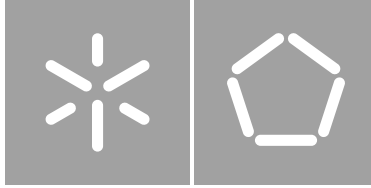


Universidade do Minho
Escola de Engenharia

Henrique Fontão Martins

**Identificação biométrica e comportamental
de utilizadores em cenários de intrusão**



Universidade do Minho
Escola de Engenharia
Departamento de Informática

Henrique Fontão Martins

**Identificação biométrica e comportamental
de utilizadores em cenários de intrusão**

Dissertação de Mestrado
Mestrado em Engenharia Informática

Trabalho realizado sob orientação de

Professor Paulo Novais

AGRADECIMENTOS

O trabalho aqui realizado teve a colaboração de diferentes pessoas. Desta forma cumpre-me endereçar um agradecimento a quem me ajudou, nomeadamente:

- Professor Paulo Novais e Eng. Ricardo Azevedo Pereira que me lideraram, motivaram e dedicaram o seu tempo em prol deste trabalho;
- Professor Henrique Santos, Ângelo Costa e restantes colegas do Laboratório de Inteligência Artificial do Departamento de Informática da Universidade do Minho pela apoio no desenvolvimento deste projeto, nomeadamente pela realização de testes ao trabalho desenvolvido;
- Ricardo Macedo, Vasco Amaral e Samuel Rodrigues meus amigos e colegas de trabalho sempre disponíveis para darem contributos nas alturas em que necessitei;

Por último reservo um especial agradecimento aos meus pais, irmão e tios que ao longo de todo o meu percurso humano e académico foram suporte indispensável nos bons e maus momentos.

RESUMO

A usurpação de contas e o roubo de identidade são problemas muito frequentes nos atuais sistemas informáticos. A facilidade de acesso à internet e a exposição das pessoas a este meio, torna muito frequente a utilização indevida e a usurpação de contas (tais como: e-mail, redes sociais, contas bancárias) por outras pessoas que não as suas legítimas proprietárias.

Atualmente o método de autenticação dominante é o da combinação nome de utilizador e palavra-chave. No entanto, este método pode não ser fiável, pois estas credenciais podem ser partilhadas, roubadas ou até esquecidas. Por outro lado podem-se combinar várias técnicas para reforçar a segurança dos sistemas. Cartões de acesso (*tokens*), certificados digitais e biometrias são algumas delas. Os cartões de acesso, por exemplo os das caixas multibanco, podem ser roubados ou duplicados, como é frequentemente noticiado em fraudes bancárias. Os certificados seguem o mesmo caminho dos *tokens* uma vez que estes podem ser distribuídos por correio eletrónico ou em dispositivos USB. As biometrias físicas (impressão digital, íris, retina ou geometria da mão por exemplo), para além de serem um pouco intrusivas, requerem a aquisição de equipamento caro. Uma possível solução para os problemas inumerados são as biometrias comportamentais.

A forma como nos comportamos e agimos num computador pode ser usada como informação biométrica. Esta informação pode ser utilizada à posteriori, geralmente complementada com mais dados, para identificar, inequivocamente, (ou pelo menos com um determinado grau de confiança) um indivíduo. A informação recolhida pode variar desde o tipo de escrita no teclado, habilidade com o rato, hábitos, cliques, número de páginas abertas, origem do acesso, etc., que depois será sujeita à utilização de algoritmos comportamentais para autenticar, de forma inequívoca, um utilizador.

Neste trabalho pretende-se implementar como reforço aos atuais sistemas de autenticação e de deteção de intrusões, a verificação de perfis comportamentais do proprietário da conta. Este sistema não irá apresentar grandes custos, já que só serão usados equipamentos básicos, e será completamente invisível para o utilizador, ou seja este será continuamente autenticado de forma silenciosa e não intrusiva.

ABSTRACT

Session hijacking and identity theft are a problem increasingly common in computer systems nowadays. With the growing usage of online services, people become more exposed to different techniques, technological or social, that can be used to easy to their personal accounts, from services such as Emails, Facebook, bank accounts, among others.

