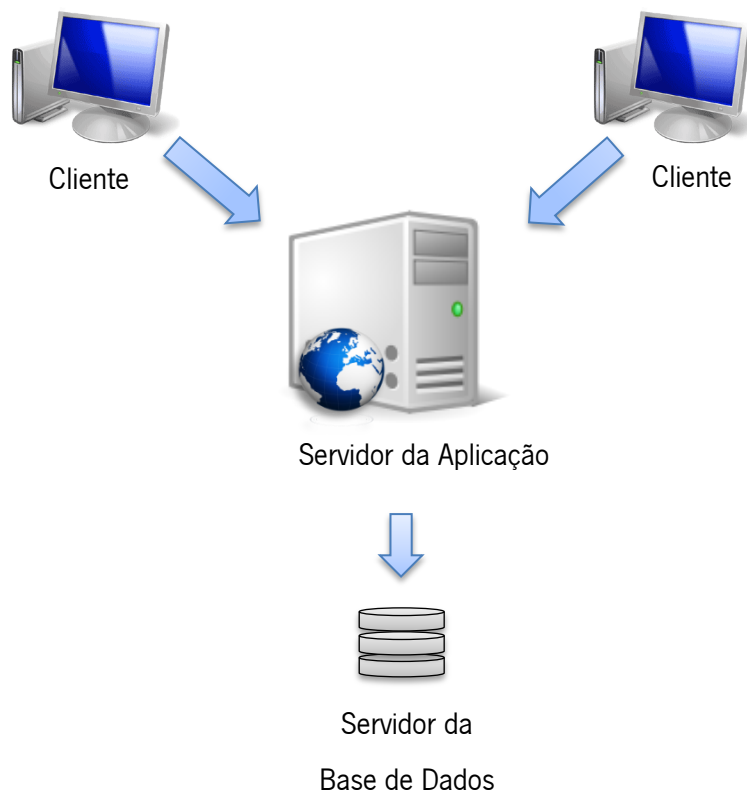


Figura 15 - Arquitetura da aplicação *Finance\$entiment*

#### 4.4.2 Ferramentas Utilizadas

Relativamente à aplicação foram utilizadas diversas ferramentas de desenvolvimento, tais como:

- *WinSCP* - utilizado para comunicar com o servidor;
- *Notepad++* - para criação de desenvolvimento de código;
- *MySQL Workbench 5.2* - para a criação da Base de Dados;
- Web Browsers (Safari, Google Chrome, Opera, Mozilla Firefox) - para teste da aplicação;
- *Facebook Developers* - para criação de uma página da aplicação no facebook.
- *Photoshop CS6* - para desenvolvimento de todo o design da aplicação;
- *mysql-admin* - para a gestão da Base de Dados;



Para desenvolvimento da documentação foram utilizadas as seguintes ferramentas:

- *Microsoft Office Visio Professional 2007* - para a descrição da aplicação através da linguagem UML;
- *Foxit Reader* - utilizado para ler documentos com extensão .pdf;
- *Microsoft Word 2010* - para desenvolvimento da documentação;
- *Microsoft Excel 2010* - utilizado na criação de tabelas relativas à análise de resultados.

## 4.5 Implementação

### 4.5.1 Descrição da Aplicação

#### 4.5.1.1 Visão do Utilizador

A aplicação é criada para obter mais facilmente as respostas dos utilizadores. O principal objetivo da aplicação é obter o máximo de respostas por parte dos utilizadores, o que implica um design simples e apelativo.

A aplicação foi partilhada entre docentes e alunos e foram utilizadas ferramentas como *LinkedIn* para a sua divulgação. A aplicação foi também incluída no *facebook* para uma maior difusão da mesma. Esta pode ser acedida através dos seguintes *links*:  
<https://financesentimentapp.dsi.uminho.pt> e  
<https://apps.facebook.com/sentencesentimentapp>.

A apresentação da aplicação pode ser vista na Figura 16. Nesta página inicial existe a possibilidade de autenticação por parte do utilizador, assim como a partilha da aplicação na rede social do *Facebook*. Para além disso existem nesta mesma página acessos importantes à política de privacidade da aplicação assim como informação relacionada com o projeto.

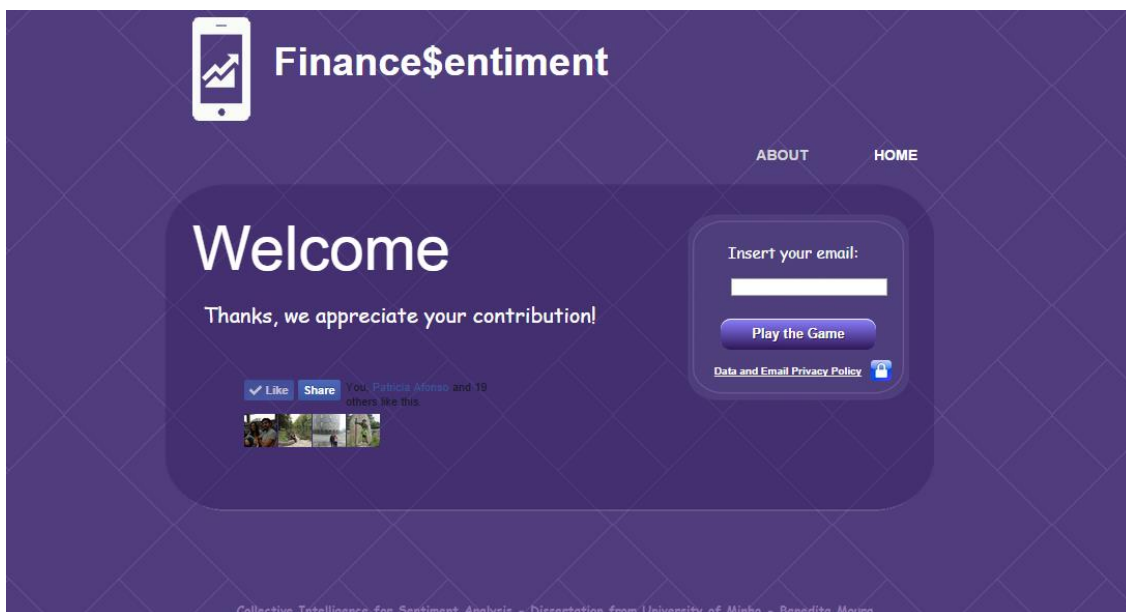


Figura 16 - Página inicial da aplicação *Finance\$entiment*

Um utilizador apenas necessita de registar o seu *email* para poder participar na aplicação. Este *email* pode mesmo ser fictício, contudo necessita de ser um *email* válido pois a aplicação realiza uma verificação de *email*. Caso o utilizador introduza um *email* inválido o sistema não permite o acesso à aplicação e mostra um aviso como se pode verificar na Figura 17. Caso o utilizador já tenha registado o seu *email*, e tente aceder novamente à aplicação, o *email* não será repetido mas entrará diretamente na sua conta. Este sistema de autenticação foi pensado de forma a simplificar o esforço do utilizador.

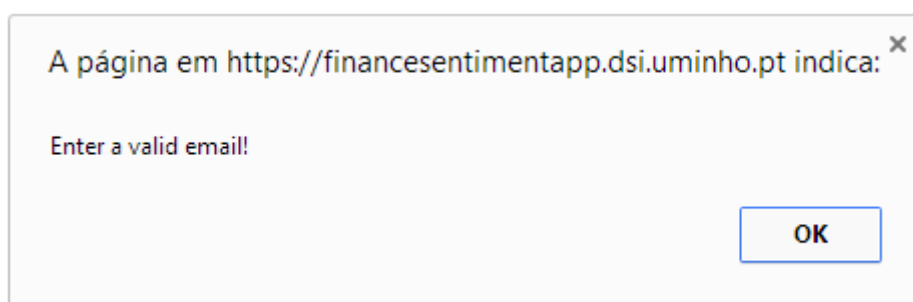


Figura 17 - Aviso de *email* inválido

Na Figura 18 pode ver-se a página relativa à política de privacidade da aplicação. Esta página tem como intuito informar os utilizadores do objetivo da sua contribuição e o fim dos dados fornecidos por estes, nomeadamente o seu *email* e as suas classificações.

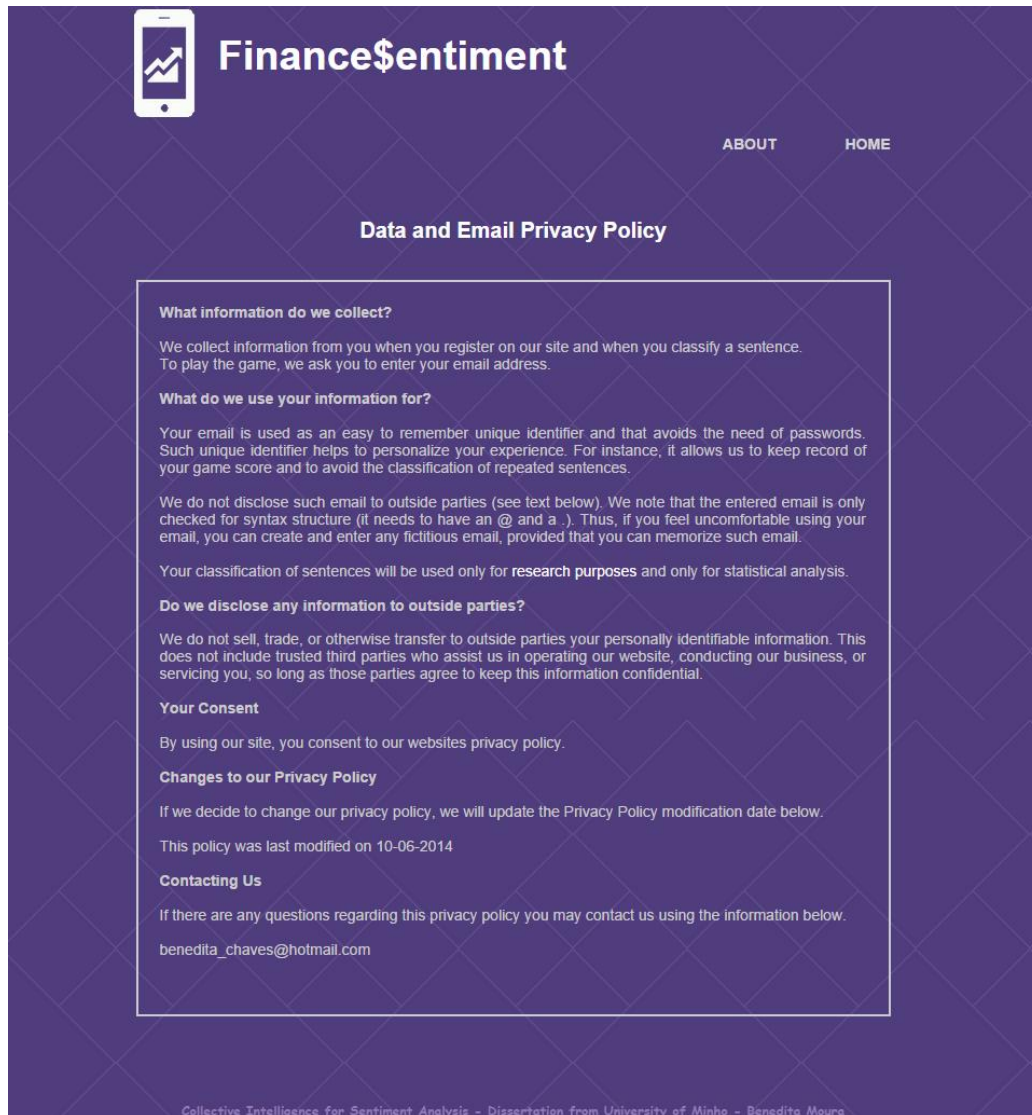


Figura 18 - Política de Privacidade

A página apresentada na Figura 19 tem como objetivo inserir o utilizador no contexto do projeto. Apresenta algumas ligações para documentos externos à aplicação que fornecem mais informação interessante.

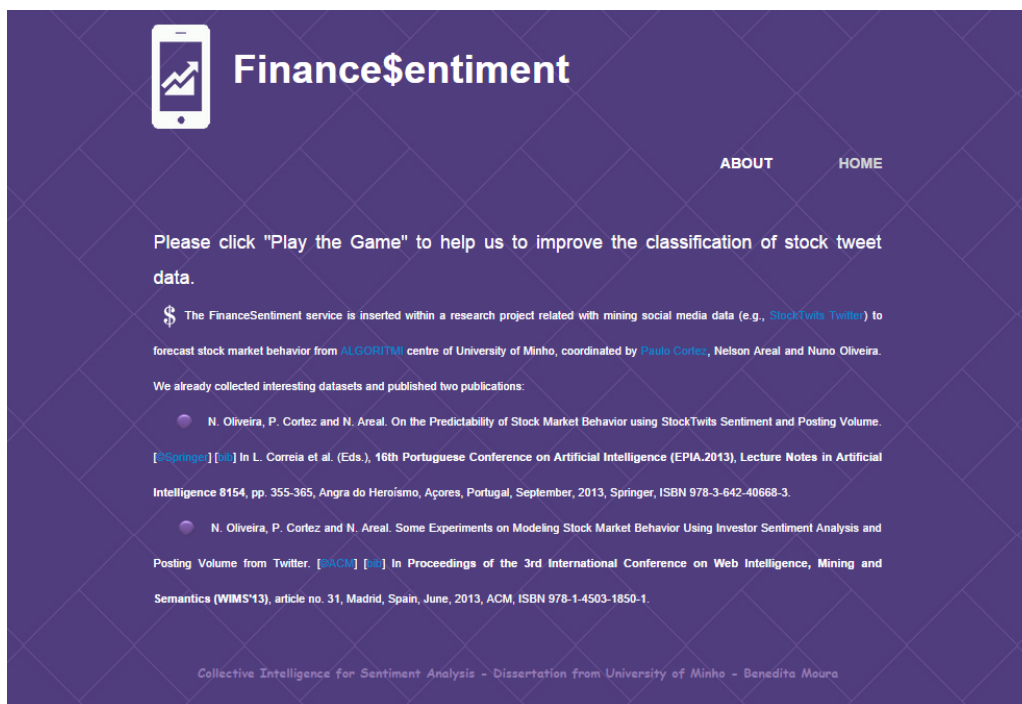


Figura 19 - Informação relativa ao âmbito do projeto

No primeiro contacto com a aplicação o utilizador é confrontado com três exemplos pormenorizadamente explicados do funcionamento da aplicação. Estes podem ser posteriormente visualizados *clikando* na “ajuda”. Estes exemplos podem ser vistos nas Figuras 20, 21 e 22.

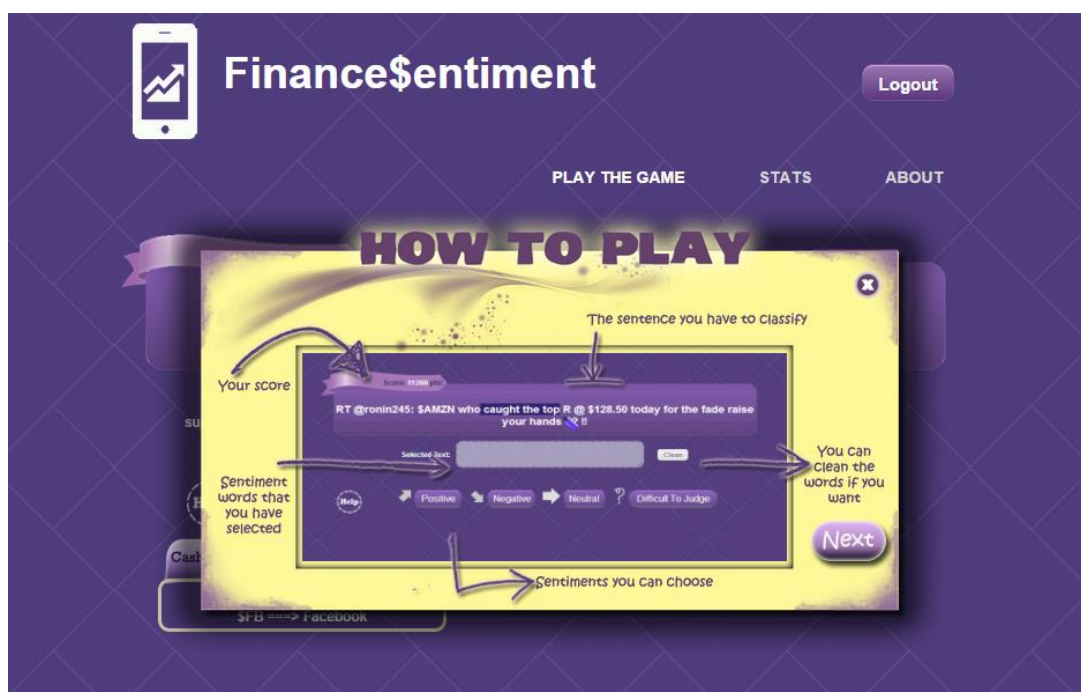


Figura 20 – Primeiro exemplo de como se joga

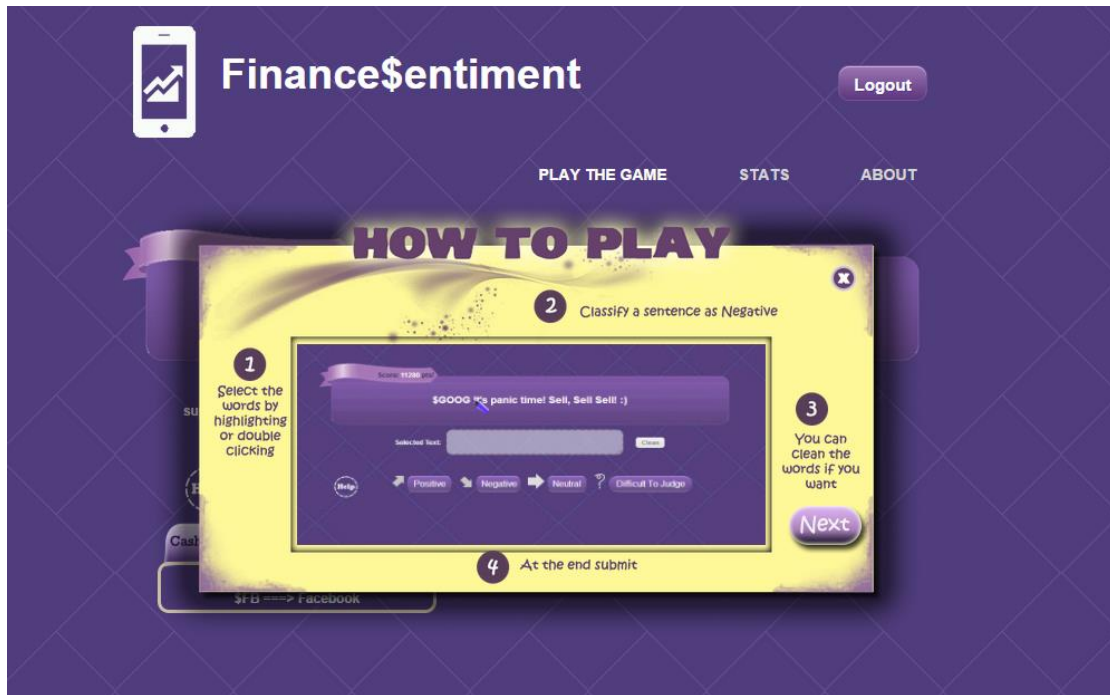


Figura 21 – Segundo exemplo de como se joga

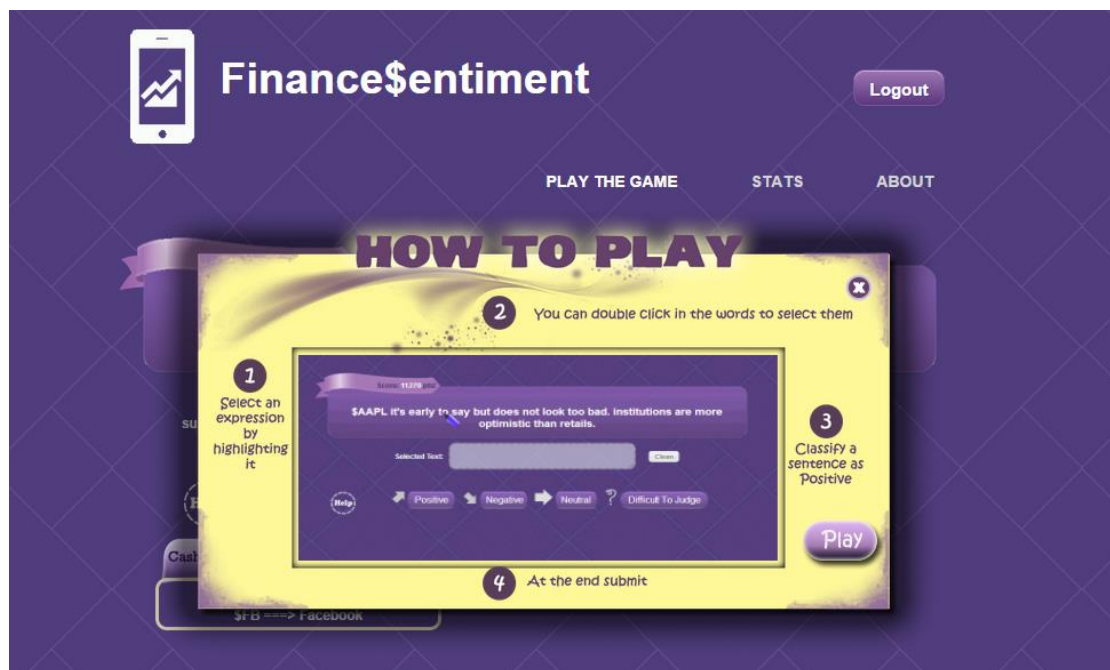


Figura 22 - 3º Exemplo de como se joga

Na Figura 23 pode visualizar-se a página de jogo, onde o utilizador pode realizar a sua contribuição. Nesta página são apresentadas as frases para classificação, uma a uma. O utilizador pode selecionar palavras ou expressões e que irão surgir na caixa de texto em baixo da

frase. Caso não tenha selecionado a palavra desejada pode sempre limpar o texto e selecionar novamente. Assim que o utilizador seleciona o sentimento expresso na frase, o botão para submeter a sua resposta surge como se pode verificar na Figura 24. Assim que a resposta do utilizador é submetida, é atribuída uma pontuação ao utilizador como se pode ver na Figura 25.