Currently, the dominant method of authentication is the combination of username and password. This method can be unreliable, because these credentials can be shared, forgotten or stolen. To offer better authentication mechanisms, other techniques are used; among then are the tokens or digital certificates and biometrics. None of them completely solve the problem once they can be duplicated or stolen. Moreover the physiological biometrics (fingerprint, iris, retina, hand geometry, etc.) are intrusive, require the purchase of expensive equipment and may not work in all the scenarios.

The way we behave and act in a computer can be used as biometric information. This information supplemented with more data (i.e. contextual data) can be used to identify unequivocally (or at least with a certain degree of confidence) an individual. The information collected may vary from the way of typing on a keyboard (keystroke dynamics), skill with the mouse (mouse dynamics), habits, clicks, number of pages open, source access, etc., which will then be subject to the use of behavioral algorithms to identify and authenticate, unequivocally, the user.

In this work we present the implementation of a system that strengthens existing authentication and intrusion detection systems, helping them by checking behavioral profiles of the account owner. This system will not be costly, since it only uses basic hardware. Additionally, it will be completely invisible to the user, i.e., it will be working in an unobtrusive way, collecting data in background mode. The aim of this paper is to present a system capable of recognizing biometric patterns and, through behavioral algorithms and complex event processing, create user profiles that are used as identification and continuously authentication to services.

ÍNDICE

1. INTRODUÇÃO.....	1
1.1. Métodos de Autenticação.....	5
1.1.1. Combinação Nome de Utilizador/Palavra-chave	6
1.1.2. Tokens	8
1.1.3. Biométrica.....	9
1.1.4. Combinação de técnicas.....	10
1.1.5. Aplicação.....	11
1.2. Sistemas de Detecção de Intrusões	12
1.3. Desafios	14
1.4. Objetivos	15
1.5. Metodologia da Investigação	17
1.6. Estrutura do Documento.....	17
2. BIOMETRIAS COMPORTAMENTAIS.....	19
2.1. Tecnologias	21
2.1.1. Audit Logs	23
2.1.2. Biometric Sketch	24
2.1.3. Call Stack	25
2.1.4. E-Mail Behaviour.....	26
2.1.5. GUI Interaction	26
2.1.6. Keystroke Dynamics	27
2.1.7. Mouse Dynamics	29
2.1.8. Network Traffic	30
2.1.9. Registry Access.....	31
2.1.10. Storage Activity	31
2.1.11. System Calls.....	32
2.1.12. Tapping.....	32
2.1.13. Web Browsing.....	32
2.2. Problemas da Análise Comportamental.....	33

2.3. Análise Comparativa	34
3. Projetos Relacionados.....	37
3.1. Keytrac.....	37
3.2. BioSig-ID	38
3.3. CVMetrics.....	39
3.4. TypeWATCH	39
3.5. Síntese	40
4. ALGORITMOS DE CLASSIFICAÇÃO.....	41
4.1. Outlier Count	41
4.2. Nearest Neighbor (Mahalanobis)	42
4.3. Algoritmo Combinado	43
4.4. Síntese	43
5. SISTEMA DE AUTENTICAÇÃO ESTÁTICA.....	45
5.1. Arquitetura do Sistema	45
5.2. Aplicação.....	46
5.3. Síntese	48
6. SISTEMA DE AUTENTICAÇÃO CONTÍNUA	49
6.1. Arquitetura do Sistema	49
6.2. Monitorização do teclado	50
6.3. Monitorização do rato	51
6.4. Aplicação.....	59
6.5. Síntese	60
7. CASOS DE ESTUDO.....	63
7.1. Caso de Estudo – Sistema de Autenticação Estática	63
7.1.1. Resultados Obtidos	64
7.1.2. Análise dos Resultados	66

7.2. Caso de Estudo – Sistema de Autenticação Contínuo	67
7.2.1. Resultados Obtidos	68
7.2.2. Análise dos Resultados	72
8. CONCLUSÕES E TRABALHO FUTURO	73
8.1. Síntese Do Trabalho	73
8.2. Trabalho Relevante Realizado.....	74
8.3. Trabalho Futuro	75
REFERÊNCIAS	77