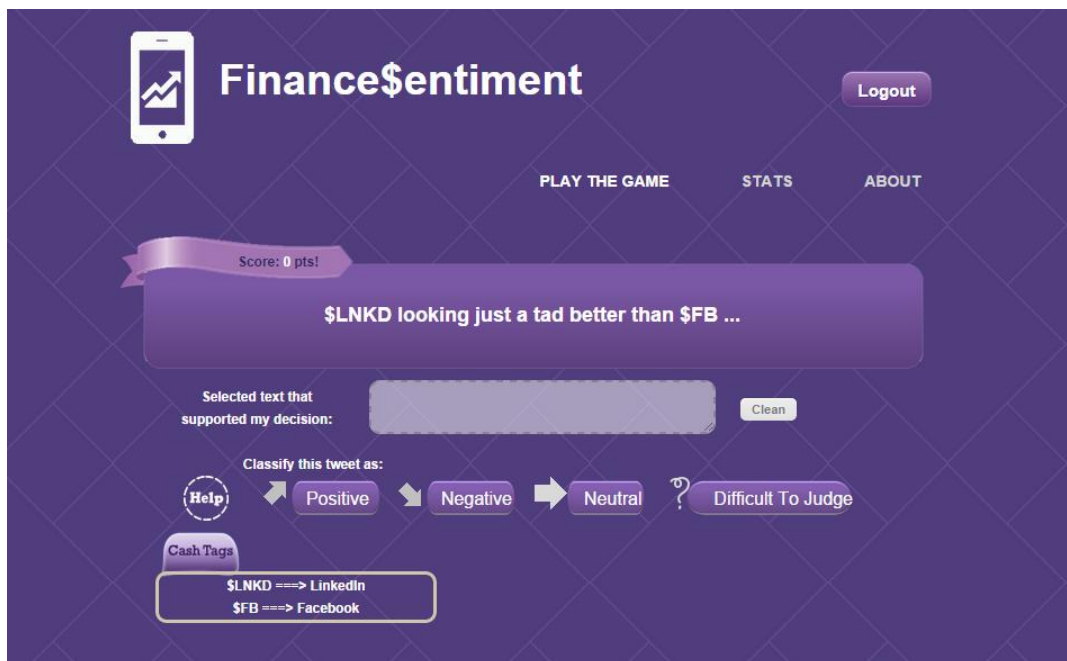


Figura 23 - Página de jogo

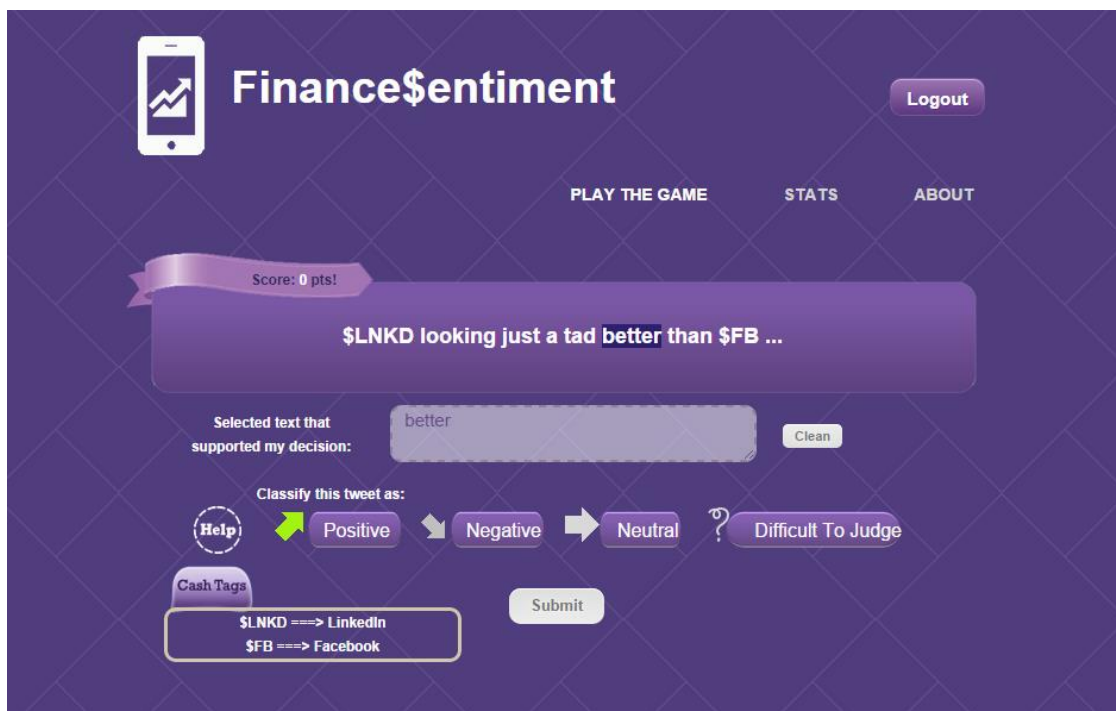


Figura 24 - Página de jogo

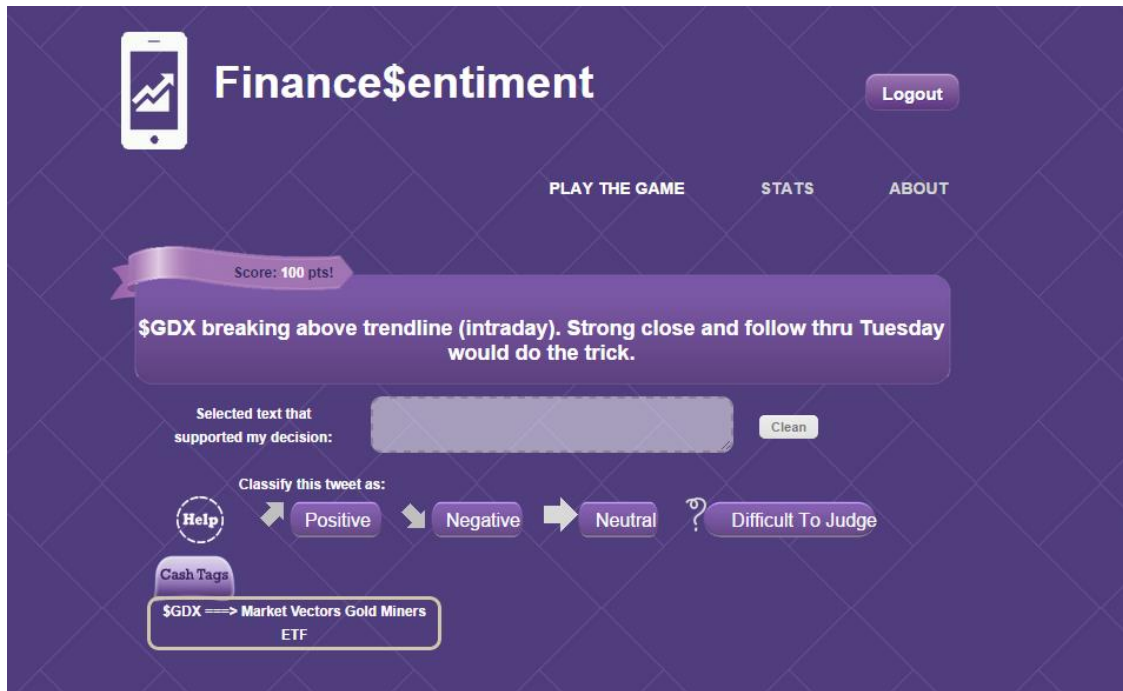


Figura 25 - Página de jogo

Para compreender melhor como funciona o jogo e a atribuição da pontuação o utilizador pode *clicar* em “Ajuda”, como se vê na Figura 26. Relativamente à explicação da pontuação esta pode ser vista na Figura 27. O utilizador ganha uma pontuação por cada frase classificada, tendo em conta se acertou ou errou na mesma. Para além disso, assim que classifica 10, 50, 100, 200, 500, 1.000 e 10.000 frases ganha uma medalha com uma pontuação associada para incentivar o utilizador a jogar. Assim que ganha uma medalha surge o aviso da Figura 28. As Figuras 29 e 30 mostram exemplos de medalhas ganhas.

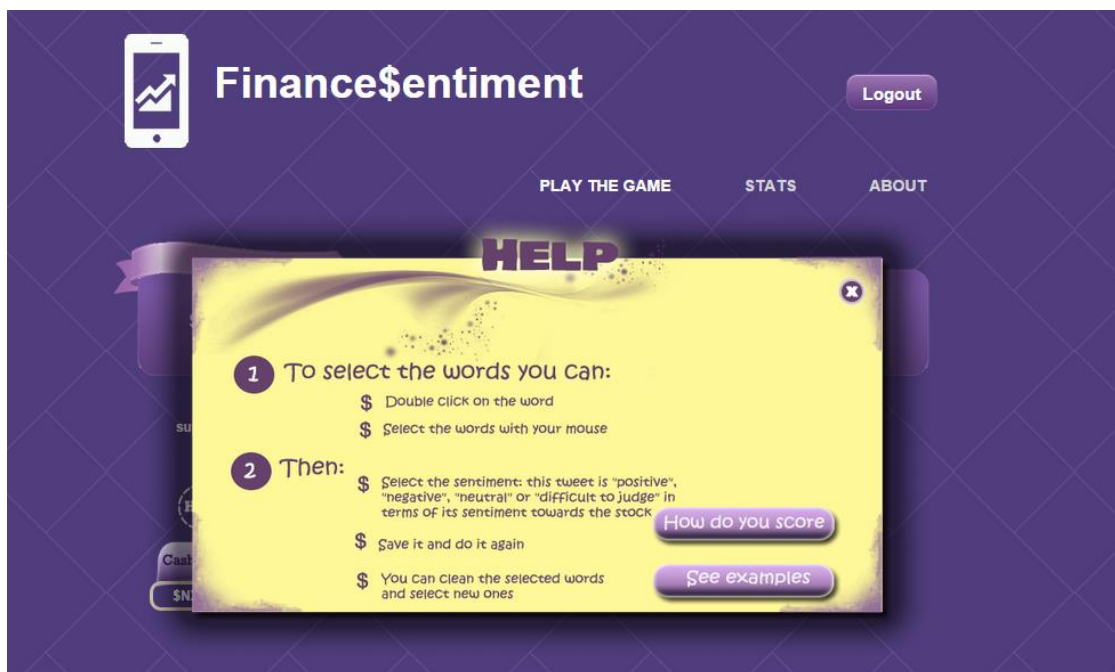


Figura 26 - Página de ajuda

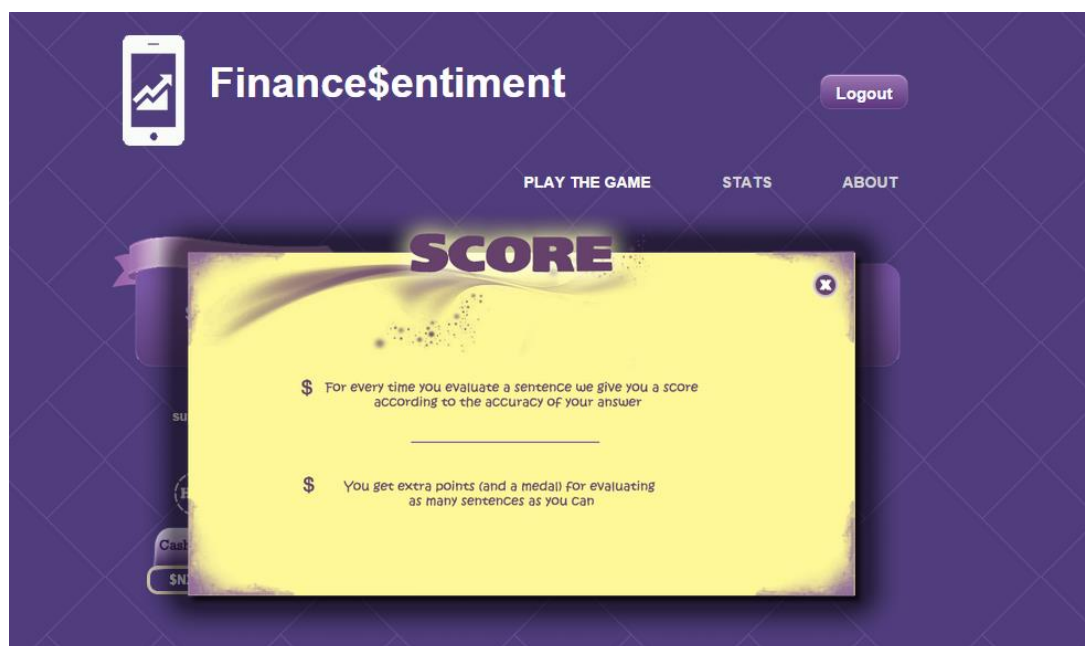


Figura 27 - Página de explicação da pontuação



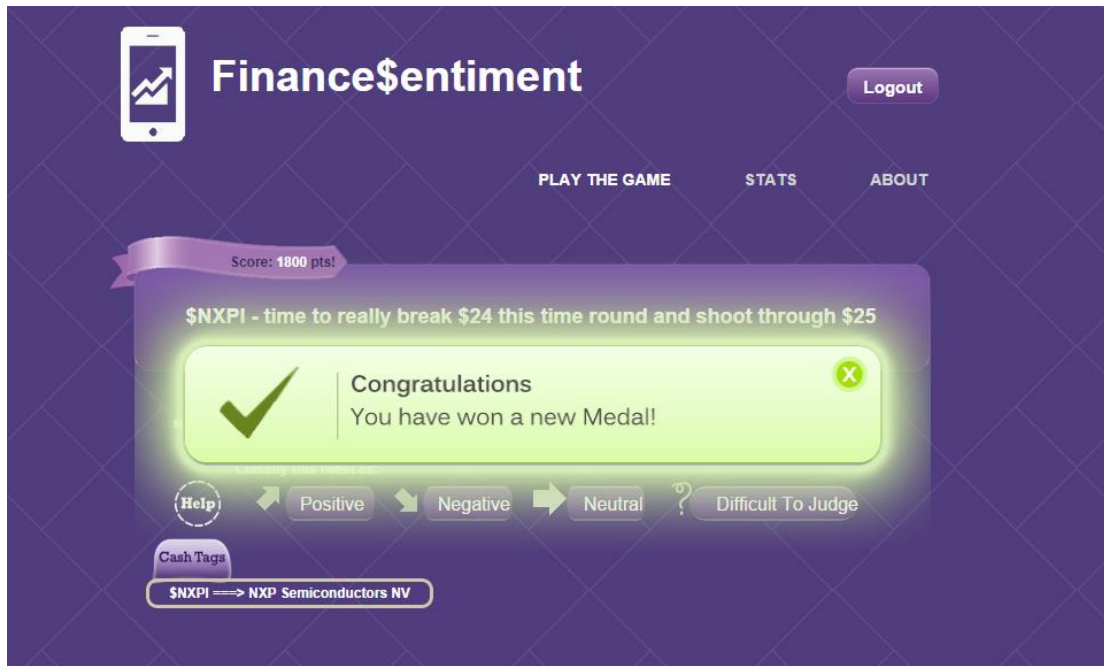


Figura 28 - Aviso de nova medalha

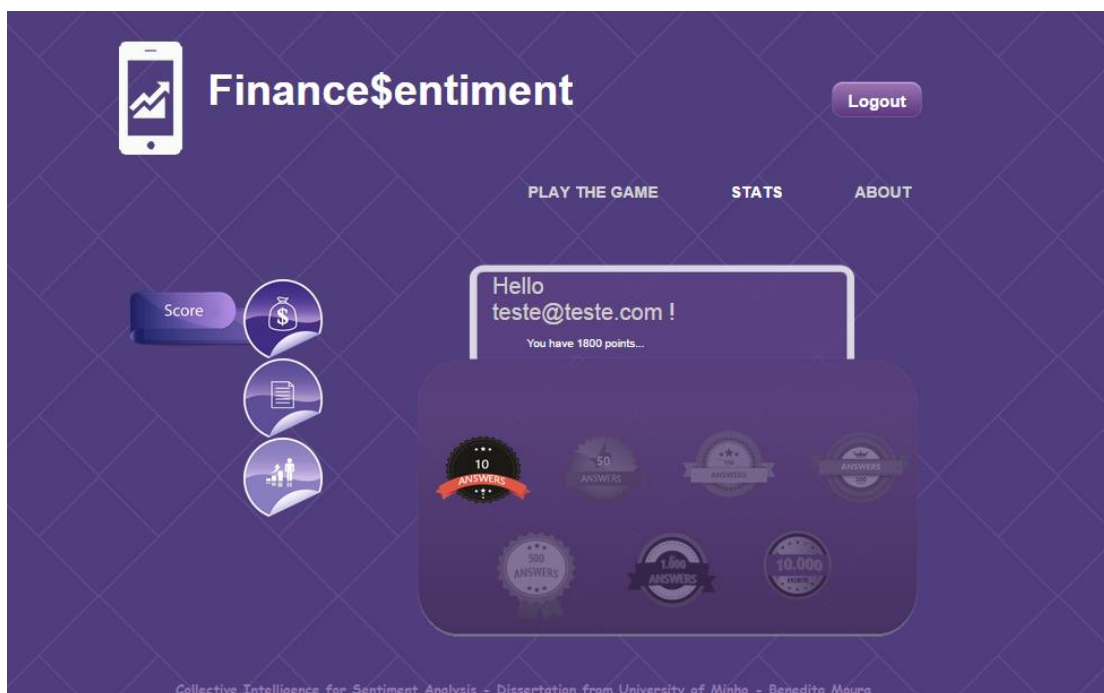


Figura 29 - Visualização das medalhas ganhas

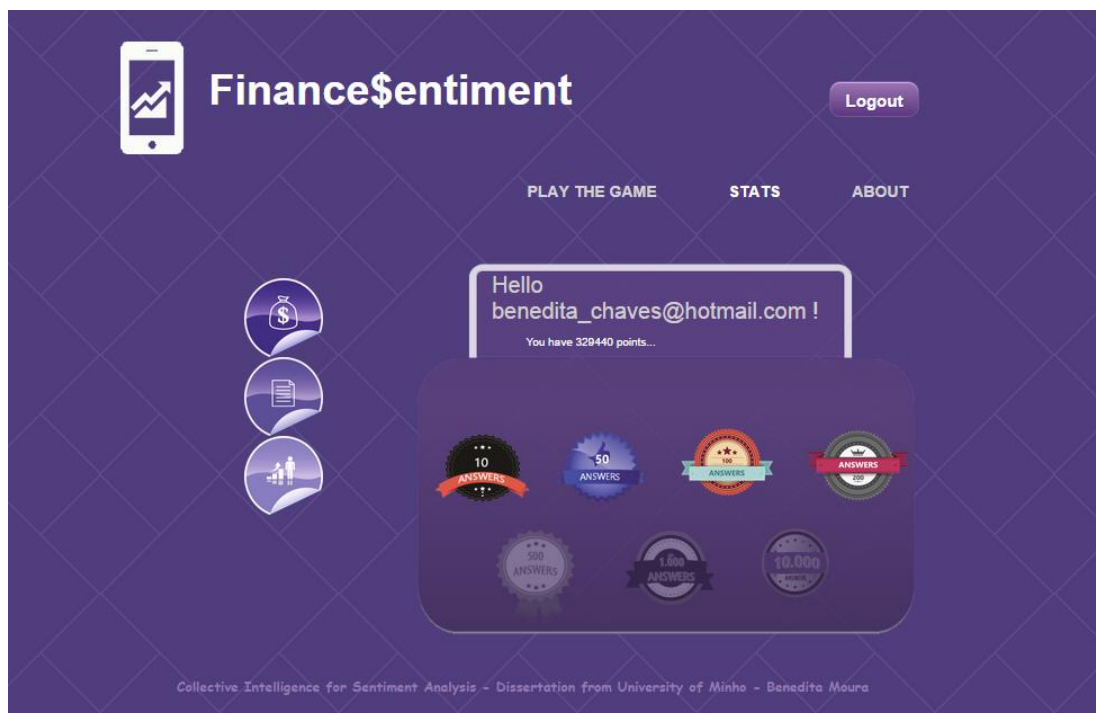


Figura 30 - Visualização das medalhas ganhas

Na Figura 31 podem ser vistas as classificações dadas pelo utilizador.

The screenshot shows the 'Answers' section of the application, displaying a table of user classifications for various financial sentences. The table has two columns: 'Sentence' and 'Sentiment'.

Sentence	Sentiment
SLNKD looking just a tad better than SFB ...	Positive
SGDX breaking above trendline (intraday). Strong close and follow thru Tuesday would do the trick.	Positive
SIBM long aug 18 200 call @ 1.65.. looking for push through 200 and hold	Positive
as long as \$spy stays above this purple desc trendline, I do plan to buy dips \$spx \$dji	Positive
SVRNG Watching closely for a break out of daily consolidation. \$3.73 is 50MA.	Positive
SVRNG Bullish Engulfing forming on the Weekly? 5 out of past 6 have resulted in higher prices.	Positive
STJX...how to stay with a strong long term trend	Positive
SCAB: wow, stumbled across this one flipping charts. i snooze i loose! hehe	Negative
SDX_F downside target to 81.30 area	Positive
\$SINA: watching for a possible dumpster diving trade, out of an ascending triangle .	Positive

Figura 31 - Classificações do utilizador

Na figura 32 pode visualizar-se o *Ranking* dos utilizadores da aplicação. O *Ranking* é realizado tendo em conta a pontuação dos participantes.

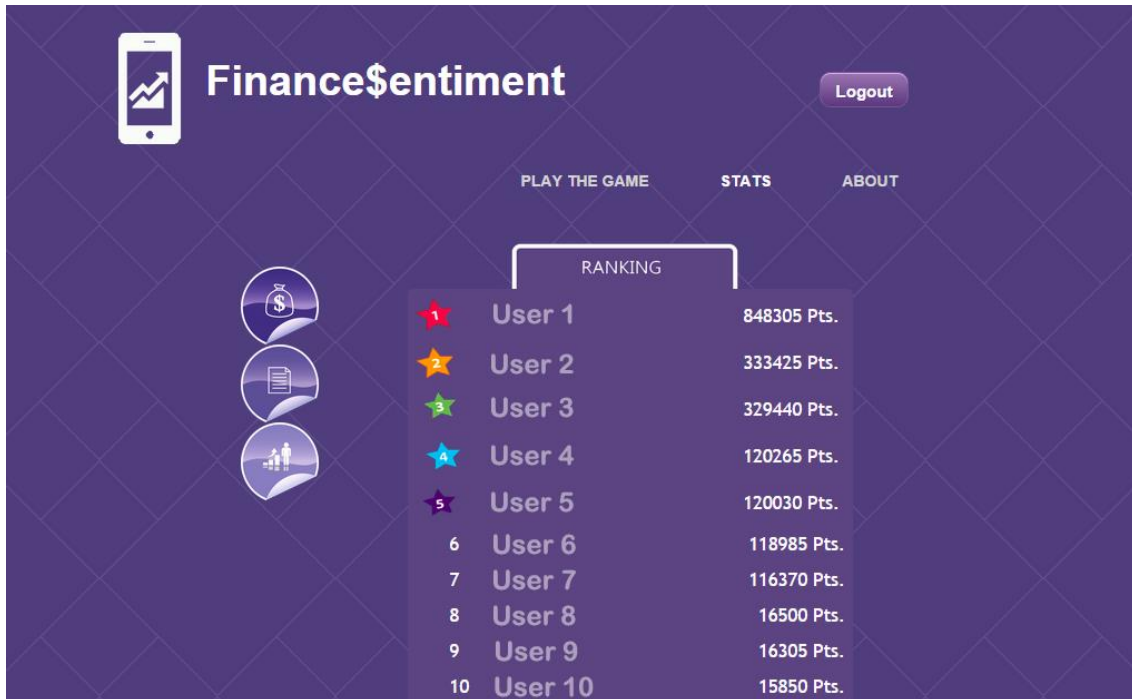


Figura 32 - *Ranking* dos utilizadores

#### 4.5.1.2 Visão do Administrador

Para haver um acompanhamento da evolução dos resultados da aplicação, foi criada uma secção apenas ao administrador da aplicação.

Esta secção é acedida a partir de uma página de *login* como se pode ver na Figura 33.

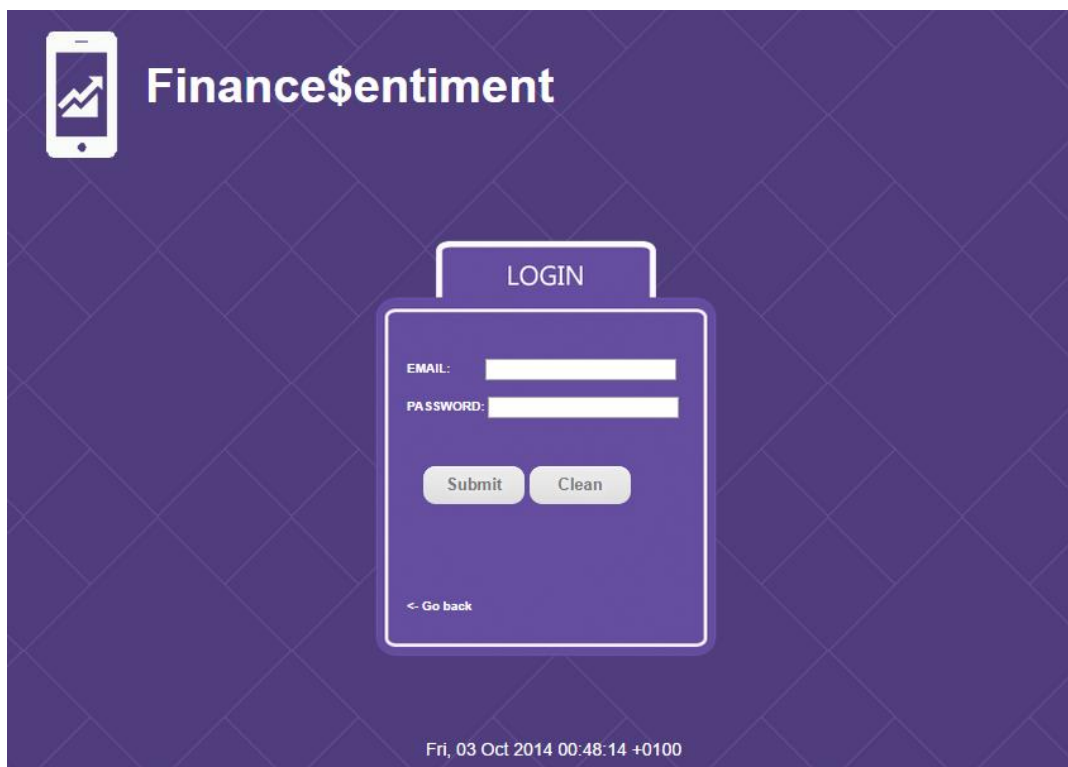


Figura 33- Página de *login* para aceder à secção do Administrador

Nesta secção podem ser visualizados diversos conteúdos, respeitantes às respostas dos utilizadores. Na Figura 34 pode ver-se todas as respostas dadas e ver os dados de cada uma. *Clicando* no número de respostas de uma dada frase, pode ver-se as classificações dadas, as palavras com sentimento seleccionadas assim como o sentimento médio. Este sentimento médio é obtido a partir de todas as classificações obtidas como pode ser visto na Figura 35.

Frase	Sentimento	Palavras	Número de respostas
\$JNPR	Positive	\$JNPR	3
\$RLTX \$IWM \$TNA \$TZA - Futures Show major shorting opportunity - pending Unemployment numbers.	Negative	\$RLTX \$IWM \$TNA \$TZA - Futures Show major shorting opportunity - pending Unemployment numbers.	3
\$FCX A close above trend line (blue line) may present the opportunity for a move towards the 36 area	Positive	\$FCX A close above trend line (blue line) may present the opportunity for a move towards the 36 area	3
\$zb_f 30 minute atr trend flips bearish	Negative	bearish	3
\$EURUSD maybe? short 1.2309 stop 1.2328 tp 1.2190. KIV FOMC ECB NFP	Negative	\$EURUSD maybe? short 1.2309 stop 1.2328 tp 1.2190. KIV FOMC ECB NFP	3
Unfortunately for the bulls...\$SPX \$SPY \$DIA \$ES_F	Negative	Unfortunately for the bulls	3
On watch list over level with volume of 20% ave full day \$ALNY	Positive	over level	3
\$DLPH showing signs of a potential reversal trend. Above 13d & 50d , positive MACD & RSI trending up	Positive	\$DLPH showing signs of a potential reversal trend. Above 13d & 50d , positive MACD & RSI trending up	3
\$IACI. Another EPS big gapper. Setting at pivot. Content player with nice accelerating sales.	Positive	\$IACI. Another EPS big gapper. Setting at pivot. Content player with nice accelerating sales.	3
\$DRYS Sorry , sending a better chart, Between 2 chairs (weekly).	Negative	\$DRYS Sorry , sending a better chart, Between 2 chairs	3

Figura 34 – Página onde se encontram todas as respostas dadas pelos utilizadores

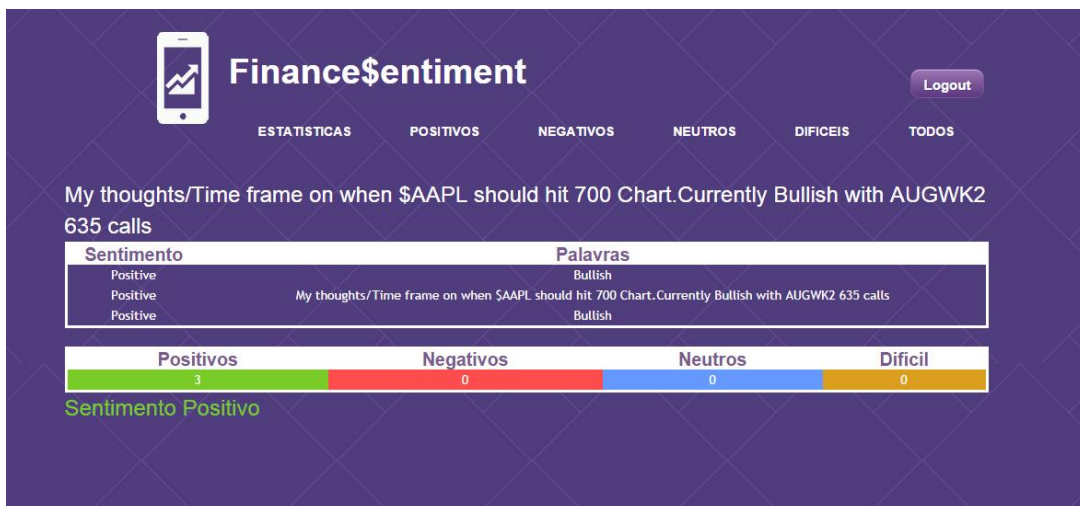
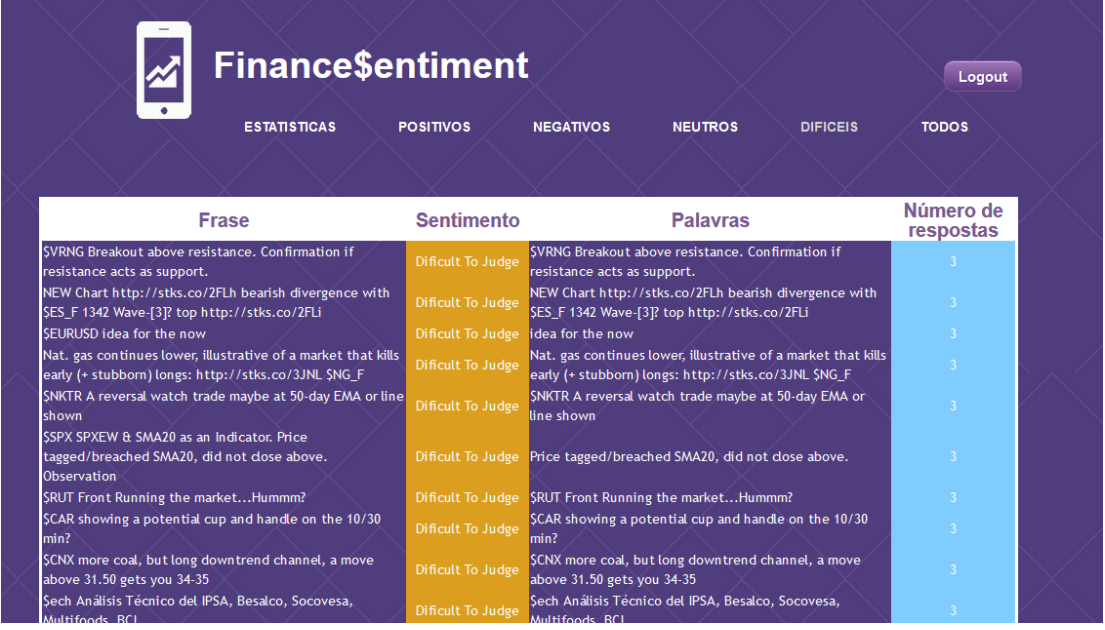


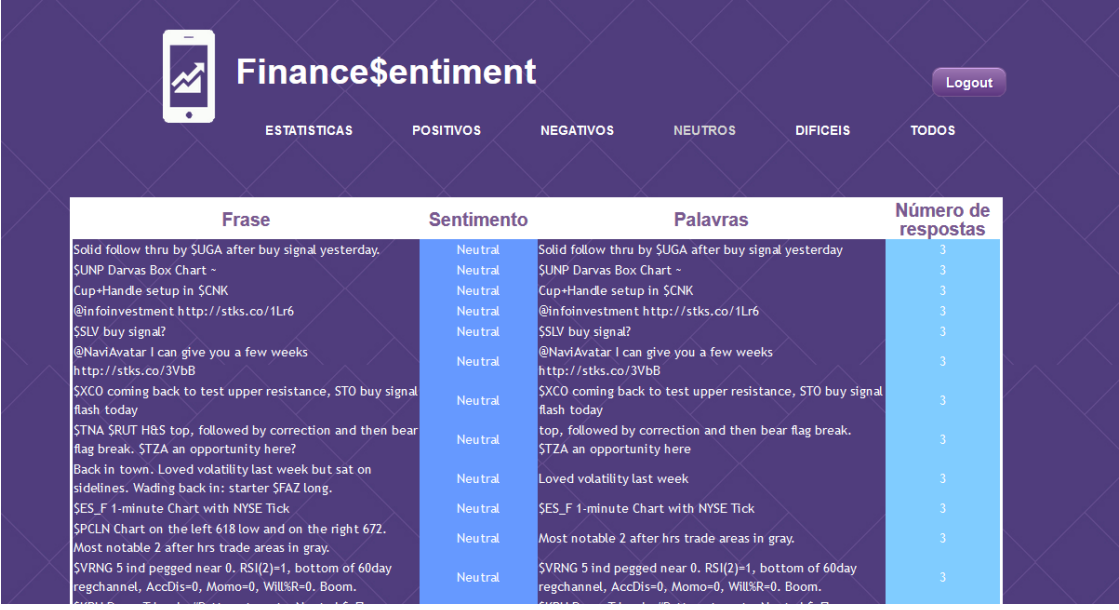
Figura 35 – Visualização das classificações e classificação média de determinada frase

As Figuras 36, 37, 38 e 39 descrevem-se exemplos de respostas cuja classificação foi “Difícil de Classificar”, “Neutro”, “Negativo” e “Positivo” respetivamente.



Frases	Sentimento	Palavras	Número de respostas
SVRNG Breakout above resistance. Confirmation if resistance acts as support.	Difícil To Judge	SVRNG Breakout above resistance. Confirmation if resistance acts as support.	3
NEW Chart <a href="http://stks.co/2FLh">http://stks.co/2FLh</a> bearish divergence with SES_F 1342 Wave-[3] top <a href="http://stks.co/2FLi">http://stks.co/2FLi</a>	Difícil To Judge	NEW Chart <a href="http://stks.co/2FLh">http://stks.co/2FLh</a> bearish divergence with SES_F 1342 Wave-[3] top <a href="http://stks.co/2FLi">http://stks.co/2FLi</a>	3
\$EURUSD Idea for the now	Difícil To Judge	Idea for the now	3
Nat. gas continues lower, illustrative of a market that kills early (+ stubborn) longs: <a href="http://stks.co/3JNL">\$NG_F</a>	Difícil To Judge	Nat. gas continues lower, illustrative of a market that kills early (+ stubborn) longs: <a href="http://stks.co/3JNL">\$NG_F</a>	3
\$NKTR A reversal watch trade maybe at 50-day EMA or line shown	Difícil To Judge	\$NKTR A reversal watch trade maybe at 50-day EMA or line shown	3
\$SPX \$PXEW & SMA20 as an Indicator. Price tagged/breached SMA20, did not close above. Observation	Difícil To Judge	Price tagged/breached SMA20, did not close above. Observation	3
\$RUT Front Running the market...Hummm?	Difícil To Judge	\$RUT Front Running the market...Hummm?	3
\$CAR showing a potential cup and handle on the 10/30 min?	Difícil To Judge	\$CAR showing a potential cup and handle on the 10/30 min?	3
\$CNX more coal, but long downtrend channel, a move above 31.50 gets you 34-35	Difícil To Judge	\$CNX more coal, but long downtrend channel, a move above 31.50 gets you 34-35	3
Sech Análisis Técnico del IPSA, Besalco, Socovesa, Multifoods- BCI	Difícil To Judge	Sech Análisis Técnico del IPSA, Besalco, Socovesa, Multifoods- BCI	3

Figura 36 – Exemplos de respostas cujo sentimento atribuído foi “Difícil de Classificar”



Frases	Sentimento	Palavras	Número de respostas
Solid follow thru by \$UGA after buy signal yesterday.	Neutral	Solid follow thru by \$UGA after buy signal yesterday	3
\$UNP Darvas Box Chart ~	Neutral	\$UNP Darvas Box Chart ~	3
Cup+Handle setup in \$CNK	Neutral	Cup+Handle setup in \$CNK	3
@infoinvestment <a href="http://stks.co/1Lr6">http://stks.co/1Lr6</a>	Neutral	@infoinvestment <a href="http://stks.co/1Lr6">http://stks.co/1Lr6</a>	3
\$SLV buy signal?	Neutral	\$SLV buy signal?	3
@NaviAvatar I can give you a few weeks <a href="http://stks.co/3VbB">http://stks.co/3VbB</a>	Neutral	@NaviAvatar I can give you a few weeks <a href="http://stks.co/3VbB">http://stks.co/3VbB</a>	3
\$XCO coming back to test upper resistance, STO buy signal flash today	Neutral	\$XCO coming back to test upper resistance, STO buy signal flash today	3
\$TNA \$RUT H&S top, followed by correction and then bear flag break. \$TZA an opportunity here?	Neutral	top, followed by correction and then bear flag break. \$TZA an opportunity here	3
Back in town. Loved volatility last week but sat on sidelines. Wading back in: starter \$FAZ long.	Neutral	Loved volatility last week	3
SES_F 1-minute Chart with NYSE Tick	Neutral	SES_F 1-minute Chart with NYSE Tick	3
\$PCLN Chart on the left 618 low and on the right 672. Most notable 2 after hrs trade areas in gray.	Neutral	Most notable 2 after hrs trade areas in gray.	3
SVRNG 5 ind pegged near 0. RSI(2)=1, bottom of 60day regchannel, AccDis=0, Momo=0, WillSR=0. Boom.	Neutral	SVRNG 5 ind pegged near 0. RSI(2)=1, bottom of 60day regchannel, AccDis=0, Momo=0, WillSR=0. Boom.	3

Figura 37 - Exemplos de respostas cujo sentimento atribuído foi “Neutro”

Frase	Sentimento	Palavras	Número de respostas
\$ALNY broke out of consolidation, more upside potential in the coming days - http://stks.co/1ZoC	Negative	broke out of consolidation	3
Massively oversold. RSI(2) = 0.34... Buying at open on green, hold for a day.	Negative	Massively oversold	3
\$EURUSD Thanks FED, Thanks ECB ...	Negative	\$EURUSD Thanks FED, Thanks ECB ...	3
\$WPRT Oversold. RSI(2) = 2.2...Would rather RSI<2, but will watch for quick scalp.	Negative	Oversold	3
\$EURUSD I am sure many are hoping the Bear Flag breaks. Patterns fail all the time. EMA8 Key.	Negative	Patterns fail all the time.	3
\$AAPL still has intermed weakness...filled gap on light vol. B tailed into close..break 601-604=Bull	Negative	intermed weakness light vol.	3
\$SPY We will just call this ECB fun this week ES 1387.50 down to 1349.50 (1346 lower channel line)	Negative	\$SPY We will just call this ECB fun this week ES 1387.50 down to 1349.50 (1346 lower channel line)	3
The grand-daddy index for the Nasdaq is bearish \$NDX \$COMPQ \$SPX \$SPY \$ES_F \$DIA	Negative	The grand-daddy index for the Nasdaq is bearish \$NDX \$COMPQ \$SPX \$SPY \$ES_F \$DIA	3
\$EXEL target 1- V	Negative	\$EXEL target 1- V	3
\$EURJPY #EURJPY Harmonics 31072012 #forex	Negative	\$EURJPY #EURJPY Harmonics 31072012 #forex	3
\$TEVA Oversold. RSI(2) = 3.6...Would rather RSI<2, but will watch for quick scalp.	Negative	\$TEVA Oversold. RSI(2) = 3.6...Would rather RSI<2, but will watch for quick scalp.	3

Figura 38 - Exemplos de respostas cujo sentimento atribuído foi “Negativo”

Frase	Sentimento	Palavras	Número de respostas
\$IBM long aug 18 200 call @ 1.65.. looking for push through 200 and hold	Positive	\$IBM long aug 18 200 call @ 1.65.. looking for push through 200 and hold	3
Last dump before new highs imo. Enjoy	Positive	Enjoy	3
ascending triangle b/o over the 200D in \$uan this AM. i'm long vs 24.50	Positive	ascending	3
\$EURUSD Secuencia alternativa para EurUsd y posible proyeccion para la proxima semana	Positive	\$EURUSD Secuencia alternativa para EurUsd y posible proyeccion para la proxima semana	3
\$SPY although its been in an upward channel, check volume. Its pretty much flat.	Positive	\$SPY although its been in an upward channel, check volume. Its pretty much flat.	3
\$BIDU: another one that i am watching dosely. if they show up, so will i.	Positive	\$BIDU: another one that i am watching dosely. if they show up, so will i.	3
\$CSTR clean break out of the Buy Zone, expect EMA100 tag next.	Positive	\$CSTR clean break out of the Buy Zone, expect EMA100 tag next.	3
\$WY - Trade Signal New - BTO @ 23.59 - Limit w/21.03 - Stop on Wkly Chart w/.5N @ Risk -	Positive	\$WY - Trade Signal New - BTO @ 23.59 - Limit w/21.03 - Stop on Wkly Chart w/.5N @ Risk -	3
Last dump before new highs imo. Enjoy	Positive	Last dump before new highs imo. Enjoy	3
\$AUDJPY daily chart explains it all. It never broke down below the consolidated area	Positive	\$AUDJPY daily chart explains it all. It never broke down below the consolidated area	3

Figura 39 - Exemplos de respostas cujo sentimento atribuído foi “Positivo”

As Figuras 40, 41 e 42 mostram algumas estatísticas, utilizadas para monitorizar a evolução das classificações. A Figura 41 mostra os resultados obtidos de acordo com cada classificação e mostra os gráficos da sua assertividade. Na Figura 42 os dados são semelhantes aos da página anterior, contudo os dados são obtidos tendo em conta a média das classificações de cada frase.



Figura 40 – Estatísticas gerais da aplicação

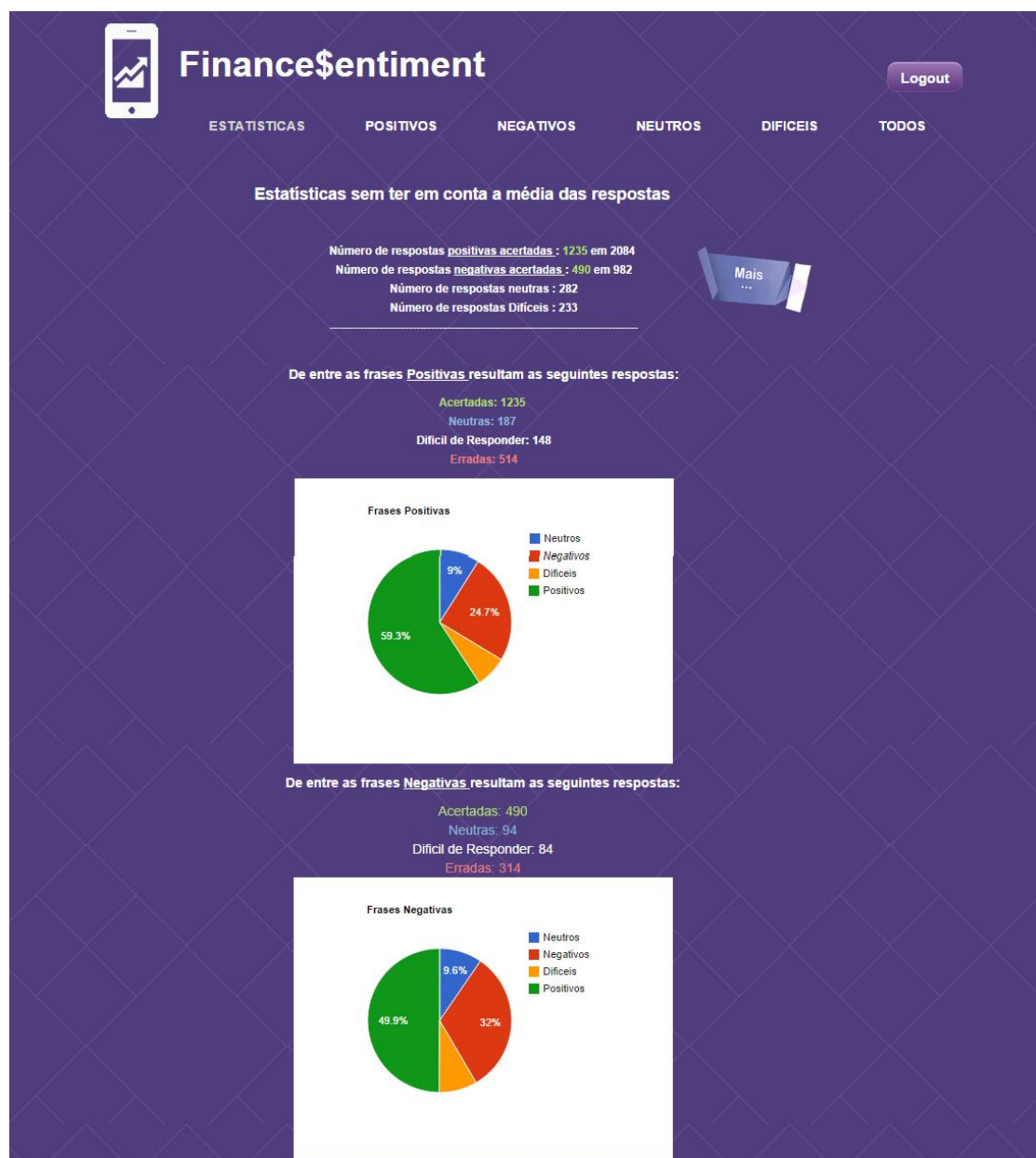


Figura 41 – Estatísticas das classificações individuais



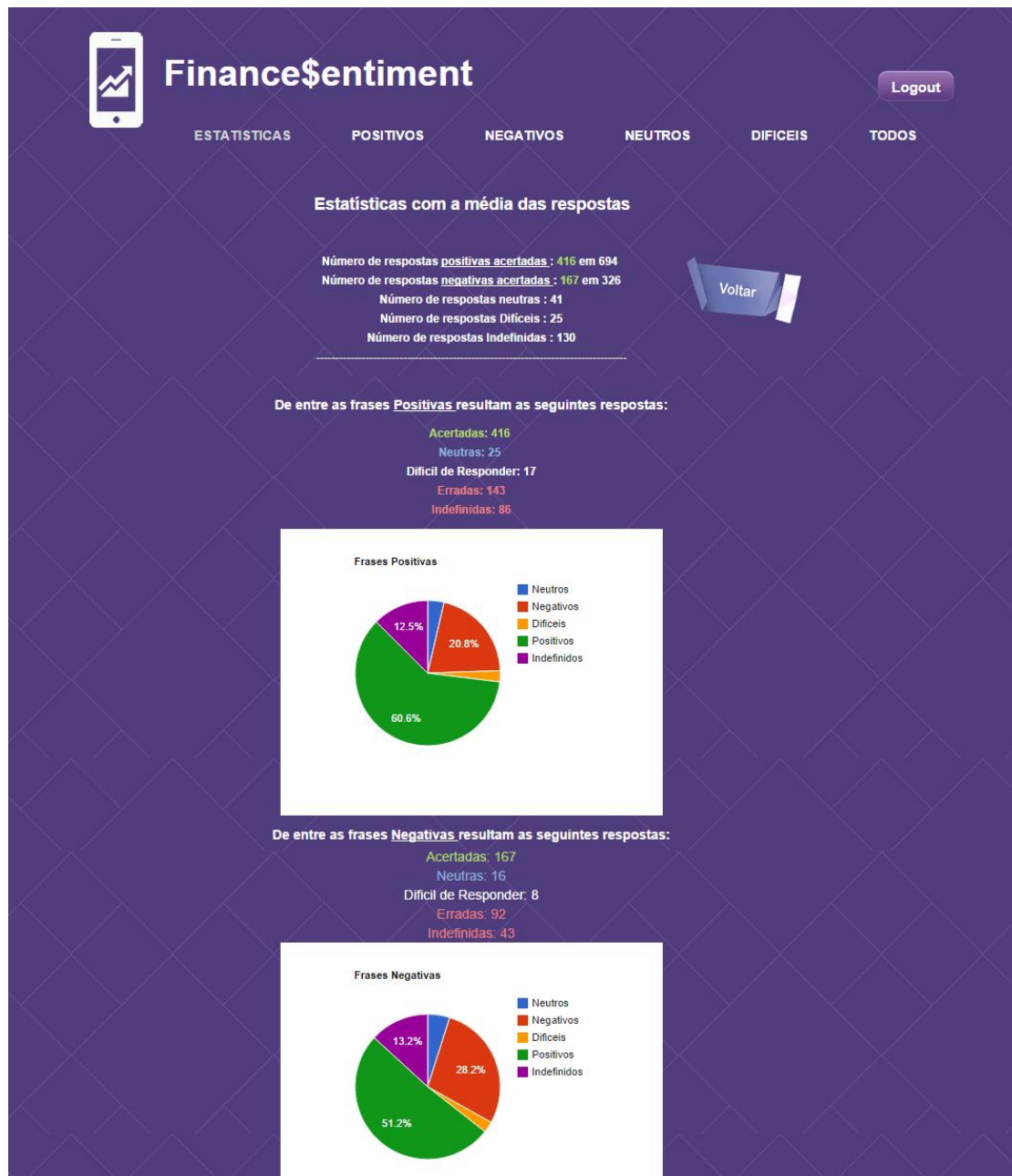


Figura 42 – Estatísticas das frases tendo em conta a média das classificações

#### 4.5.1.3 Descrição das Funcionalidades Implementadas

Na Figura 43 são descritos os casos de uso para a aplicação *Finance\$entiment* assim como os seus utilizadores.