ÍNDICE DE FIGURAS

Figura 1 Número de identidades violadas até Agosto 2012 (milhões).....	2
Figura 2 Causas da violação de identidade até Agosto 2012.....	3
Figura 3 Número de ataques de violação de dados até Agosto 2012 (milhões)	3
Figura 4 Tipo de informação capturada em ataques de violação de dados	4
Figura 5 Tipos de erro	35
Figura 6 Utilização do BioSig-ID	38
Figura 7 Arquitetura do Sistema de Autenticação Estática.....	45
Figura 8 Tempo de Pressão	46
Figura 9 Latência.....	46
Figura 10 Funcionamento Sistema Autenticação Estática	47
Figura 11 Arquitetura do Sistema de Identificação.....	49
Figura 12 Direção do movimento	51
Figura 13 Distância do movimento e a sua Velocidade	53
Figura 14 Comparação entre sessões do mesmo utilizador	54
Figura 15 Comparação entre sessões de utilizadores diferentes	55
Figura 16 Histograma das direções de movimento	56
Figura 17 Velocidade média do movimento por direção.....	57
Figura 18 Histograma dos tipos de ação	58
Figura 19 Velocidade média por tipo de ação	59
Figura 20 Funcionamento Sistema Autenticação Contínua.....	60
Figura 21 Percentagem de erro do algoritmo Outlier Count.....	64
Figura 22 Percentagem de erro do algoritmo Mahalanobis	65
Figura 23 Percentagem de erro do algoritmo Combinado	65
Figura 24 Percentagem de erro do algoritmo Outlier Count para o teclado.....	68
Figura 25 Percentagem de erro do algoritmo Mahalanobis para o teclado.....	69
Figura 26 Percentagem de erro do algoritmo Combinado para o teclado.....	69
Figura 27 Percentagem de erro do algoritmo Outlier para o rato.....	70
Figura 28 Percentagem de erro do algoritmo Mahalanobis para o rato.....	71
Figura 29 Percentagem de erro do algoritmo Combinado para o rato	71

ÍNDICE DE TABELAS

Tabela 1 Tipos de Autenticação	6
Tabela 2 Principais ataques em 2012	7
Tabela 3 Tempo de descoberta de uma palavra-chave de acordo com o seu tamanho	7
Tabela 4 Combinação de técnicas de autenticação.....	11
Tabela 5 Classificação e propriedades das biometrias	21
Tabela 6 Dados possíveis de obter usando Audit Logs	24
Tabela 7 Exemplo da utilização de Biometric Sketch.....	25
Tabela 8 Características do Mouse Dynamics	30
Tabela 9 Comparação de resultados entre biometrias	35
Tabela 10 Comparação de algoritmos de classificação	41
Tabela 11 Exemplo de N-Grafo.....	50
Tabela 12 Comparação de algoritmos.....	66
Tabela 13 Comparação de algoritmos para os diferentes dispositivos	72

1. INTRODUÇÃO

A usurpação de contas, o roubo de identidade e o acesso não autorizado a dados confidenciais são problemas cada vez mais habituais nos sistemas informáticos atuais (ENISA). Com a maior facilidade de acesso à internet e exposição das pessoas neste meio, é muito frequente a utilização e usurpação de contas (e-mail, redes sociais, contas bancárias, etc.) por outras pessoas que não as suas proprietárias.