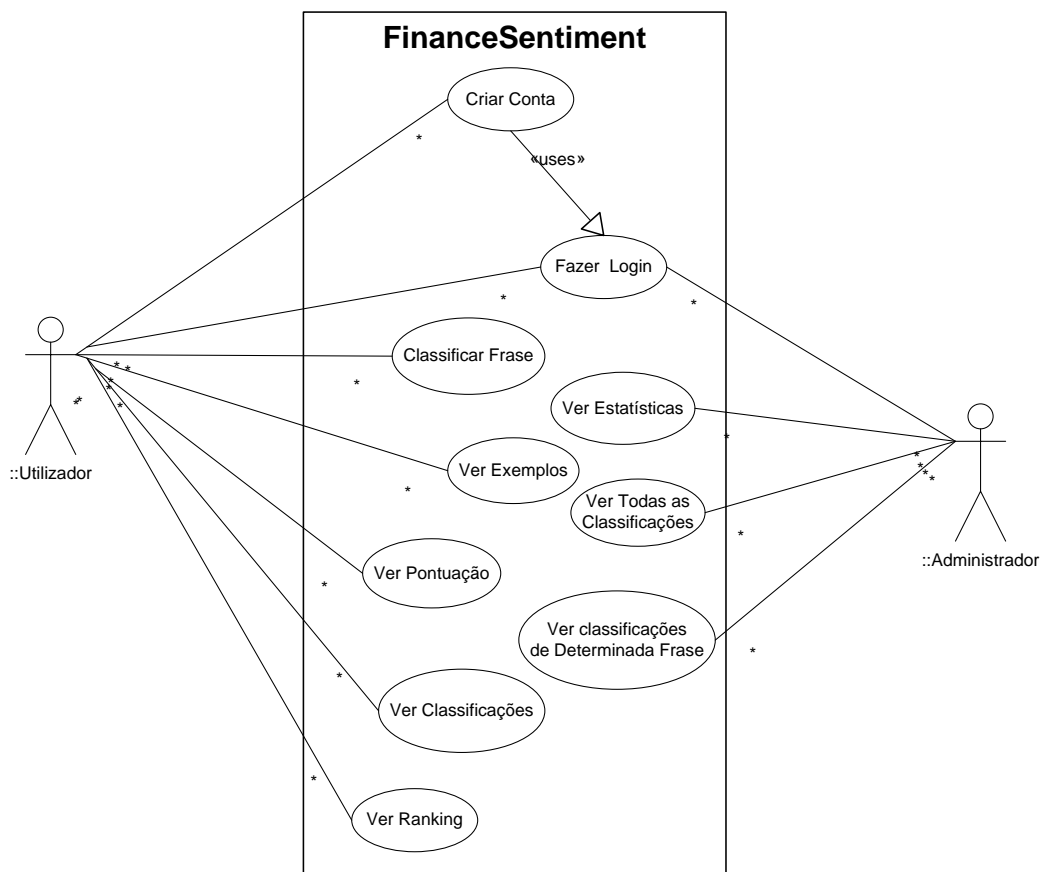


Figura 43 - Casos de Uso (UML) da aplicação *FinanceSentiment*

Na aplicação existem dois tipos de utilizadores, o utilizador comum, cuja função é classificar frases e o administrador, cuja função é verificar a evolução das classificações.

O utilizador comum necessita de criar uma conta e fazer *login* no sistema para poder aceder à aplicação. Para um utilizador que nunca se registou estas duas atividades são realizadas em simultâneo. O utilizador só se regista uma vez e a partir desse momento pode aceder à sua conta sempre que quiser com a simples introdução do *email* que registou.

O utilizador comum possui atividades como classificar frases, ver exemplos, ver a sua pontuação, ver as classificações realizadas e visualizar a sua posição no *Ranking* geral da aplicação. Sendo a atividade de classificar frases a principal da aplicação, as outras servem de apoio e incentivo para a realização desta.

No caso do administrador, este deve fazer *login* para aceder ao sistema através do seu *username* e *password*.

Este possui atividades de controlo como visualizar estatísticas, ver todas as classificações, assim como visualizar todas as classificações feitas a determinada frase e qual o sentimento medio resultante.

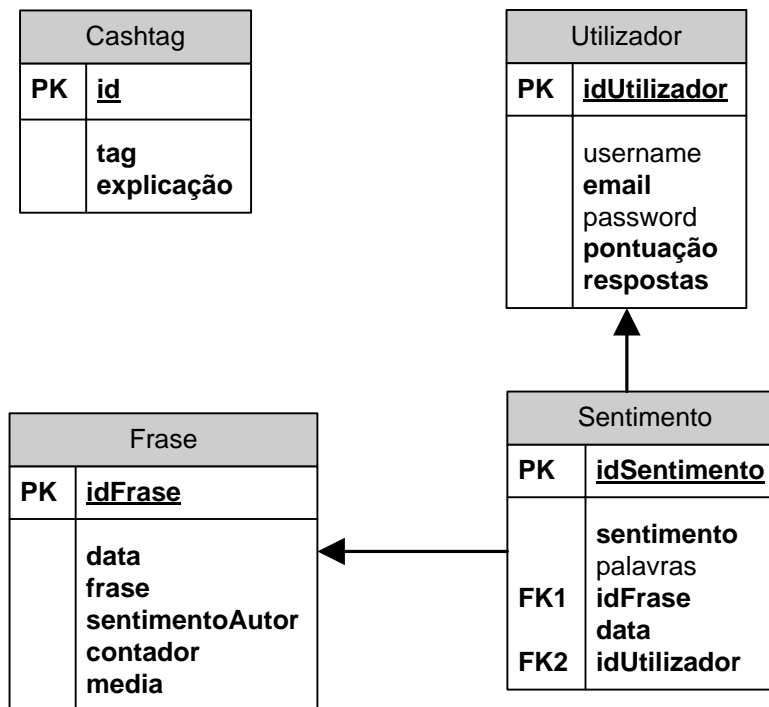


Figura 44 - Estrutura da Base de Dados da aplicação *Finance\$entiment*

Na Figura 44 são descritas as tabelas da Base de Dados da aplicação. A tabela “Frase” armazena todas as frases a ser classificadas. Esta possui um total de 4.179 registos, contudo apenas 1.028 foram classificadas e 1.019 apresentam três classificações. Este possui campos como a “data” em que foi escrita, a “frase” propriamente dita, o “sentimentoAutor”, campo que como o próprio nome indica, armazena o sentimento que o autor da frase lhe atribuiu. Este campo é importante pois é através do mesmo que se verifica a assertividade das classificações realizadas pelos utilizadores. Esta tabela possui mais dois campos, o “contador” e a “média”. O primeiro armazena o número de vezes que dada frase foi classificada, pois cada frase deve ser classificada três vezes. A “média” armazena o sentimento médio de cada frase, obtido tendo em conta as três classificações realizadas.

A tabela “Cashtag” possui 8.948 registos relativos à explicação das *tags* financeiras encontradas nas frases. Para além do seu identificador “id”, possui os campos “tag” e

“explicação”. Quando é apresentada uma dada frase ao utilizador, que possui uma dada *tag*, a respetiva explicação é apresentada. Esta explicação é uma pequena ajuda para os utilizadores que têm menos conhecimentos financeiros.

A tabela “Utilizador” armazena todos os utilizadores da aplicação, assim como os seus dados. Esta possui 121 registos, embora apenas 89 desses utilizadores possuam respostas. Os campos “username” e “password” armazenam a informação do administrador, pois para o registo e *login* de um utilizador comum é apenas necessário o campo “email”. O campo “pontuação” diz respeito à pontuação obtida ao longo das classificações e “respostas” indica o número de classificações realizadas.

A tabela “Sentimento” é aquela que guarda as classificações realizadas pelos utilizadores, com um total de 3.057 registos. Esta apresenta chaves estrangeiras de duas tabelas, ou seja apresenta o “idUtilizador” da tabela “Utilizador” e “idFrase” da tabela “Frase”. Assim cada classificação realizada está associada a uma frase e também a um utilizador. Cada registo possui um “sentimento” e “data” e pode apresentar “palavras” ou não. Este campo representa as palavras com sentimento selecionadas pelo utilizador.

#### 4.5.2 Divulgação

A aplicação foi divulgada no ambiente universitário, por um período de tempo de dois meses entre Julho e Agosto de 2014. Este período foi o necessário para obter os dados mínimos necessários para a investigação. Esta foi divulgada por meio de redes sociais como o *LinkedIn*, através da comunicação direta dos professores Paulo Cortez e Nelson Areal aos seus alunos e colegas de trabalho. Contudo a partilha desta aplicação teve algumas restrições devido à utilização de dados da plataforma *StockTwits*, e apesar do seu uso ser apenas académico, não se obteve permissão para uma maior partilha. Assim, não se tenha obtido um valor elevado de adesão. De qualquer modo, conseguiu-se obter um conjunto de avaliações consideradas interessantes para análise dentro deste trabalho.

## 4.6 Análise de Resultados

### 4.6.1 Análise dos Utilizadores

Relativamente aos utilizadores da aplicação, obtiveram-se os dados gerais apresentados Figura 43. Dos 121 utilizadores registados, apenas 89 realizaram classificações. Destes 89, 48 realizaram menos de 10 classificações, 30 realizaram entre 10 e 50 classificações, 5 possuem entre 50 e 100 classificações e apenas 6 possuem mais de 100 classificações. Com 3.057 classificações realizadas, ou seja, 1.019 frases classificadas, os utilizadores realizaram em média 34 classificações.

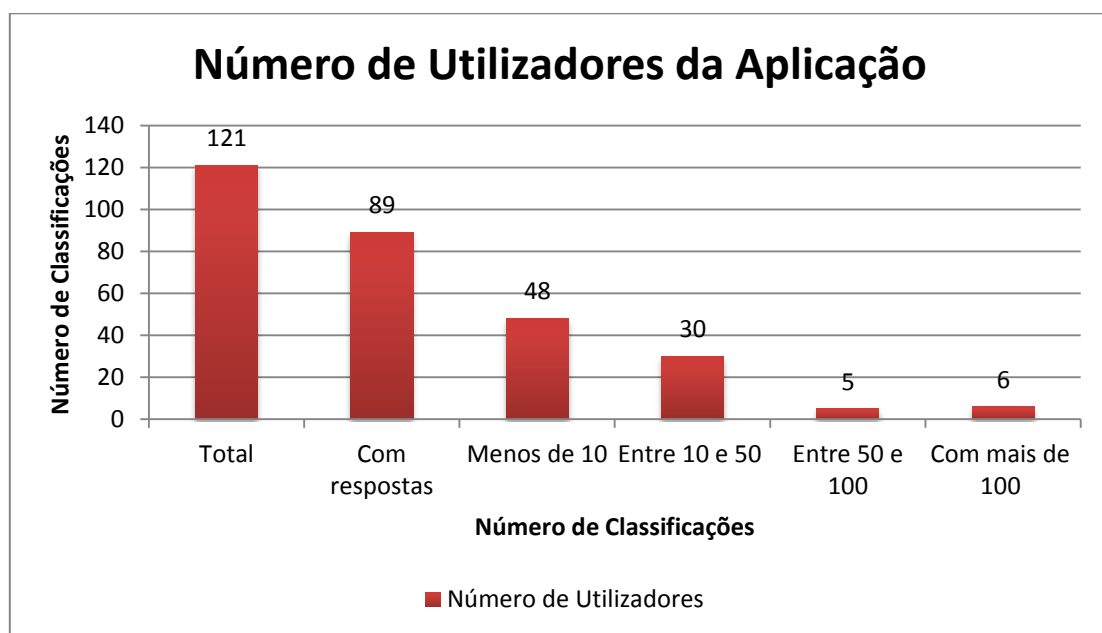


Figura 45- Relação entre número de Utilizadores e número de classificações

#### 4.6.1.1 Utilizadores com Menos de 10 Respostas

Relativamente aos utilizadores com menos de 10 respostas obtiveram-se os seguintes dados expostos nas Figuras 46 e 47.

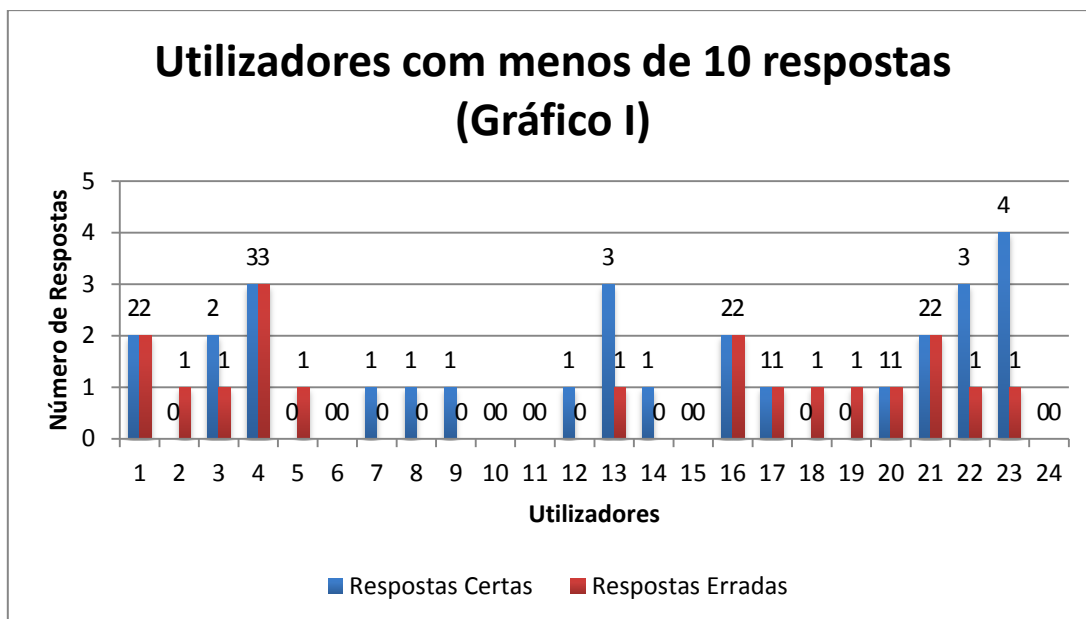


Figura 46 -Utilizadores com menos de 10 respostas (Gráfico I)

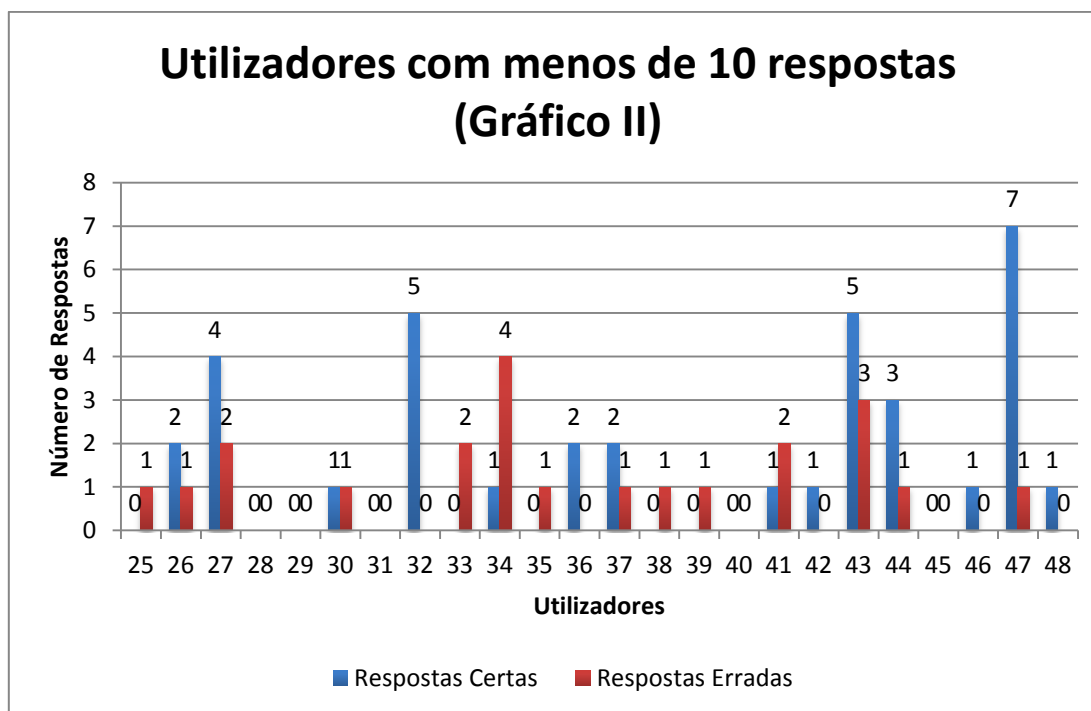


Figura 47 - Utilizadores com menos de 10 respostas (Gráfico II)

Relativamente a estes resultados pouco se pode concluir, pois estes utilizadores realizam poucas classificações. Contudo pode verificar-se que na sua maioria os utilizadores acertam mais classificações do que aquelas que erram. Alguns utilizadores não possuem nenhuma

resposta correta nem nenhuma resposta errada pois realizaram as classificações “Neutra” ou “Difícil de Classificar”, que não são consideradas nem certas nem totalmente erradas.

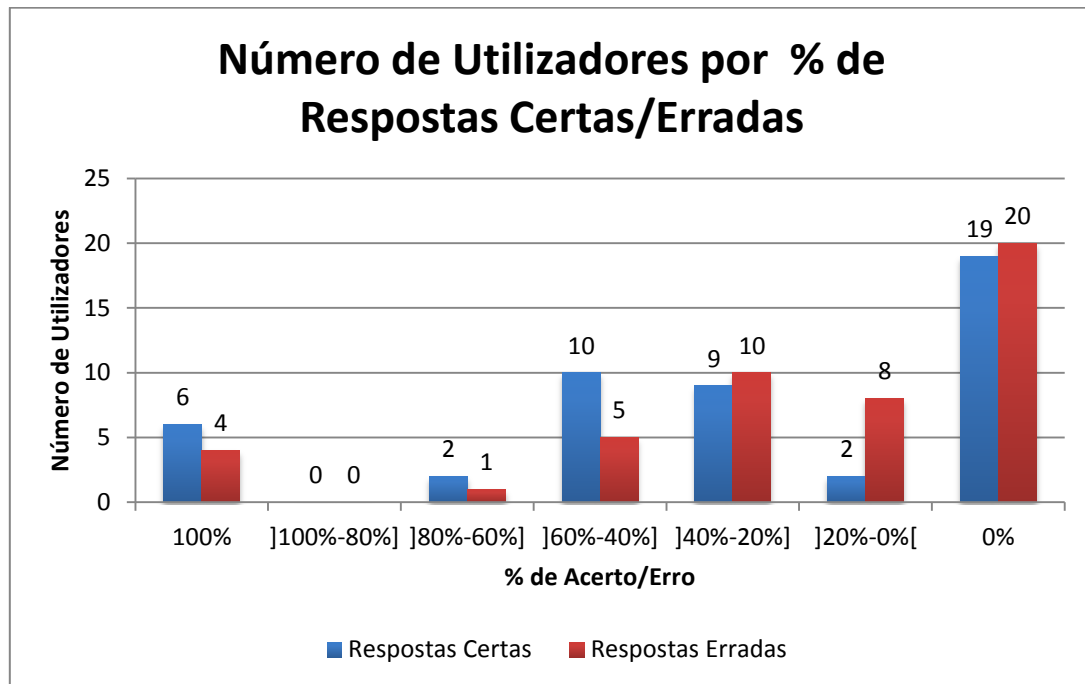


Figura 48 - Número de Utilizadores por % de Respostas Certas/Erradas

A Figura 48 diz respeito à quantidade de utilizadores (com menos de 10 classificações) que conseguiram acertar/errar as percentagens apresentadas, das classificações realizadas. Pode verificar-se que houve um grande número de utilizadores a obter uma percentagem de acerto de 0%, contudo o contrário também se verifica, pois o número de utilizadores com uma percentagem de erro igual a 0% também é elevado. A maior parte dos utilizadores teve uma percentagem de erro menor que 40% enquanto a percentagem de acerto se encontra um pouco mais distribuída.

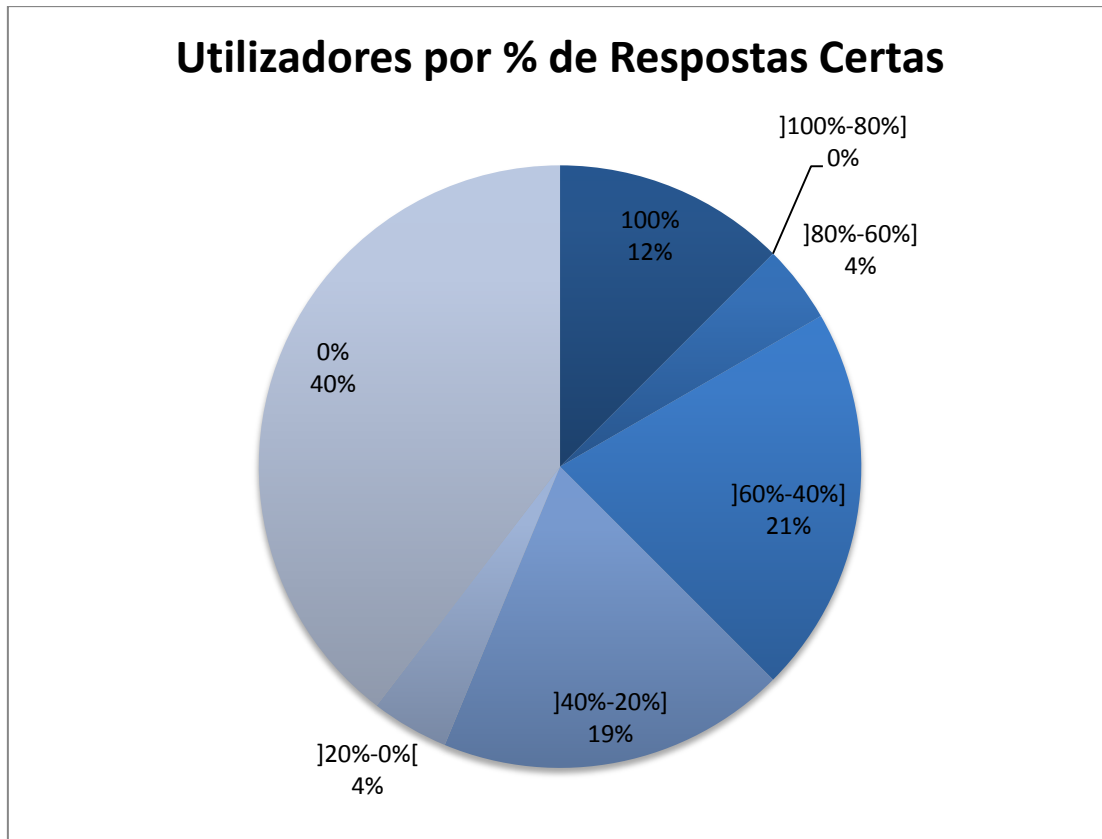


Figura 49 - Utilizadores por % de Respostas Certas

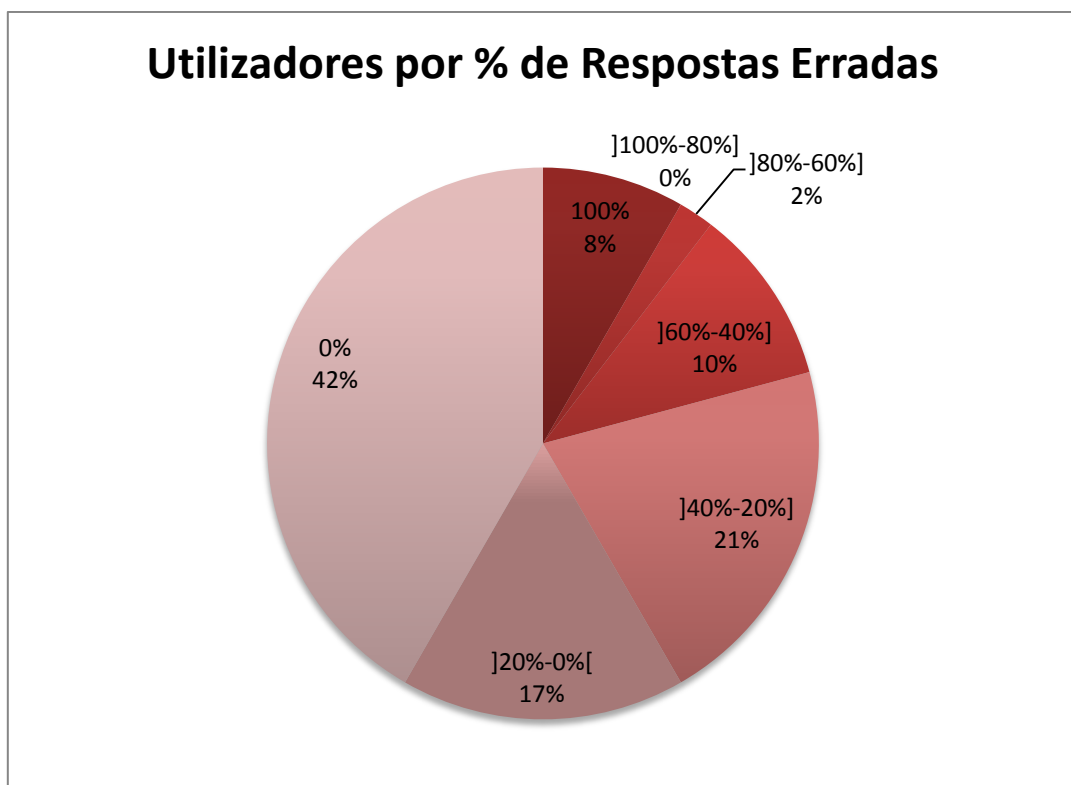


Figura 50 - Utilizadores por % de Respostas Erradas



As Figuras 49 e 50 pode verificar-se novamente quantos utilizadores obtiveram determinada percentagem de acerto e erro separadamente. Verifica-se novamente que há um maior número de utilizadores que acertam em 0% das suas classificações, apesar da sua média de acerto se encontrar bastante distribuída. Para além disso a maioria dos utilizadores apresenta uma percentagem de acerto menor que 40% e uma percentagem de erro menor que 20%

#### 4.6.1.2 Utilizadores com Mais de 10 e Menos de 50 Respostas

Os utilizadores com mais de 10 e menos de 50 respostas obtiveram os seguintes dados representados nas Figuras 49 e 50.

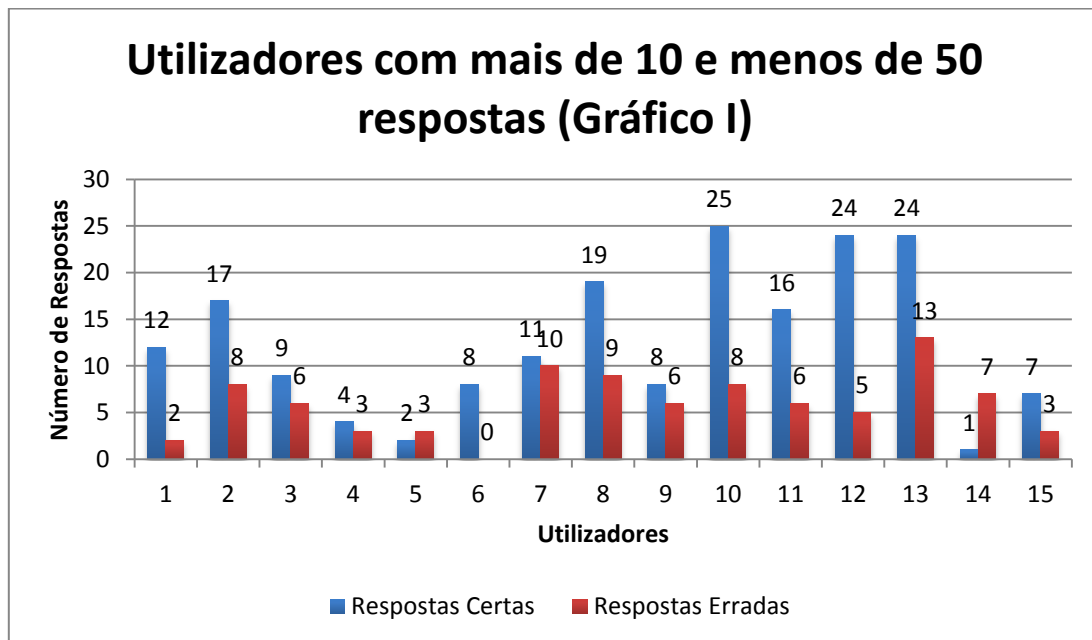


Figura 51 - Utilizadores com mais de 10 e menos de 50 respostas (Gráfico I)

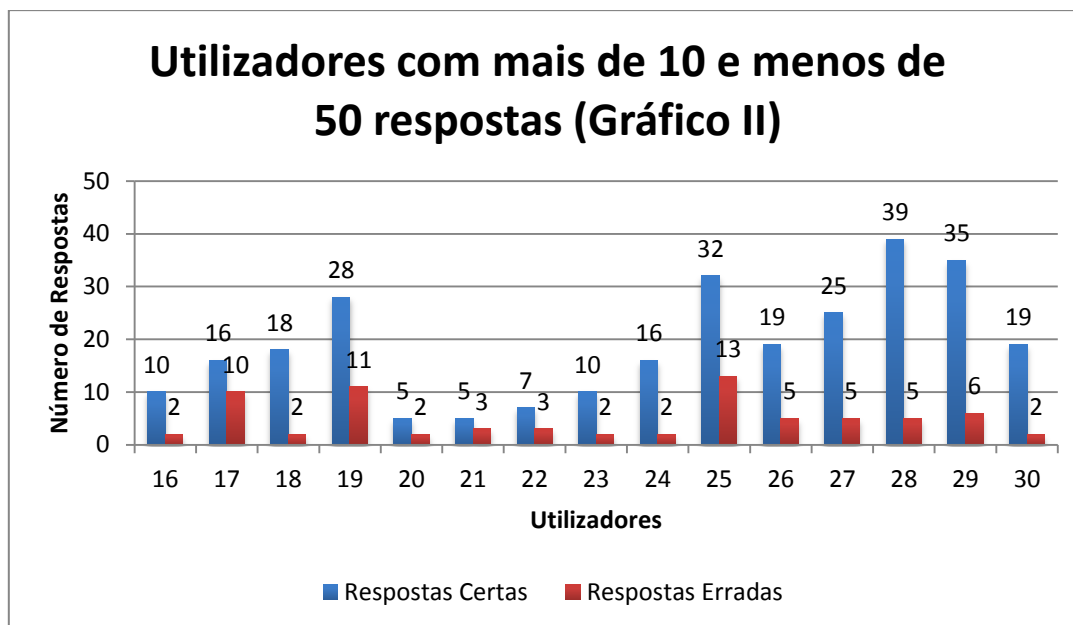


Figura 52 - Utilizadores com mais de 10 e menos de 50 respostas (Gráfico II)

Dos gráficos apresentados nas Figuras 51 e 52, pode concluir-se que apenas 2 utilizadores em 30 obtiveram mais classificações erradas do que certas. Em média quase todos os utilizadores acertaram relativamente bastante mais do que erraram.

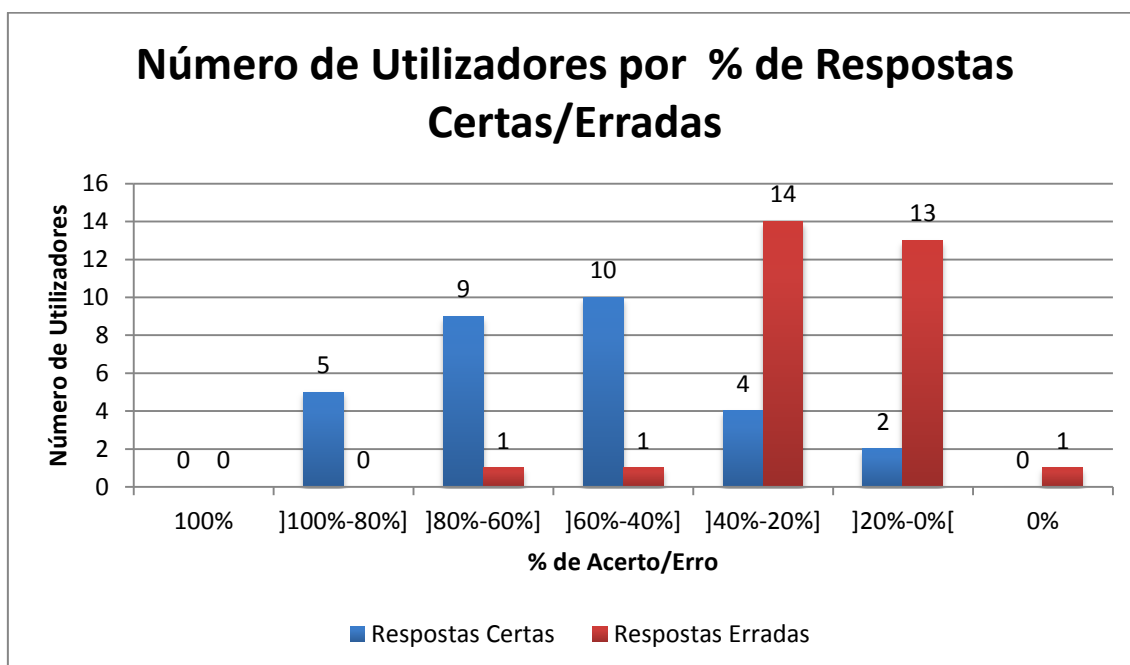


Figura 53 - Número de Utilizadores por % de Respostas Certas/Erradas



Figura 54 - Utilizadores por % de Respostas Certas

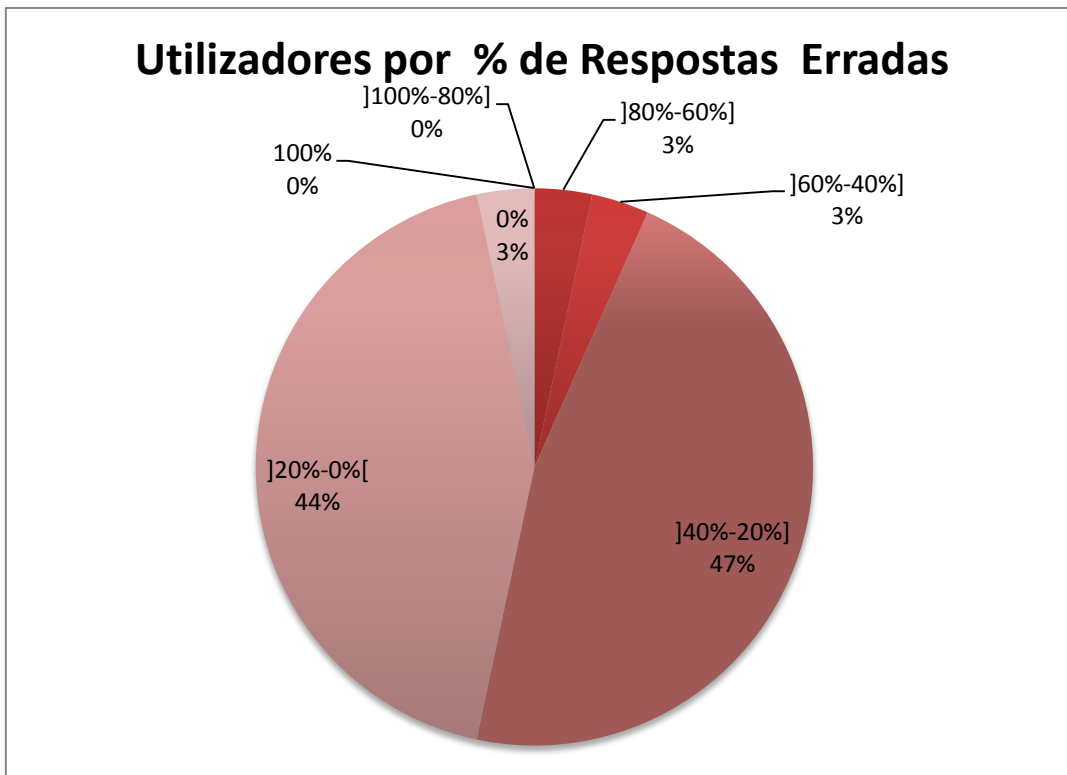


Figura 55 - Utilizadores por % de Respostas Erradas

Relativamente aos gráficos expostos nas Figuras 53, 54 e 55, pode concluir-se que em média as percentagens de acerto com mais utilizadores são entre os 80% e os 40%. No caso dos utilizadores com mais de 10 e menos de 50 classificações pode verifica-se que existe um número muito reduzido de utilizadores com menos de 20% de respostas certas. Para além disso os utilizadores possuem na sua maioria de 40% a 0% de respostas erradas. Os utilizadores com percentagens de erro fora deste intervalo são mínimos.

#### 4.6.1.3 Utilizadores com Mais de 50 e Menos de 100 Respostas

Dos utilizadores com mais de 50 e menos de 100 respostas obtiveram os seguintes dados expostos na Figura 56.

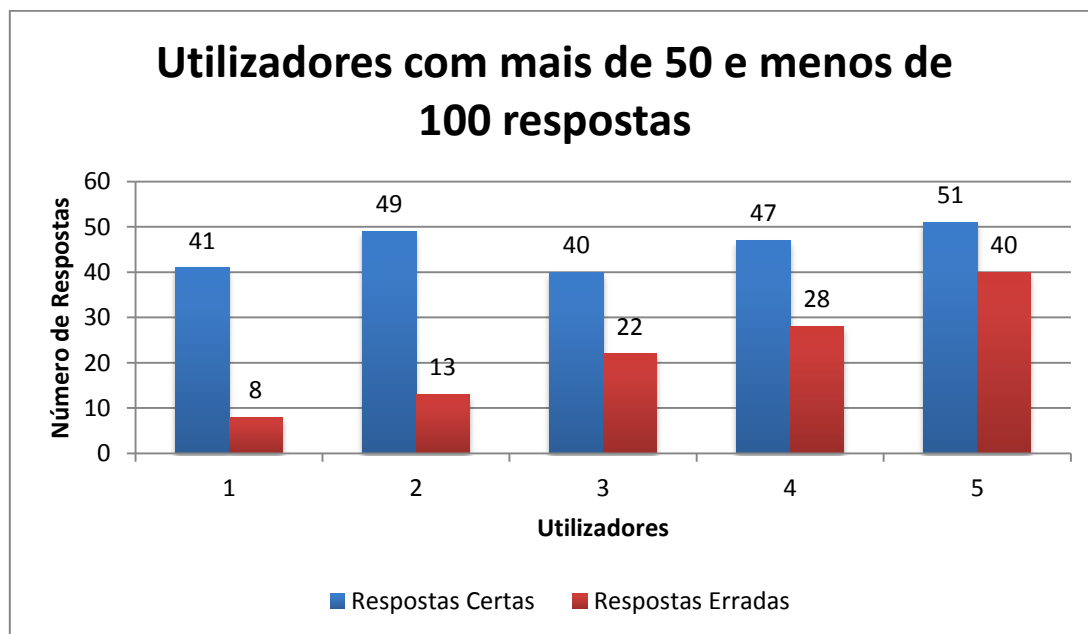


Figura 56 - Utilizadores com mais de 50 e menos de 100 respostas Utilizadores com mais de 50 e menos de 100 respostas

Relativamente a este grupo de utilizadores, pode constatar-se que o número de classificações certas é sempre mais elevado do que o de classificações erradas. Contudo alguns utilizadores apresentam um número de classificações certas e erradas mais próximos que outros.

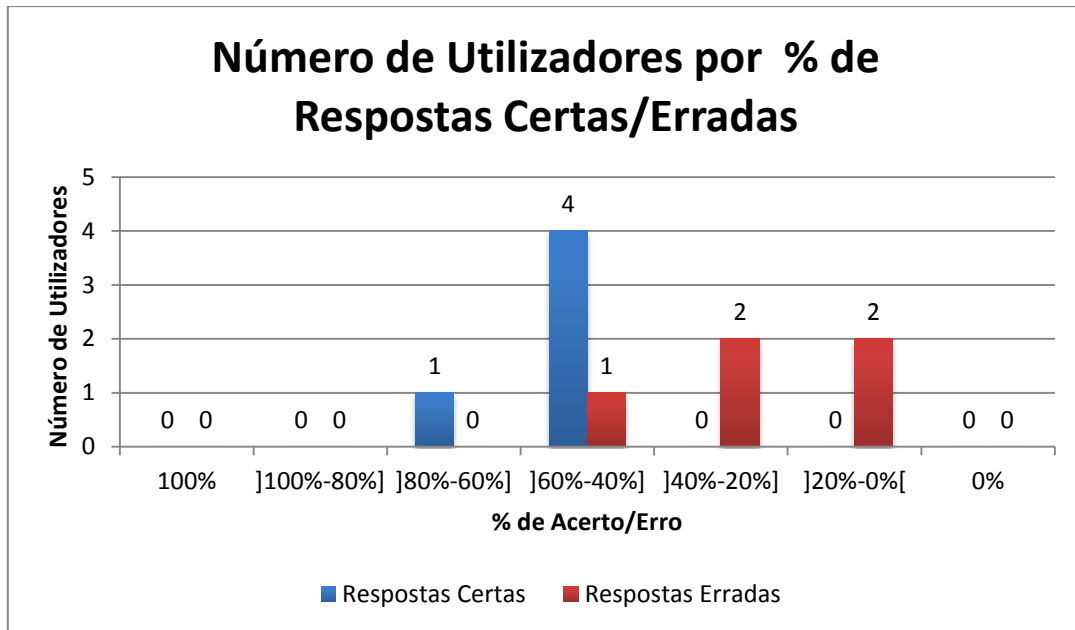


Figura 57 - Número de Utilizadores por % de Respostas Certas/Erradas

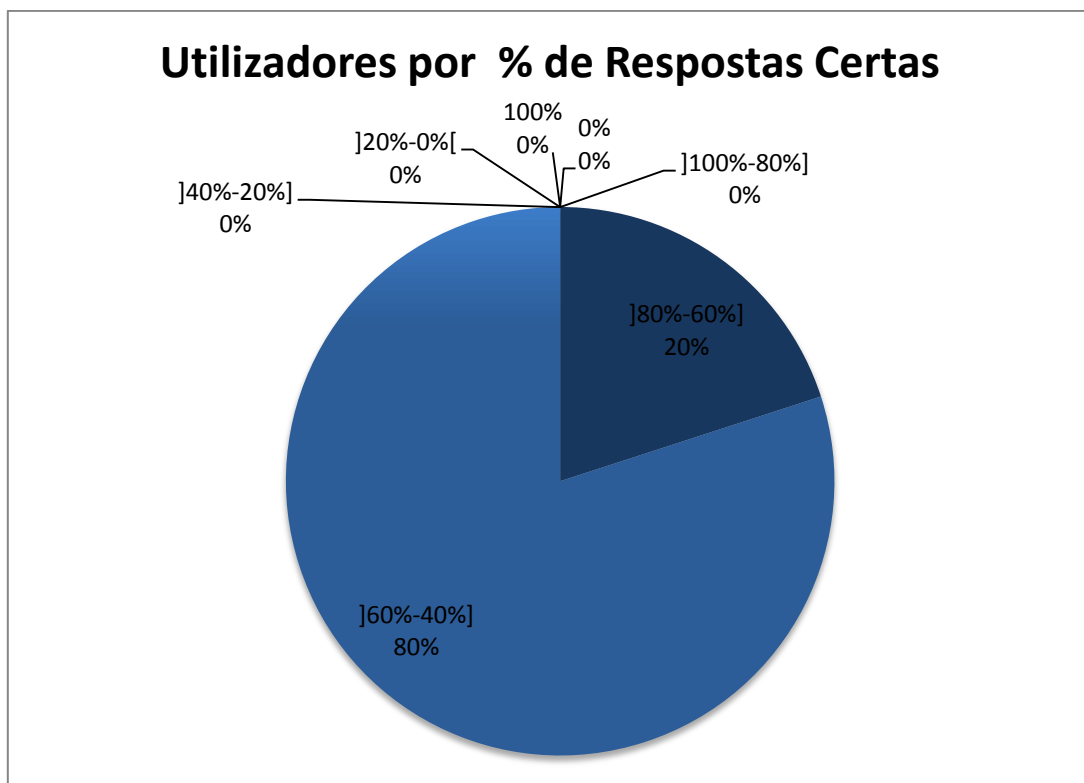


Figura 58 - Utilizadores por % de Respostas Certas

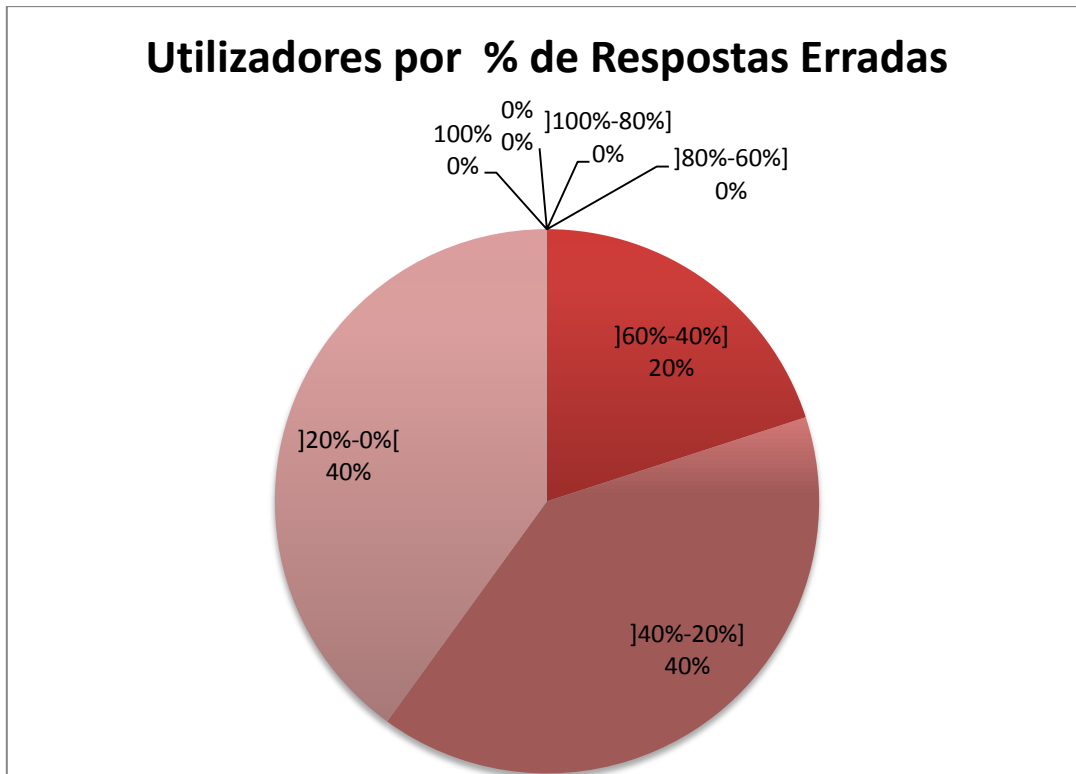


Figura 59 - Utilizadores por % de Respostas Erradas

Nas Figuras 57, 58 e 59, pode verifica-se que as percentagens de acerto apenas se encontram entre os 80% e os 40%. Contudo a maioria dos utilizadores apresentam uma percentagem de assertividade entre os 60% e os 40%. Relativamente às respostas erradas, pode verificar-se que estas apenas apresentam utilizadores com uma percentagem de erro entre os 60% e 0%, contudo a maioria destes apresenta-se entre os 40% e 0%.

#### 4.6.1.4 Utilizadores com Mais de 100 Respostas

Para os utilizadores com mais de 100 respostas obtiveram-se os seguintes resultados expostos na Figura 60.

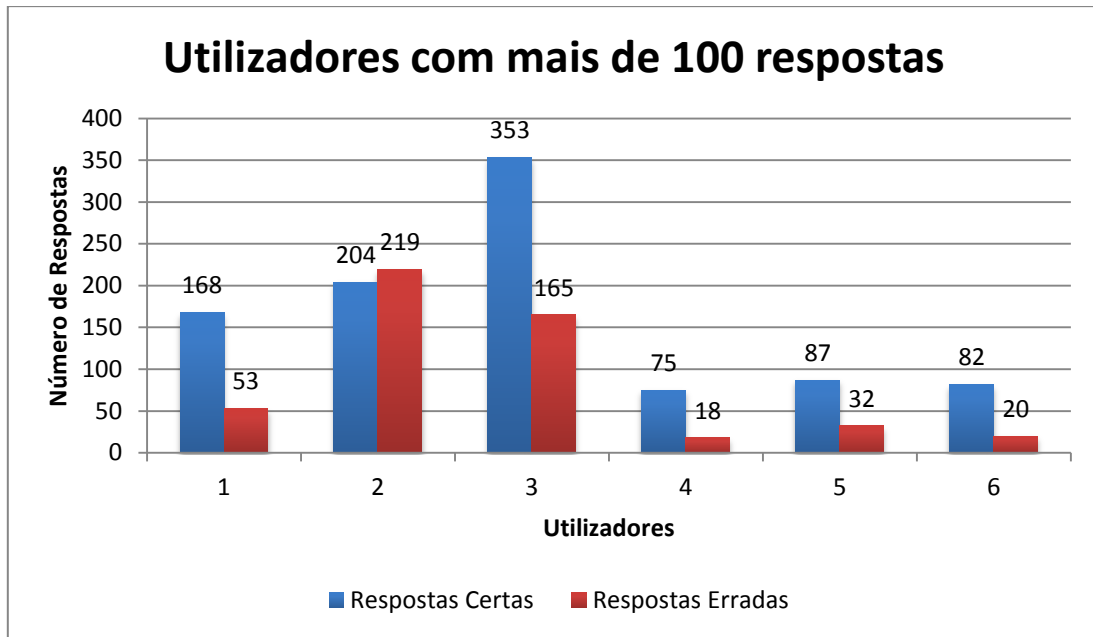


Figura 60 - Utilizadores com mais de 100 respostas

Relativamente a este grupo de utilizadores, pode constatar-se que apesar de poucos, apresentam uma importante parcela das classificações. Existe um utilizador que apresenta mais classificações erradas do que classificações certas. Dos restantes a diferença entre classificações certas e erradas é bastante significativa.

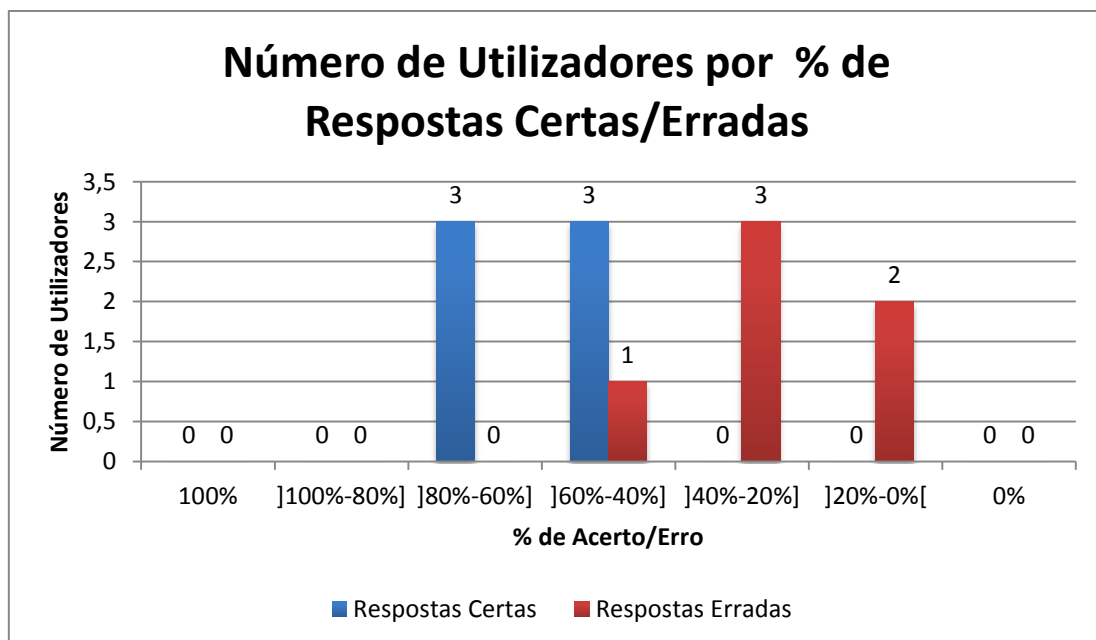


Figura 61 - Número de Utilizadores por % de Respostas Certas/Erradas



Figura 62 - Utilizadores por % de Respostas Certas

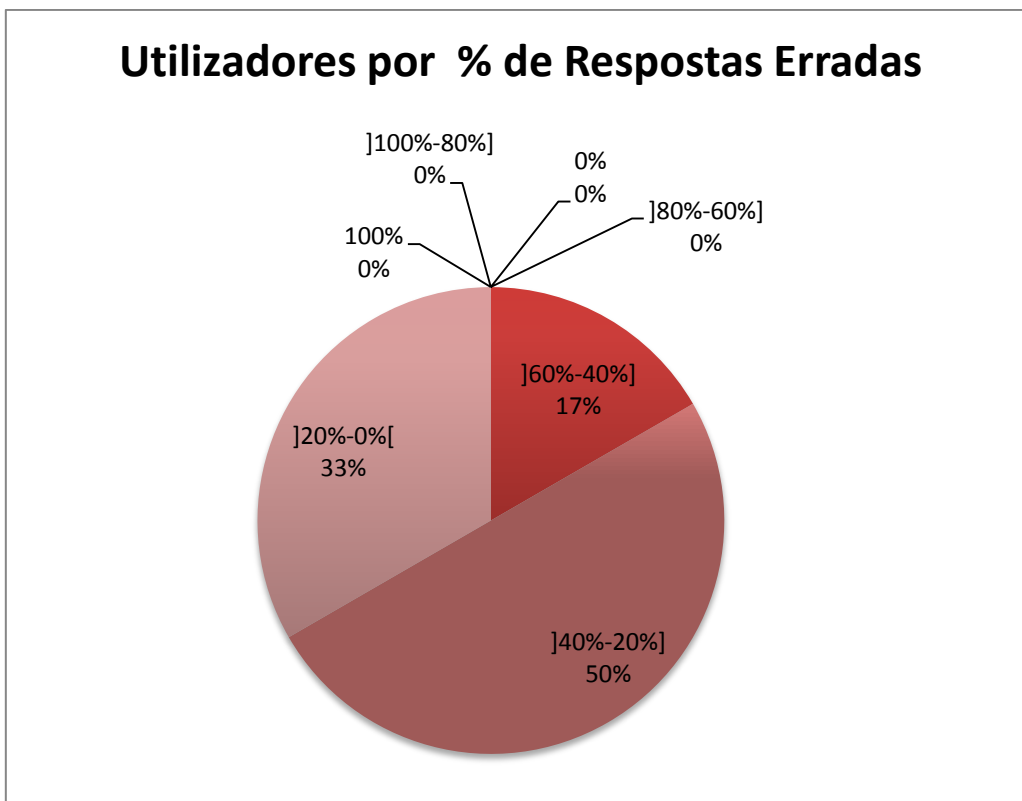


Figura 63 - Utilizadores por % de Respostas Erradas



Nas Figuras 61, 62 e 63 podem analisar-se mais pormenorizadamente as percentagens de acerto e de erro dos utilizadores. Todos os utilizadores possuem uma percentagem de acerto entre os 80% e 40%, metade entre 80% e 60% e a outra metade entre 60% e 40%. Em termos de respostas erradas os utilizadores possuem uma percentagem de erro entre os 60% e 0%, contudo a maioria possui menos de 40%.

#### 4.6.1.5 Dados Gerais

De forma geral obtiveram-se os seguintes resultados representados nas Figuras 64, 65 e 66.

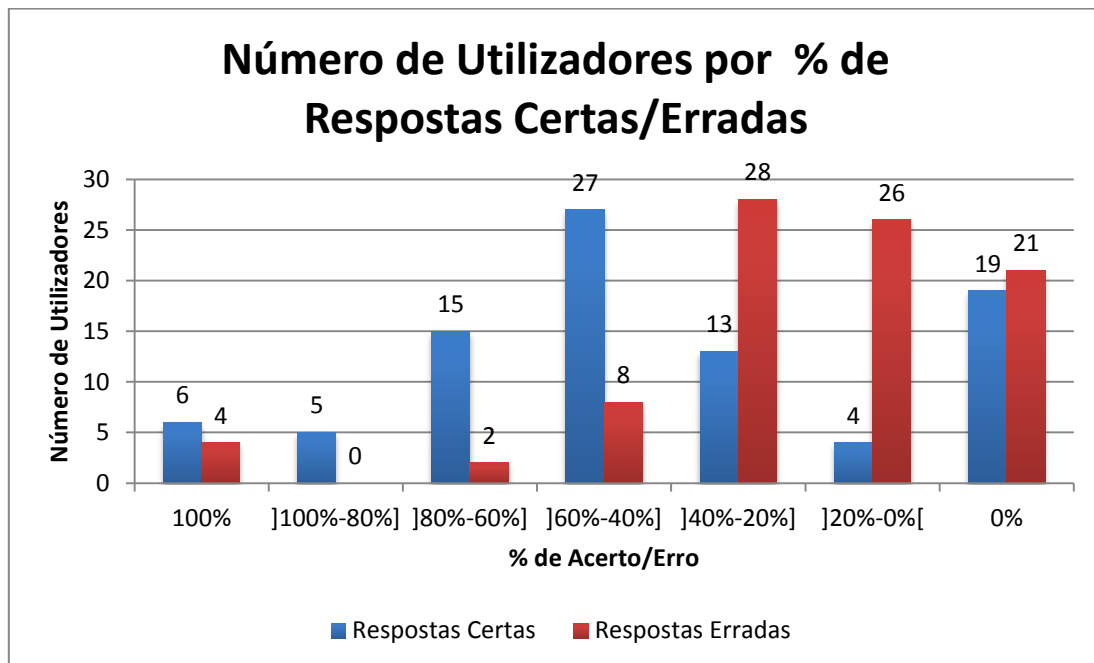


Figura 64 - Número de Utilizadores por % de Respostas Certas/Erradas

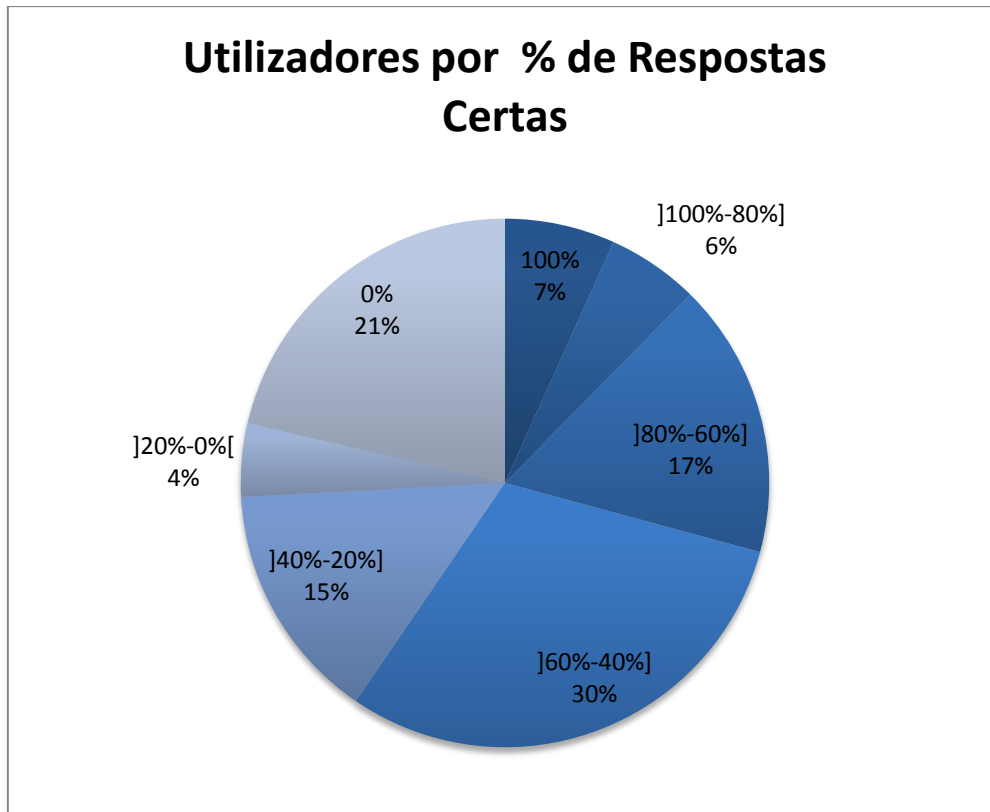


Figura 65 - Utilizadores por % de Respostas Certas

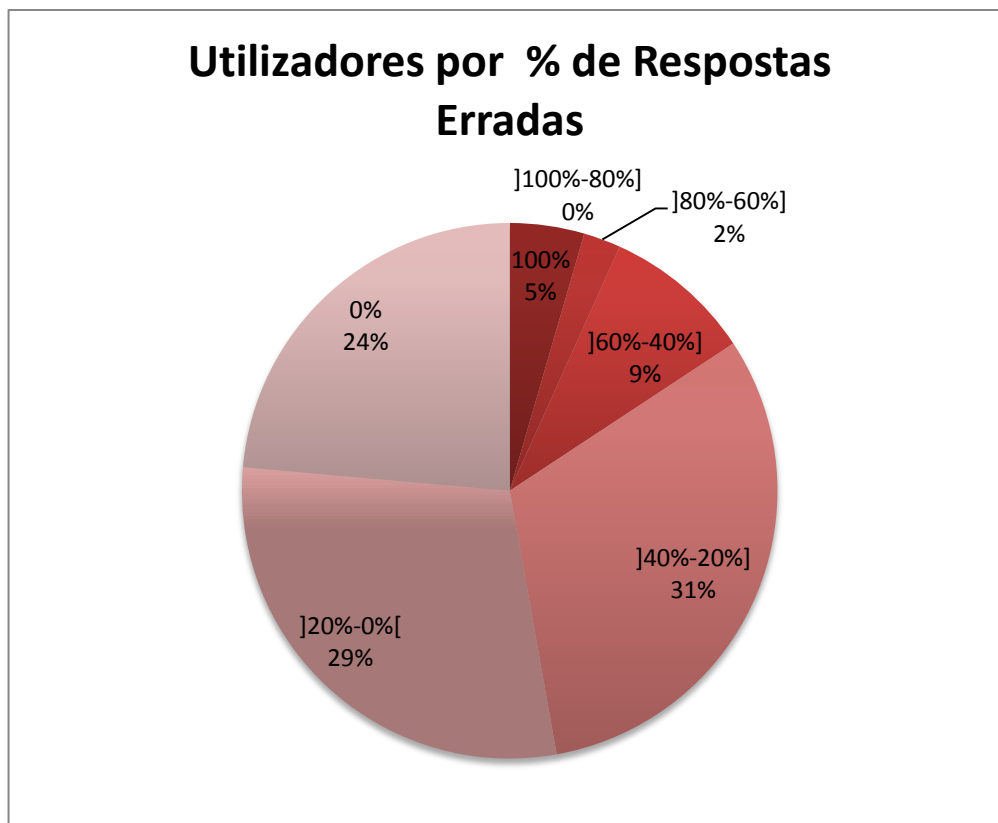


Figura 66 - Utilizadores por % de Respostas Erradas

Relativamente aos gráficos gerais pode verificar-se que na sua grande maioria os utilizadores apresentam uma percentagem de acerto entre os 80% e os 20%, apesar de também existir uma grande quantidade de utilizadores com percentagem de acerto igual a 0%. As percentagens de erro com mais utilizadores são de 40% a 0%.

#### 4.6.2 Análise das Classificações Individuais

Neste subcapítulo são descritos os dados relativos a todas as classificações realizadas pelos utilizadores da aplicação. Das 3.057 Classificações realizadas, 1.547 foram classificadas como Positivas, 1.004 como Negativas, 277 como Neutras e 229 como Díficeis de Classificar. Das 1.547 classificações Positivas, 1.242 foram acertadas e das 1.004 Negativas apenas 489 foram acertadas. Ao todo, as frases positivas são 2.061 e negativas 978 como se pode ver na Tabela 6.

Tabela 6 - Classificações individuais

	Respostas dadas	Acertadas	Numero total de Positivas	Numero total de Negativas
Respostas Positivas	1547	1242	2061	-
Respostas Negativas	1004	489	-	978
Respostas Neutras	277	-	-	-
Respostas Díficeis de Classificar	229	-	-	-
Total de Respostas	3057	-	-	-

Das frases positivas obtiveram-se os seguintes resultados apresentados na Tabela 7 e na Figura 67.

Tabela 7 - Classificações atribuídas às frases Positivas

Frases Positivas				
Acertadas	Erradas			Total
	Neutras	Díficeis de Classificar	Negativas	
1242	182	144	493	2061
60%	9%	7%	24%	100%

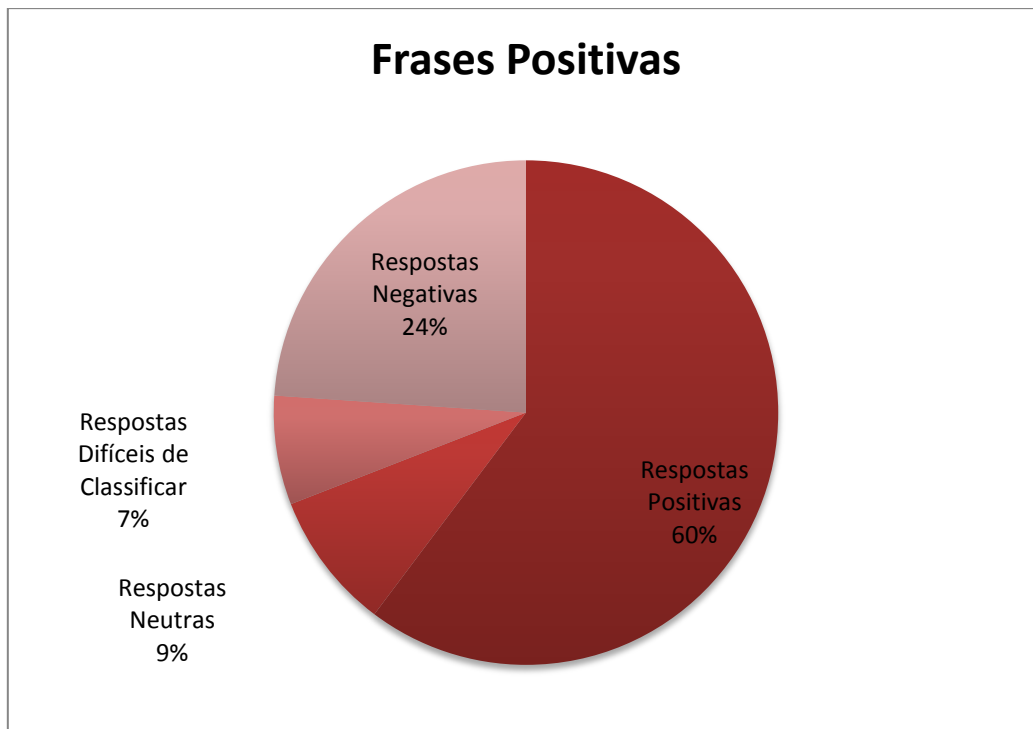


Figura 67 - Classificações atribuídas às Frases Positivas

Pode verificar-se que a maior parte das classificações foram acertadas (60%) e a percentagem de respostas Negativas foi de apenas 24%. 7% das classificações foram consideradas Difíceis de Classificar e 9% Neutras.

Das frases negativas obtiveram-se os resultados representados na Tabela 8 e na Figura 68.

Tabela 8 - Classificações atribuídas às Frases Negativas

Frases Negativas				
Acertadas	Erradas			Total
	Neutras	Difíceis de Classificar	Positivas	
489	94	84	311	978
50%	10%	9%	32%	100%

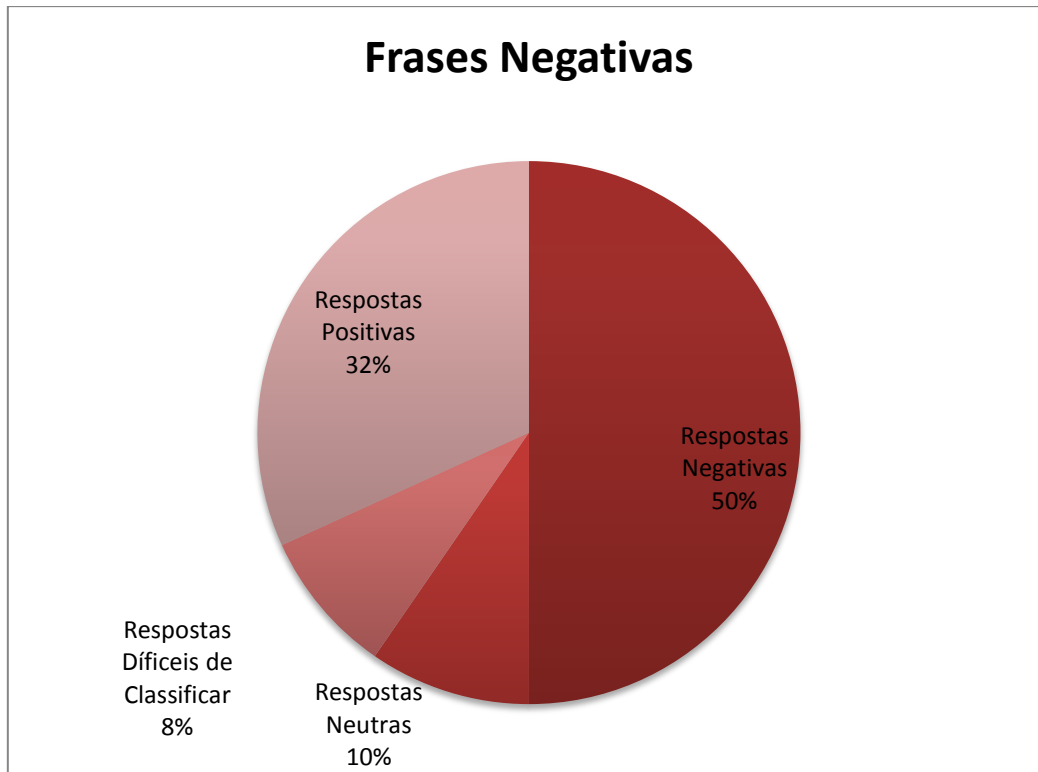


Figura 68 - Classificações atribuídas às Frases Negativas

A percentagem de acerto no caso das frases negativas foi 10% mais baixa que nas frases positivas. Apesar disso esta percentagem é ainda de 50%, o que implica que metade das classificações foi acertada e 32% destas foram classificadas como positivas. As classificações Neutras tiveram uma percentagem de 10% e as Dífceis de Classificar 8%, percentagens semelhantes às obtidas nas frases positivas.

#### 4.6.3 Análise das Classificações Médias

Neste subcapítulo são descritos os dados relativos à média das classificações realizadas a cada frase. Foram analisadas apenas as frases com exatamente 3 classificações de forma a obter um sentimento em maioria. Contudo, para casos onde nenhum sentimento está em maioria, as frases foram classificadas como indefinidas. Das classificações médias analisadas, 514 foram dadas como positivas e 423 destas, estão corretas, 309 foram dadas como negativas, mas apenas 167 foram acertadas. Para além destas, 41 foram dadas como neutras, 27 como dífceis de classificar e 128 foram consideradas indefinidas. De entre as frases

classificadas, 687 são positivas, 326 são negativas e com um total de 590 classificações acertadas, a percentagem de acerto médio é de 58%.

Tabela 9 - Classificações Médias

	Respostas dadas	Acertadas	Numero total de Positivas	Numero total de Negativas
Respostas Positivas	514	423	687	-
Respostas Negativas	309	167	-	326
Respostas Neutras	41	-	-	-
Respostas Difíceis de Classificar	27	-	-	-
Respostas Indefinidas	128	-	-	-
Total de Respostas	1019	590	-	-
% Acerto Médio	58%	-	-	-

Relativamente às frases positivas obtiveram-se os resultados apresentados na Tabela 10 e na Figura 69.

Tabela 10 - Classificação Média das Frases Positivas

Frases Positivas					
Acertadas	Erradas				Total
	Neutras	Difíceis de Classificar	Negativas	Indefinidas	
423	25	19	136	84	687
61%	4%	3%	20%	12%	100%

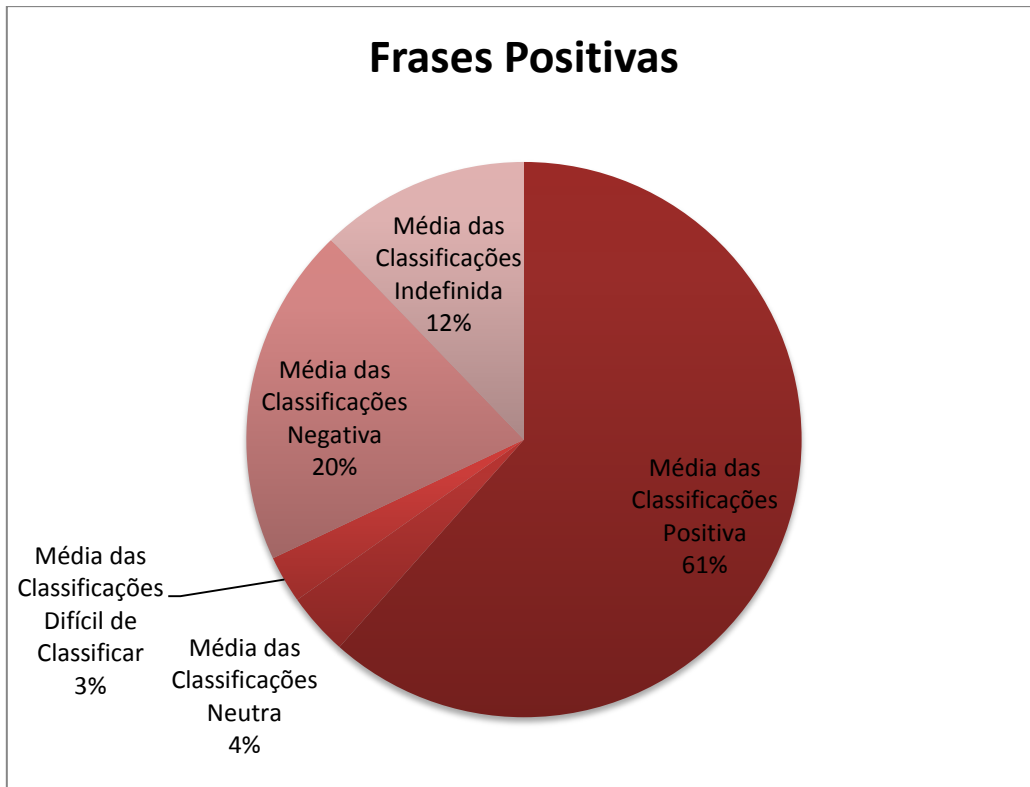


Figura 69 - Classificação Média das Frases Positivas

A maior parte das classificações médias foi considerada positiva, mais precisamente 61%. 20% das classificações médias foram consideradas negativas, 3% difíceis de Classificar e 4% neutras. Para além destes resultados, 12% das classificações médias foram consideradas indefinidas pois não apresentaram consenso quanto ao sentimento exposto.

No que diz respeito às frases negativas obtiveram-se os dados da Tabela 11 e da Figura 70.

Tabela 11 - Classificação Média das Frases Negativas

Frases Negativas					
Acertadas	Erradas				Total
	Neutras	Difíceis de Classificar	Positivas	Indefinidas	
167	16	8	92	43	326
51%	5%	2%	28%	13%	100%

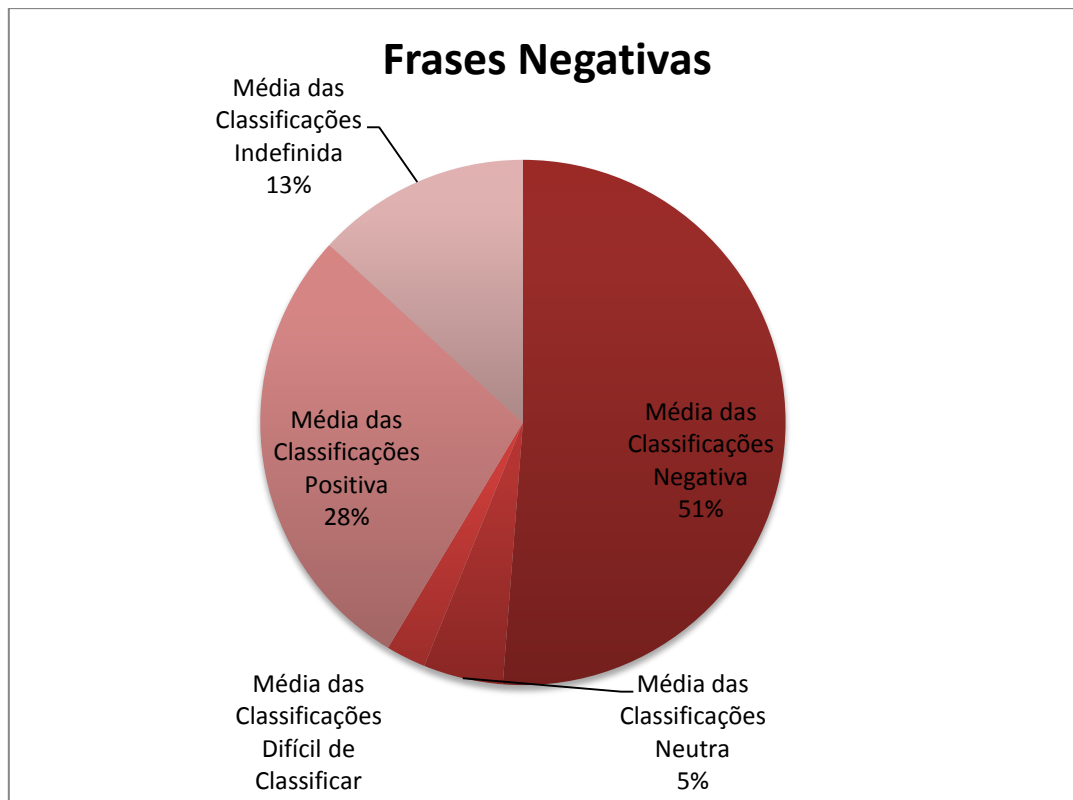


Figura 70 - Classificação Média das Frases Negativas

Das frases negativas analisadas, 51% das classificações médias estão certas, contudo 28% expressam o sentimento oposto (Positivo), 3% são consideradas difíceis de classificar, 5% neutras e 13% indefinidas. Relativamente às frases Positivas, as classificações médias Negativas têm menos percentagem de acerto e mais de erro (classificação positiva quando é negativa).

Nas Figuras 71 e 72 são apresentados os gráficos de comparação dos resultados entre as percentagens obtidas com as classificações individuais e com a média destas.



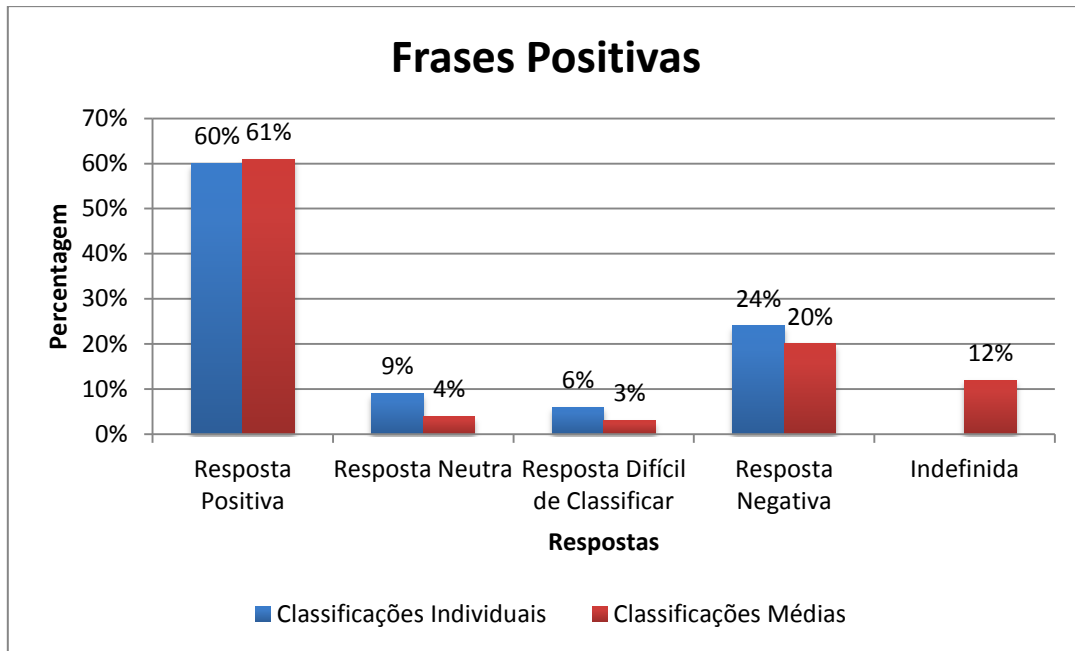


Figura 71 - Comparação das Classificações Individuais e da Média (Frases Positivas)

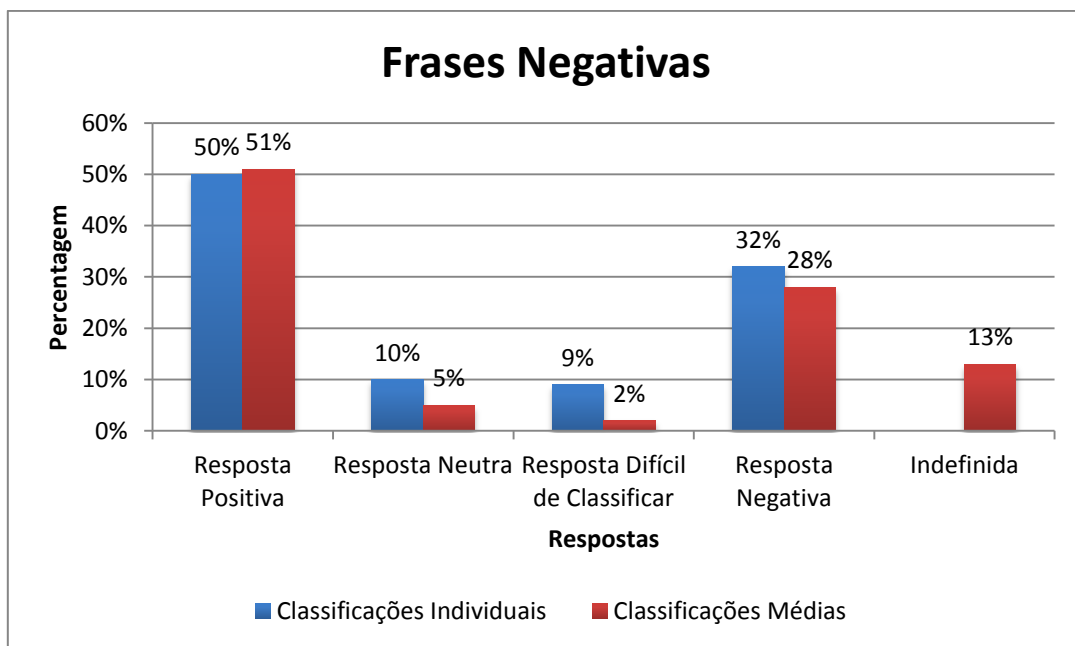


Figura 72 - Comparação das Classificações Individuais e da Média (Frases Negativas)

Das Figuras 71 e 72 pode concluir-se que de uma forma geral as classificações individuais apresentam uma percentagem de acertos mais baixo e de erro mais elevadas que a média. Relativamente à percentagem de acerto, os valores têm uma diferença pequena de 1%, mas em termos de percentagem de erro variam em 4%. A classificação Indefinida é apenas considerada

na Média. Relativamente às classificações neutras e difíceis de classificar, os valores também são mais reduzidos nas classificações médias.

#### 4.6.4 Comparação dos Resultados com Abordagens Automáticas

Nesta secção são descritos e comparados os resultados obtidos a partir das médias das classificações realizadas na aplicação, um algoritmo de Análise de Sentimento e os léxicos FIN e SWN. Na Tabela 11 são expostos os resultados e a percentagem de respostas acertadas através de cada método.

Tabela 12 - Comparação de Resultados da Média das Classificações, do algoritmo e dos léxicos FIN e SWN

	Média	Algoritmo	FIN	SWN
Respostas Certas	589	771	168	591
Respostas Erradas	429	248	851	428
% de respostas certas	58%	76%	16%	58%

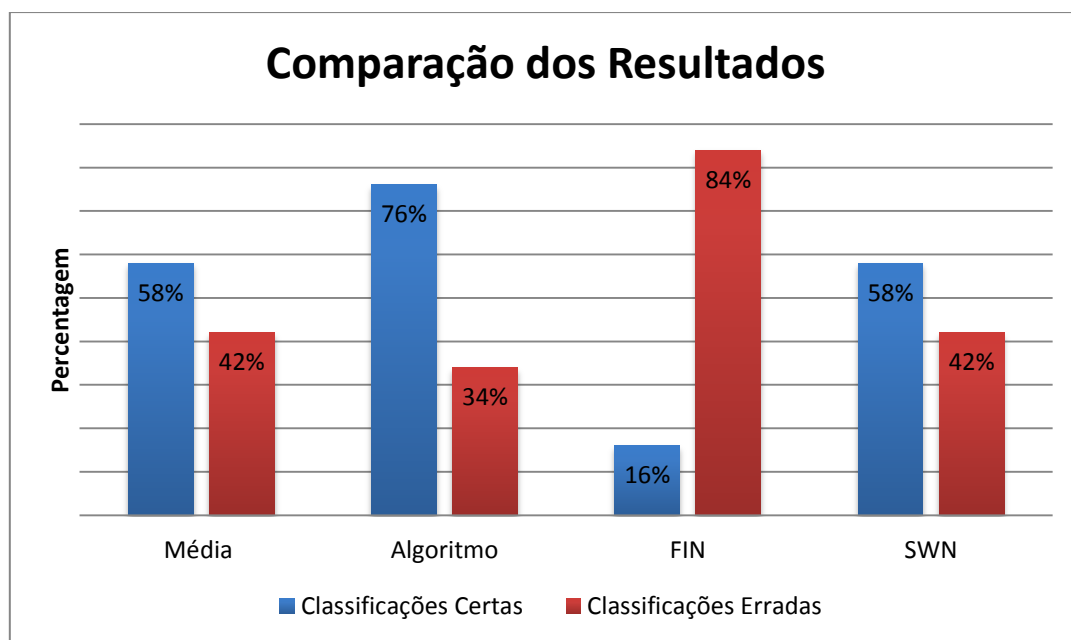


Figura 73 - Comparação de Resultados da Média das Classificações, do algoritmo e dos léxicos FIN e SWN

Através da análise do gráfico da Figura 73 pode verificar-se que o algoritmo de Análise de Sentimento (desenvolvido pelo aluno de doutoramento Nuno Oliveira e que se baseia num novo

léxico especificamente criado a partir de mensagens *StockTwits*) obtém uma maior percentagem de acerto na classificação correta da frase, seguido da média obtida a partir das classificações da aplicação deste projeto assim como do Léxico SWN que apresenta os mesmos valores. Em último lugar surge o Léxico FIN com uma grande percentagem de erro nos dados analisados.

#### 4.6.5 Sumário

Nesta fase do projeto, foi desenvolvida uma aplicação *Web* para verificar a capacidade que a inteligência coletiva possui para realizar análise de sentimentos na área Financeira. As frases para análise foram obtidas a partir da plataforma *StockTwits*, por possuírem um âmbito financeiro e também um sentimento já definido pelo seu autor. Desta forma é possível comparar os resultados obtidos com os corretos.

A aplicação permite aos seus utilizadores classificar uma frase em “Positiva”, “Negativa”, “Neutra” e “Difícil de Classificar”. Para obter dados mais concretos, cada frase é classificada por 3 utilizadores diferentes para que possa ser realizada uma agregação das respostas obtidas (via uma maioria das respostas) e que neste trabalho foi designada por “média” para simplificar a leitura deste documento.

Foram analisados os dados dos utilizadores, em relação ao número classificações realizadas e à sua percentagem de acerto e erro, de modo a verificar a fiabilidade dos mesmos em grupo e individualmente.

Os utilizadores com menos de 10 classificações possuem poucas classificações, porém é possível verificar que são poucos aqueles que possuem mais respostas erradas do que certas. Existem porém alguns utilizadores cujo número de classificações corretas e erradas é igual a 0 pois realizaram as classificações “Neutra” ou “Difícil de Classificar”, que não são consideradas nem certas nem totalmente erradas. Este grupo apresenta uma percentagem de acerto baixa, e uma percentagem de erro mais baixa ainda. Sendo a sua média de percentagem de acerto menor que 40% e a média de percentagem de erro menor que 20%.

Os utilizadores com mais de 10 e menos de 50 classificações acertaram mais do que erraram, tendo apenas dois destes utilizadores fugido ao padrão. De entre estes utilizadores existe um número muito reduzido que apresenta percentagem de acerto menor do que 20%. A

percentagem média de acerto encontra-se entre os 40% e 80% e de erro entre os 20% e 40%. Assim a percentagem de acerto é média-alta e a percentagem de erro é baixa.

O grupo de utilizadores com mais de 50 e menos de 100 respostas, apresenta apenas 5 participantes. Apesar disso todos eles apresentam um número de classificações certas sempre mais elevado do que o de classificações erradas. As suas médias percentuais de acerto encontram-se entre os 40% e 60%, e as de erro entre os 0% e 40%. Assim pode concluir-se que as percentagens médias de acerto são médias e as percentagens de erro são baixas.

Os utilizadores com mais de 100 classificações são apenas 6, contudo este é um grupo importante pois representam uma grande parcela das classificações realizadas. Com a exceção de um utilizador, todos eles apresentam mais classificações certas do que erradas. Assim a percentagem média de acerto encontra-se entre os 40% e 80% e de erro entre os 0% e 40%, sendo a percentagem média de acerto média-alta e de erro baixa.

De uma forma geral, a maior parte dos utilizadores possui uma percentagem média de acerto entre os 20% e 80% e com alguns valores iguais a 0%, e uma percentagem média de erro entre os 0% e 40%. O grupo com melhor média percentual de erro é o dos utilizadores com menos de 10 classificações, pois esta é menor que 20%. Os grupos com melhor média percentual de acerto foram aqueles com mais de 10 e menos de 50 classificações e com mais de 100 respostas pois ambos possuem uma percentagem média entre os 40% e 80%.

Por sua vez os dados são também analisados ao nível das classificações individuais de modo a obter todos os resultados possíveis e também ao nível das classificações médias. A maior parte das classificações individuais positivas foram acertadas (60%) e a percentagem de respostas Negativas foi de apenas 24%. 7% das classificações foram consideradas Difíceis de Classificar e 9% Neutras. A percentagem de acerto no caso das frases negativas foi 10% mais baixa que nas frases positivas, apesar de possuir ainda uma percentagem de 50%, o que implica que metade das classificações foi acertada e 32% destas foram classificadas como positivas. As classificações Neutras tiveram uma percentagem de 10% e as Difíceis de Classificar, 8%, percentagens semelhantes às obtidas nas frases positivas.

Relativamente à média das classificações, há a introdução de uma nova classificação, “indefinida”, quando não há nenhum sentimento em maioria. A percentagem de acerto médio é

de 58%. No caso das frases positivas, o acerto foi de 61%, 20% foram consideradas negativas, 3% difíceis de Classificar, 4% neutras e 12% indefinidas. Das frases negativas analisadas, 51% das classificações médias estão corretas, contudo 28% expressam o sentimento oposto (Positivo), 3% são consideradas difíceis de classificar, 5% neutras e 13% indefinidas. Relativamente às frases Positivas, as classificações médias Negativas têm menos percentagem de acerto e mais de erro (classificação positiva quando é negativa).

De uma forma geral as classificações individuais apresentam uma percentagem de acertos mais baixo e de erro mais elevadas que a média. Relativamente à percentagem de acerto os valores têm uma diferença pequena de 1%, mas em termos de percentagem de erro variam em 4%. Relativamente às classificações neutras e difíceis de classificar, os valores também são mais reduzidos nas classificações médias.

Por fim, foram comparados os dados obtidos a partir da aplicação com a classificação dos mesmos dados por um algoritmo de Análise de Sentimento e dois Léxicos *FIN* e *SWN*.

O algoritmo obtém uma maior percentagem de acerto na classificação correta da frase, seguido da média obtida a partir das classificações da aplicação deste projeto assim como do Léxico *SWN* que apresenta os mesmos valores. Em último lugar surge o Léxico *FIN* com uma grande percentagem de erro nos dados analisados.



## 5 Conclusões

### 5.1 Síntese

Com a expansão da *Internet* e o surgimento da *Web 2.0* surgiram inúmeras possibilidades de cooperação. Em particular, elevou-se o potencial da Inteligência Coletiva, que utiliza uma colaboração entre um conjunto de indivíduos para alcançar um dado objetivo. Por outro lado, a massificação da *Internet* a nível mundial aumentou de modo exponencial a quantidade de texto disponível em formato eletrónico, contribuindo assim para um desenvolvimento da área de *text mining*, que utiliza métodos computacionais para extrair conhecimento útil a partir de textos. A Análise de Sentimento é uma área particular do *text mining* e visa a extração automática de opiniões (e.g., gosto/não gosto) sobre uma frase ou texto.

Este projeto foca-se na avaliação do desempenho da inteligência coletiva quando orientada para a Avaliação de Sentimentos de mensagens de âmbito financeiro do *microblog StockTwits* (stocktwits.com). Para o efeito, desenvolveu-se um protótipo de uma aplicação, designada de *Finance\$entiment*, e que foi disponibilizada num servidor da Universidade do Minho (<https://financesentimentapp.dsi.uminho.pt/>). Esta aplicação foi divulgada em língua portuguesa a docentes e alunos do ensino superior, bem como redes sociais (e.g. *Linked-in*), tendo sido utilizada durante um período de dois meses. Cada mensagem foi classificada como “Positiva”, “Negativa”, “Neutra” ou “Difícil de Classificar”. No final do período de disponibilização da aplicação foram obtidas 3057 classificações, o equivalente a 1019 mensagens classificadas por três utilizadores distintos. Como método de agregação (das três respostas), optou-se pelo valor de sentimento que estivesse em maioria (i.e., “Positiva”, “Negativa”, “Neutra” ou “Difícil de Classificar”). Foi ainda definida a categoria “Indefinida”, que é utilizada quando não existe uma maioria de respostas comuns.

No que diz respeito às classificações individuais obtidas, obteve-se uma percentagem de acerto de 60% para as mensagens positivas e 50% para as mensagens negativas. Quando se agregam as três respostas, os valores de acerto rondam os 61% para as mensagens positivas e 51% para os *tweets* negativos. Em termos globais, a Inteligência Coletiva obteve uma taxa de acerto de 58%.

Os resultados obtidos foram comparados com três métodos automáticos relativos ao léxico SWN, FIN e um algoritmo de *text mining* especificamente proposto para mensagens do serviço *Stocktwits*. O algoritmo de *text mining* obteve o melhor desempenho (taxa de acerto global de 76%), seguindo-se a Inteligência Coletiva e léxico SWN (ambos com 58%), sendo que o léxico FIN obteve o pior desempenho (16% de acerto).

## 5.2 Discussão

Este trabalho é considerado inovador, uma vez que não se conhecem outros trabalhos que tenham aplicado o conceito da Inteligência Coletiva para a classificação da Análise de Sentimento de *tweets* da área financeira. Para aplicar e avaliar tal conceito, foi necessária a execução de diversas etapas, incluindo o desenvolvimento de uma aplicação Web, a divulgação e monitorização de tal aplicação, e a posterior análise dos resultados obtidos.

Durante a condução dos trabalhos existiram limitações que afetaram a fase de divulgação e monitorização da aplicação. Primeiro, porque os dados analisados não são totalmente públicos, sendo que o acordo de cedência dos dados que foi estabelecido com a empresa *StockTwits* dizia sobretudo respeito a sua análise para fins de investigação. Assim sendo, a divulgação da aplicação *Finance\$entiment* ocorreu num ambiente mais reservado, com divulgação em língua portuguesa e num contexto nacional. Segundo, porque dados os limites temporais para a execução deste trabalho, só foi possível utilizar a aplicação durante um período curto de dois meses. Apesar destas limitações, conseguiu-se reunir um conjunto interessante de respostas, com um total de 3057 classificações e 1019 mensagens classificadas por três indivíduos distintos.

Os resultados obtidos são considerados interessantes e úteis dentro de um projeto mais alargado que está relacionado com a previsão de indicadores financeiros a partir da análise automática de mensagens de *microblogs* sobre mercados financeiros. É certo que o resultado obtido pela Inteligência Coletiva não tem um desempenho elevado, equiparando-se ao método automático que utiliza o léxico SWN e tendo obtido um acerto de 61% para as mensagens positivas, 51% para as mensagens negativas, resultando num acerto global de 58%. Tal resultado pode dever-se a múltiplos fatores, nomeadamente:



- os *tweets* estão escritos em língua inglesa, sendo que a divulgação da aplicação foi efectuada em língua portuguesa e ambiente nacional;
- os *tweets* têm um elevado número de termos específicos do serviço de *microblogging* e da área financeira, sendo que procurou-se uma divulgação da aplicação por alunos e docentes do ensino superior que estivessem mais ligados à área da Economia e Finanças, mas não é garantido que todos utilizadores tivessem esse perfil;
- muitos dos *tweets* têm um contexto associado (e.g., saber o que aconteceu ontem ou na semana anterior, imagem gráfica que compara 2 ou mais ações da bolsa) que não é possível discernir pela análise estrita do texto escrito no *tweet*.

Sobretudo, o resultado evidencia que a Análise de Sentimento em serviços de *microblogging* e em particular de *tweets* sobre ações financeiras é uma tarefa não trivial. Como tal, valorizam o método automático *de texto mining* que foi desenvolvido pelo aluno de doutoramento Nuno Oliveira, e que tendo somente acesso aos *tweets* (e não à informação de contexto), conseguiu um desempenho notável de 76%. Tal resultado é relevante pois existem outros *microblogs* que disponibilizam de modo público os seus *tweets* sobre a área financeira (e.g., *Twitter*) mas que não têm uma classificação (em termos de “bullish” ou “bearish”) dos mesmos pelo próprio autor. Quando à fraca taxa de acerto do léxico FIN, convém referir que apesar deste ser vocacionado para a área financeira, foi criado utilizando documentos financeiros, não estando por isso adaptado ao *microblogging* (que tende a utilizar poucas palavras e caracteres).

Importa realçar ainda que a empresa *StockTwits* já manifestou interesse em conduzir um projeto similar de uso de Inteligência Coletiva, sendo que o mesmo poderá ser disponibilizado no próprio portal da empresa.

### 5.3 Trabalho Futuro

O trabalho desenvolvido nesta dissertação é considerado como o primeiro da área e por isso tendo um forte teor exploratório. Assim, existem diversas perspetivas de melhoria ao trabalho desenvolvido, destacando-se aqui três possibilidades interessantes de trabalho futuro:

- Melhoria da aplicação, com adaptação da mesma às necessidades da empresa *StockTwits*, de modo a que (tal como foi já referido) se possa testar a Inteligência Coletiva diretamente no portal da *StockTwits*. Tal permitiria obter um maior e melhor *feedback*, ou seja, obter-se-ia um maior número de classificações vindas de utilizadores especializados.
- Transformação da aplicação num serviço de verificação para aceder a determinado conteúdo (sobretudo de serviços financeiros), de modo similar ao que é executado pelo projeto *reCAPTCHA*, e que permitiria uma maior adesão e um maior número de classificações. Tal como o *reCAPTCHA*, o sistema de verificação poderia exigir N pares de questões, sendo que cada par implicaria uma questão cuja resposta é conhecida e outra não, ou seja, um *tweet* em que o sentimento já esteja classificado pelo autor e outro para classificar.
- Adaptação da aplicação desenvolvida para classificar o sentimento de *tweets* de outras áreas que não a financeira (e.g., turismo, com a análise de opiniões sobre visitas a cidades ou países).

## Referências Bibliográficas

Abrich, R., Berbenetz, V., & Thorpe, M. (2011). Distinguishing between Humans and Robots on the Web. 1-7.

Aggarwal, C. C., & Zhai, C. (2012). *Mining Text Data*, USA, Springer.

Alag, S. (2009). *Collective Intelligence in Action*. United States of America: Manning Publications Co.

*Alias-i*. (04 de Abril de 2013). Obtido de LingPipe 4.1.0: <http://www.alias-i.com/lingpipe>

Antweiler, W., & Frank, M. Z. (2004). Is All That Talk Just Noise? The Information Content of Internet Stock Message Boards. *The Journal of Finance*, 1259-1294.

Atal, K., Arora, A., & Sachan, D. S. (2013). reCAPTCHA assisted OCR for Devanagiri Texts. *Proceedings of the 1st Indian Workshop on Machine*.

Baccianella, S., Esuli, A., & Sebastiani, F. (2010). SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining., (pp. 2200-2204). Italy.

Barbosa, L., & Feng, J. (2010). Robust Sentiment Detection on Twitter from Biased and Noisy Data. *Coling 2010: Poster Volume*, (pp. 36-44). Beijing.

Brabham, D. C. (2008). Crowdsourcing as a Model for Problem Solving. *Convergence: The International Journal of Research into New Media Technologies* (pp. 75-90). London, Los Angeles, New Delhi and Singapore: Sage Publications.

Beneti, A., Hammoumi, W., Hielscher, E., Müller, M., & Persons, D. (22, Março, 2006). *Automatic Generation of Fine-Grained Named Entity Classifications*.

Boiy, E., Hens, P., Deschacht, K., & Moens, M.-F. (2007). Automatic Sentiment Analysis in On-line Text. *In Proceedings of ELPUB2007 Conference on Electronic Publishing*, pp. 349-360, Vienna, Austria.

Bonabeau, E. (2009). Decisions 2.0: The Power of Collective Intelligence. *MIT SLOAN MANAGEMENT REVIEW*, pp. 45-53.

Bravo-Marquez, F., Mendoza, M., & Poblete, B. (11 de August de 2013). Combining Strengths, Emotions and Polarities for Boosting Twitter Sentiment Analysis. *WISDOM*.

- Carenini, G., Ng, R. T., & Zwart, E. (2005). Extracting Knowledge from Evaluative Text.
- Carpenter, B. (2006). Character Language Models for Chinese Word Segmentation and Named Entity Recognition. p. 4.
- Carpenter, B. (2007). LingPipe for 99.99% Recall of Gene Mentions. Valencia, Spain.
- Chan, K., Omokore, J., & Miller, R. K. (2009). *Practical CakePHP Projects*. United States of America: Apress.
- Cunningham, H., Tablan, V., Roberts, A., & Bontcheva, K. (Fevereiro, 2013). Getting More Out of Biomedical Documents with GATE's Full Lifecycle Open Source Text Analytics. *PLOS Computational Biology*, 1-16.
- Das, S. R., & Chen, M. Y. (9 de Setembro de 2007). Yahoo! for Amazon: Sentiment Extraction from Small Talk on the Web. *Management Science*, pp. 1375-13388.
- Davidov, D., Tsur, O., & Rappoport, A. (2010). Enhanced Sentiment Learning Using Twitter Hashtags and Smileys. *Proceedings of Coling 2010*.
- de Albornoz, J. C., Plaza, L., & Gervás, P. (2012). SentiSense: An easily scalable concept-based affective lexicon for sentiment analysis. pp. 3562-3567.
- Ding, X., Liu, B., & Yu, P. S. (2008). A Holistic Lexicon-Based Approach to Opinion Mining. *Proceedings of the Conference on Web Search and Web Data Mining*.
- Doan, A., Franklin, M. J., Kossmann, D., & Kraska, T. (2011). Crowdsourcing Applications and Platforms: A Data Management Perspective. *VLDB 2011*.
- Doan, a., RamaKRishnan, R., & haLeVy, a. y. (2011). Crowdsourcing systems on the World-Wide Web., (pp. 86-96).
- Fader, A., Soderland, S., & Etzioni, O. (2011). *Identifying Relations for Open Information Extraction*. Seattle.
- Feldman, R. (Abril de 2013). Techniques and Applications for Sentiment Analysis. *communications of the ACM*, pp. 82-89.

GAO, H., WANG, W., FAN, Y., QI, J., & LIU, X. (2014). The Robustness of “Connecting Characters Together” CAPTCHAs. *JOURNAL OF INFORMATION SCIENCE AND ENGINEERING*, 347-369.

Glenn, J. C. (5 de October de 2013). Collective intelligence systems and an application by The Millennium Project for the Egyptian Academy of Scientific Research and Technology. *Elsevier*, p. 8.

Godbole, N., Srinivasaiah, M., & Skiena, S. (2007). Large-Scale Sentiment Analysis for News and Blogs. *ICWSM'2007*. Boulder , Colorado, USA.

Gonçalves, R. B. (2012). Utilizacao da Ferramenta RapidMiner no Processo de Analise de Sentimentos. *ENCONTRO DE COMPUTAÇÃO E INFORMÁTICA DO TOCANTINS*. Palmas.

González-Ibáñez, R., Muresan, S., & Wacholder, N. (2011). Identifying Sarcasm in Twitter: A Closer Look. *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics:shortpapers*, (pp. 581–586). Portland, Oregon: Association for Computational Linguistics.

Günther, H. (Maio-Agosto de 2006). Pesquisa Qualitativa versus Quantitativa Esta é a questão? *Psicologia: Teoria e Pesquisa*, pp. 201-210.

Hagenau, M., Liebmann, M., Hedwing, M., & Neumann, D. (2012). Automated news reading: Stock Price Prediction based on Financial News Using Context-Specific Features. *International Conference on System Science* (pp. 1040-1049). Hawaii: IEEE.

Herbrich, R., & Graepel, T. (2010). *HANDBOOK OF NATURAL LANGUAGE PROCESSING*. Cambridge, UK: Chapman & Hall/CRC.

Herzog, C., Luger, M., & Herzog, M. (2007). Combining Social and Semantic Metadata for Search in a Document Repository. *International Workshop at the 4th European Semantic Web Conference*, (pp. 14-21). Innsbruck, Austria.

Howe, J. ( June de 2006). The Rise of Crowdsourcing. *Wired Magazine*.

Informatics Research and Development Unit, Public Health Informatics & Technology Program Office, Office of Surveillance, Epidemiology and Laboratory Services, US Centers for Disease Control and Prevention. (2010). *Open Souce Data Mining Software Evaluation*.

- Kerns, G. (2010). Introduction to Probability and Statistics Using R (Vol. First Edition).
- Kittur, A., & Kraut, R. E. (February de 2010). Beyond Wikipedia: Coordination and Conflict in On-line Production Groups. *ACM*, pp. 6–10.
- Kittur, A., Smus, B., & Kraut, R. E. (1 de February de 2011). CrowdForge: Crowdsourcing Complex Work. *IIS*, p. 13.
- Leimeister, J. M. (2010). Collective Intelligence. *Business & Information Systems Engineering*, pp. 245-249.
- Lesser, E., Ransom, D., Shah, R., & Pulver, B. (2012). *Collective Intelligence Capitalizing on the crowd*. United States of America.
- Liu, B. (2012). *Sentiment Analysis and Opinion Mining*. Graeme Hirst, Series Editor.
- Maynard, D., Bontcheva, K., & Rout, D. (s.d.). Challenges in developing opinion mining tools for social media. UK.
- McMillen, C., & Veloso, M. (2008). Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence. *Unknown Rewards in Finite-Horizon Domains*, 963-968.
- Mejova, Y. (2009). Sentiment Analysis: An overview. Iowa.
- Michalsky, S., Mamani, E. Z., & Gerosa, M. A. (2010). A Inteligência Coletiva na Web: Uma Análise de Domínio para o Jornalismo On-line. p. 4.
- Mizumoto, K., Yanagimoto, H., & Yoshioka, M. (2012). Sentiment Analysis of Stock Market News with Semi-supervised Learning. *11th International Conference on Computer and Information Science* (pp. 325-328). Japão: IEEE/ACIS.
- Mohammad, S. M., Kiritchenko, S., & Zhu, X. (2013). NRC-Canada: Building the State-of-the-Art in Sentiment Analysis of Tweets. *Proceedings of the seventh international workshop on Semantic Evaluation Exercises*. Atlanta, Georgia, USA.
- Moreo, A., Romero, M., Castro, J., & Zurita, J. (2012). Lexicon-based Comments-oriented News Sentiment Analyser system. *Elsevier, 39*(Expert Systems with Applications), 9166-9180.
- Muenchen, R. A. (2010). *R for SAS and SPSS users*.
- Muenchen, R. A. (2013). The Popularity of Data Analysis Software.

Nasukawa, T., & Yi, J. (2003). Sentiment Analysis: Capturing Favorability Using Natural Language Processing.

O'Reilly, T., & Battelle, J. (2009). *Web Squared: Web 2.0 Five Years On*. O'Reilly.

Oliveira, N., Cortez, P., & Areal, N. (07 de July de 2014). Automatic Creation of Stock Market Lexicons for Sentiment Analysis Using StockTwits Data. *IDEAS*, pp. 115-123.

Osimo, D., & Mureddu, F. (2011). Research Challenge on Opinion Mining and Sentiment Analysis.

Palanisamy, P., Yadav, V., & Elchuri, H. (2013). Serendio: Simple and Practical lexicon based approach to Sentiment Analysis. India.

Paltoglou, G., Gobron, S., Skowron, M., Thelwall, M., & Thalmann, D. (2010). Sentiment analysis of informal textual communication in cyberspace. *ENGAGE 2010*. Switzerland.

Pang, B., & Lee, L. (2008). *Opinion Mining and Sentiment Analysis* (Vol. 2). Foundations and Trends.

Pawar, S. E., & Bauskar, M. M. (2013). CAPTCHA: A SECURITY MEASURE AGAINST SPAM ATTACKS. *IJRET: International Journal of Research in Engineering and Technology*, 854-857.

Payal, N., Chaudhary, N., & Astya, P. N. (2012). JigCAPTCHA: An Advanced Image-Based CAPTCHA Integrated with Jigsaw Piece Puzzle using AJAX. *International Journal of Soft Computing and Engineering*, 180-185.

Prekopcsák, Z., Makrai, G., Henk, T., & Gáspár-papanek, C. (2011). Radoop: Analysing Big Data with RapidMiner and Hadoop.

Sano, S., Otsuka, T., & Okuno, H. G. (2014). HMM-based CAPTCHA Breaker for Overlapped Symbols. *Information Processing Society of Japan*, 585-586.

Saxena, V. (2013). A Study on User Friendly Approach : CAPTCHA. *International Journal of Engineering & Management Technology*, 1-8.

Schlaikjer, A. (2007). A Dual-Use Speech CAPTCHA: Aiding Visually Impaired Web Users while Providing Transcriptions of Audio Streams. Pittsburgh, PA.

- Segaran, T. (2007). *Programming Collective Intelligence*. United States of America: O'Reilly.
- Shanker, D., Gupta, P., & Jaiswal, A. (2013). Hybrid Collage CAPTCHA. *INTERNATIONAL JOURNAL OF SCIENTIFIC & ENGINEERING RESEARCH*.
- Sharma, P., Tyagi, N., & Singhal, D. (2013 ). CAPTCHAs: VULNERABILITY TO ATTACKS. *International Journal of Emerging Trends & Technology in Computer Science*, 73-78.
- Shi, G., Ying, Y., & Yin, Y. (2013). The reCaptcha Helper: A Machine Learning Study.
- Sutherland, C. (2012). Usability and Security of Text-based CAPTCHAs. *UMM CSci Senior Seminar Conference*. USA.
- Tablan, V., Roberts, I., Cunningham, H., & Bontcheva, K. (2011). *GATECloud.net: a Platform for Large-Scale, Open-Source Text Processing on the Cloud*. Sheffield, United Kingdom: The Royal Society.
- Taboada, M., Brooke, J., Tofiloski, M., Voll, K., & Stede, M. (2011). Lexicon-Based Methods for Sentiment Analysis. Association for Computational Linguistics.
- Turney, P. D. (2002). Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews. *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, (pp. 417-424). Philadelphia.
- Vaishakh, B. N., & Harish, G. (2011). CAPTCHAS: SURVEY OF EXISTING TECHNIQUES AND A NEW APPROACH. *National Conference on Recent Trends in Computer Technology Technology*, 70–73.
- Vargas, J. (2001). *Sociologia*. Porto: Porto Editora.
- Wang, C.-J., Tsai, M.-F., Liu, T., & Chang, C.-T. (2013). Financial Sentiment Analysis for Risk Prediction. *International Joint Conference on Natural Language Processing*, (pp. 802–808). Nagoya, Japan.
- Weerkamp, W., & de Rijke, M. (2012). Activity Prediction: A Twitter-based Exploration. Portland.
- Wiebet, J. M., Bruce, R. F., & O'Hara, T. P. (1999). Development and Use of a Gold-Standard Data Set for Subjectivity Classifications. *Proceedings of the Association for Computational Linguistics*, (pp. 246-253).



Wilson, T., Wiebe, J., & Hoffmann, P. (2008). Recognizing Contextual Polarity: An Exploration of Features for Phrase-Level Sentiment Analysis. *Association for Computational Linguistics*, 400-433.

Xu, F. (2012). Data Mining in Social for Stock Market Prediction. Halifax, Nova Scotia.

Yu, J., Zha, Z.-J., MengWang, & Chua, T.-S. (2011). Aspect Ranking: Identifying Important Product Aspects from On-line Consumer Reviews. *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics* (pp. 1496–1505). Portland, Oregon: Association for Computational Linguistics.

Zhai, Z., Liu, B., Xu, H., & Jia, P. (2010). Grouping Product Features Using Semi-Supervised Learning with Soft-Constraints. *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, (pp. 1272–1280). Beijing.

Zhang, L., Ghosh, R., Dekhil, M., Hsu, M., & Liu, B. (2011). *Combining Lexicon-based and Learning-based Methods for Twitter Sentiment Analysis*. Hewlett-Packard Development Company.

Zhao, J., Dong, L., Wu, J., & Xu, K. (12 de Dezembro de 2012). MoodLens: An Emoticon-Based Sentiment Analysis System of Chinese Tweets. pp. Beijing, China.





Utilizador 32	0	0	0	0	0	0	0	0	0	0	1	1	1
Utilizador 33	0	0	0	0	0	0	0	0	0	1	0	1	1
Utilizador 34	1	0	1	0	0	0	0	0	0	0	0	0	1
Utilizador 35	2	1	3	1	0	1	0	1	1	1	0	1	6
Utilizador 36	1	0	1	0	0	0	0	0	0	0	0	0	1
Utilizador 37	7	4	11	6	4	10	1	3	4	0	0	0	25
Utilizador 38	0	0	0	0	0	0	0	0	0	0	1	1	1
Utilizador 39	2	0	2	0	2	2	1	2	3	2	1	3	10
Utilizador 40	1	0	1	1	0	1	0	0	0	0	0	0	2
Utilizador 41	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 42	0	0	0	1	0	1	0	0	0	0	0	0	1
Utilizador 43	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 44	0	0	0	0	1	1	1	0	1	0	1	1	3
Utilizador 45	0	1	1	0	1	1	0	0	0	0	0	0	2
Utilizador 46	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 47	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 48	1	1	2	2	0	2	0	0	0	2	0	2	6
Utilizador 49	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 50	3	0	3	1	0	1	2	0	2	1	0	1	7
Utilizador 51	4	0	4	0	1	1	2	0	2	3	0	3	10
Utilizador 52	9	10	19	8	1	9	6	1	7	0	0	0	35
Utilizador 53	0	0	0	0	0	0	0	0	0	1	0	1	1
Utilizador 54	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 55	2	0	2	1	0	1	0	0	0	0	1	1	4
Utilizador 56	260	93	353	97	68	165	10	1	11	5	1	6	535
Utilizador 57	8	0	8	5	1	6	1	1	2	5	2	7	23
Utilizador 58	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 59	0	0	0	1	0	1	5	0	5	4	0	4	10
Utilizador 60	21	4	25	7	1	8	3	1	4	3	0	3	40
Utilizador 61	14	2	16	4	2	6	4	0	4	5	1	6	32
Utilizador 62	32	8	40	19	3	22	12	2	14	0	0	0	76
Utilizador 63	4	0	4	2	0	2	0	0	0	0	0	0	6
Utilizador 64	0	0	0	0	0	0	0	0	0	1	1	2	2
Utilizador 65	21	3	24	4	1	5	2	0	2	0	0	0	31
Utilizador 66	24	0	24	0	13	13	0	0	0	0	0	0	37
Utilizador 67	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 68	1	0	1	7	0	7	0	1	1	1	1	2	11
Utilizador 69	5	2	7	1	2	3	6	0	6	0	1	1	17
Utilizador 70	9	1	10	1	1	2	0	1	1	0	2	2	15
Utilizador 71	0	0	0	0	0	0	0	0	0	3	0	3	3
Utilizador 72	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 73	6	10	16	6	4	10	3	2	5	5	4	9	40
Utilizador 74	0	1	1	0	1	1	0	0	0	0	0	0	2

Utilizador 75	0	0	0	0	0	0	0	1	1	3	2	5	6
Utilizador 76	8	10	18	0	2	2	0	2	2	0	0	0	22
Utilizador 77	2	3	5	0	0	0	0	0	0	0	0	0	5
Utilizador 78	23	5	28	7	4	11	2	2	4	2	2	4	47
Utilizador 79	48	27	75	8	10	18	6	3	9	4	4	8	110
Utilizador 80	4	1	5	0	2	2	3	1	4	1	1	2	13
Utilizador 81	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 82	69	18	87	22	10	32	0	0	0	0	0	0	119
Utilizador 83	0	0	0	1	1	2	0	0	0	1	1	2	4
Utilizador 84	4	1	5	2	1	3	2	0	2	2	1	3	13
Utilizador 85	5	2	7	3	0	3	1	1	2	1	1	2	14
Utilizador 86	0	1	1	1	3	4	1	0	1	0	0	0	6
Utilizador 87	3	7	10	2	0	2	1	0	1	3	0	3	16
Utilizador 88	28	19	47	15	13	28	7	9	16	4	5	9	100
Utilizador 89	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 90	0	0	0	0	1	1	0	0	0	0	0	0	1
Utilizador 91	5	11	16	1	1	2	0	2	2	9	8	17	37
Utilizador 92	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 93	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 94	1	1	2	0	0	0	0	0	0	0	3	3	5
Utilizador 95	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 96	2	0	2	0	1	1	0	1	1	1	1	2	6
Utilizador 97	0	0	0	0	1	1	0	0	0	0	0	0	1
Utilizador 98	26	25	51	24	16	40	0	0	0	0	0	0	91
Utilizador 99	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 100	0	0	0	1	0	1	2	1	3	0	1	1	5
Utilizador 101	0	0	0	0	0	0	0	0	0	2	0	2	2
Utilizador 102	0	1	1	1	1	2	1	2	3	1	0	1	7
Utilizador 103	1	0	1	0	0	0	1	0	1	0	0	0	2
Utilizador 104	4	1	5	2	1	3	0	0	0	1	0	1	9
Utilizador 105	1	2	3	1	0	1	3	0	3	1	0	1	8
Utilizador 106	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 107	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 108	0	0	0	0	0	0	0	0	0	0	1	1	1
Utilizador 109	1	0	1	0	0	0	1	1	2	1	0	1	4
Utilizador 110	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 111	70	12	82	14	6	20	23	6	29	11	3	14	145
Utilizador 112	7	0	7	1	0	1	1	0	1	1	0	1	10
Utilizador 113	27	5	32	9	4	13	3	1	4	1	0	1	50
Utilizador 114	17	2	19	3	2	5	0	0	0	0	0	0	24
Utilizador 115	22	3	25	3	2	5	1	0	1	0	0	0	31
Utilizador 116	36	3	39	3	2	5	0	2	2	0	0	0	46
Utilizador 117	33	2	35	1	5	6	0	0	0	0	0	0	41

Utilizador 118	18	1	19	1	1	2	1	0	1	0	0	0	22
Utilizador 119	8	0	1	0	0	0	0	0	0	0	0	0	1
Utilizador 120	0	0	0	0	0	0	0	0	0	0	0	0	0
Utilizador 121	0	0	0	0	0	0	0	0	0	0	0	0	0

## Anexo 2 – Tabela das percentagens de assertividade dos utilizadores

	Percentagens			
	Acertadas	Erradas	Neutras	Díficeis
Utilizador 1	46%	9%	20%	25%
Utilizador 2	58%	18%	16%	9%
Utilizador 3	33%	33%	17%	17%
Utilizador 4	0%	100%	0%	0%
Utilizador 5	-	-	-	-
Utilizador 6	-	-	-	-
Utilizador 7	47%	50%	2%	1%
Utilizador 8	-	-	-	-
Utilizador 9	-	-	-	-
Utilizador 10	-	-	-	-
Utilizador 11	52%	9%	35%	4%
Utilizador 12	55%	26%	10%	10%
Utilizador 13	-	-	-	-
Utilizador 14	64%	17%	12%	7%
Utilizador 15	45%	30%	10%	15%
Utilizador 16	31%	23%	15%	31%
Utilizador 17	-	-	-	-
Utilizador 18	15%	23%	15%	46%
Utilizador 19	-	-	-	-
Utilizador 20	33%	17%	33%	17%
Utilizador 21	33%	33%	0%	33%
Utilizador 22	0%	50%	50%	0%
Utilizador 23	-	-	-	-
Utilizador 24	0%	0%	100%	0%
Utilizador 25	67%	0%	17%	17%
Utilizador 26	-	-	-	-
Utilizador 27	33%	0%	33%	33%
Utilizador 28	100%	0%	0%	0%
Utilizador 29	-	-	-	-
Utilizador 30	-	-	-	-

Utilizador 31	100%	0%	0%	0%
Utilizador 32	0%	0%	0%	100%
Utilizador 33	0%	0%	0%	100%
Utilizador 34	100%	0%	0%	0%
Utilizador 35	50%	17%	17%	17%
Utilizador 36	100%	0%	0%	0%
Utilizador 37	44%	40%	16%	0%
Utilizador 38	0%	0%	0%	100%
Utilizador 39	20%	20%	30%	30%
Utilizador 40	50%	50%	0%	0%
Utilizador 41	-	-	-	-
Utilizador 42	0%	100%	0%	0%
Utilizador 43	-	-	-	-
Utilizador 44	0%	33%	33%	33%
Utilizador 45	50%	50%	0%	0%
Utilizador 46	-	-	-	-
Utilizador 47	-	-	-	-
Utilizador 48	33%	33%	0%	33%
Utilizador 49	-	-	-	-
Utilizador 50	43%	14%	29%	14%
Utilizador 51	40%	10%	20%	30%
Utilizador 52	54%	26%	20%	0%
Utilizador 53	0%	0%	0%	100%
Utilizador 54	-	-	-	-
Utilizador 55	50%	25%	0%	25%
Utilizador 56	66%	31%	2%	1%
Utilizador 57	35%	26%	9%	30%
Utilizador 58	-	-	-	-
Utilizador 59	0%	10%	50%	40%
Utilizador 60	63%	20%	10%	8%
Utilizador 61	50%	19%	13%	19%
Utilizador 62	53%	29%	18%	0%
Utilizador 63	67%	33%	0%	0%
Utilizador 64	0%	0%	0%	100%
Utilizador 65	77%	16%	6%	0%
Utilizador 66	65%	35%	0%	0%
Utilizador 67	-	-	-	-
Utilizador 68	9%	64%	9%	18%
Utilizador 69	41%	18%	35%	6%
Utilizador 70	67%	13%	7%	13%
Utilizador 71	0%	0%	0%	100%
Utilizador 72	-	-	-	-
Utilizador 73	40%	25%	13%	23%

Utilizador 74	50%	50%	0%	0%
Utilizador 75	0%	0%	17%	83%
Utilizador 76	82%	9%	9%	0%
Utilizador 77	100%	0%	0%	0%
Utilizador 78	60%	23%	9%	9%
Utilizador 79	68%	16%	8%	7%
Utilizador 80	38%	15%	31%	15%
Utilizador 81	-	-	-	-
Utilizador 82	73%	27%	0%	0%
Utilizador 83	0%	50%	0%	50%
Utilizador 84	38%	23%	15%	23%
Utilizador 85	50%	21%	14%	14%
Utilizador 86	17%	67%	17%	0%
Utilizador 87	63%	13%	6%	19%
Utilizador 88	47%	28%	16%	9%
Utilizador 89	-	-	-	-
Utilizador 90	0%	100%	0%	0%
Utilizador 91	43%	5%	5%	46%
Utilizador 92	-	-	-	-
Utilizador 93	-	-	-	-
Utilizador 94	40%	0%	0%	60%
Utilizador 95	-	-	-	-
Utilizador 96	33%	17%	17%	33%
Utilizador 97	0%	100%	0%	0%
Utilizador 98	56%	44%	0%	0%
Utilizador 99	-	-	-	-
Utilizador 100	0%	20%	60%	20%
Utilizador 101	0%	0%	0%	100%
Utilizador 102	14%	29%	43%	14%
Utilizador 103	50%	0%	50%	0%
Utilizador 104	56%	33%	0%	11%
Utilizador 105	38%	13%	38%	13%
Utilizador 106	-	-	-	-
Utilizador 107	-	-	-	-
Utilizador 108	0%	0%	0%	100%
Utilizador 109	25%	0%	50%	25%
Utilizador 110	-	-	-	-
Utilizador 111	57%	14%	20%	10%
Utilizador 112	70%	10%	10%	10%
Utilizador 113	64%	26%	8%	2%
Utilizador 114	79%	21%	0%	0%
Utilizador 115	81%	16%	3%	0%
Utilizador 116	85%	11%	4%	0%



Utilizador 117	85%	15%	0%	0%
Utilizador 118	86%	9%	5%	0%
Utilizador 119	100%	0%	0%	0%
Utilizador 120	-	-	-	-
Utilizador 121	-	-	-	-