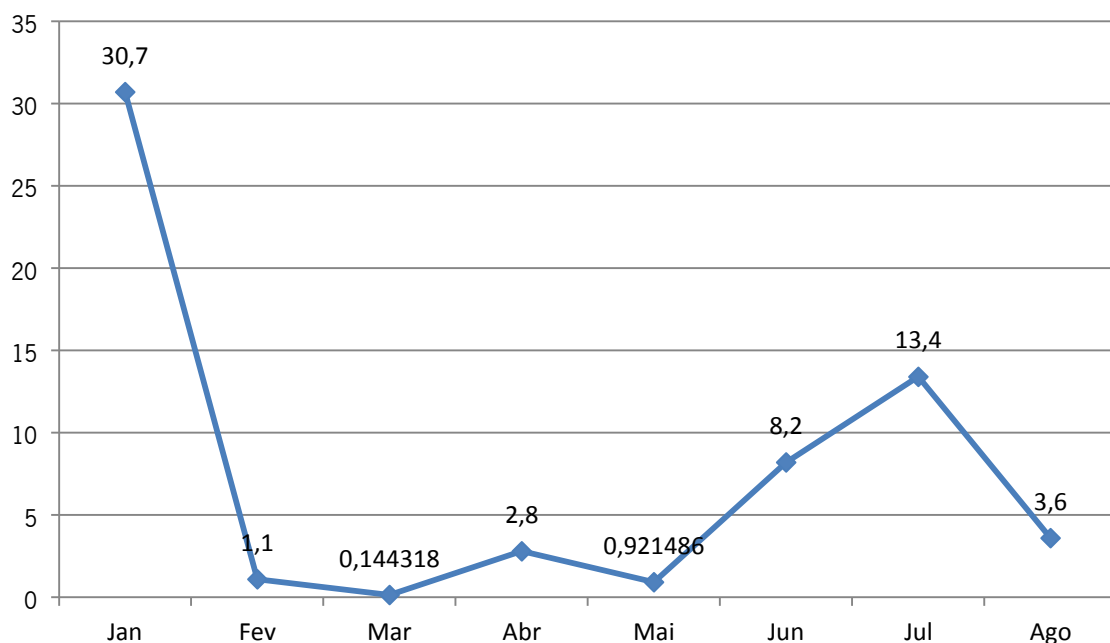
Num mundo cada vez mais ligado, a identidade de um utilizador é a única peça de informação que o torna distinguível. Esta informação é normalmente um par de credenciais (nome de utilizador/palavra-chave) ou outro tipo de informação confidencial como o número de identificação ou número de cartão de crédito. O roubo de identidade é um ataque que ocorre quando um criminoso rouba as credenciais de um utilizador utilizando-as para objetivos maliciosos, principalmente relacionados com fraudes financeiras ou acesso a documentos confidenciais (Marinos & Sfakianakis, 2012).

Segundo o relatório da CLUSIT, (CLUSIT, 2012) o roubo de dados pessoais, e de identidade em geral, é agora o crime mais importante na internet devido ao seu número, em relação aos outros tipos de crime. Os dados capturados servem para alimentar crimes como o acesso não autorizado aos sistemas de informação, fraude de informação e a disseminação de programas que visam danificar ou influenciar o comportamento de outros sistemas. Também no mesmo relatório é referido que as identidades são frequentemente roubadas através de engenharia social, que são técnicas utilizadas para manipular pessoas até estas fornecerem, de forma involuntária, informações confidenciais.

Não estamos perante um problema novo, este tipo de crimes sempre aconteceu, mesmo antes da utilização massiva dos meios eletrónicos. Mas, no passado, o agressor e a vítima tinham que estar sempre em estreita proximidade geográfica. Com o uso generalizado da internet, a situação mudou radicalmente e hoje em dia, normalmente, não há ligação geográfica entre o agressor e a vítima. Além disso, um criminoso pode obter dados de centenas ou milhares de vítimas com muito pouco esforço (Federal Office for Information Security, 2011).

A violação de dados refere-se a uma quebra de informação que ocorreu de forma intencional ou de divulgação de informação involuntária por agentes internos ou externos. Este tipo de ameaça tem como principal alvo vários setores como organizações públicas, organizações de saúde, organizações não-governamentais, empresas de pequenas ou médias dimensões, grandes empresas, etc. Os ataques de violação de dados são geralmente realizadas através de *hackers*, *malware*, ataques físicos, engenharia social ou uso indevido de privilégios (Marinos & Sfakianakis, 2012).

Figura 1 Número de identidades violadas até Agosto 2012 (milhões)



Os ataques de roubo de identidade e de violação de dados estão intimamente relacionados. Normalmente os dados acedidos neste tipo de ataques são informações de contas, números de identificação e nomes de pessoas. Por isso é necessário estudá-los em conjunto de modo a tentar encontrar a melhor forma de evitar e prevenir este tipo de ataques. Nas próximas figuras, adaptadas de (Symantec Intelligence, 2012), vemos o número de identidades violadas em 2012 até Agosto (figura 1) e as principais causas desses ataques no mesmo período de tempo (figura 2).

Figura 2 Causas da violação de identidade até Agosto 2012

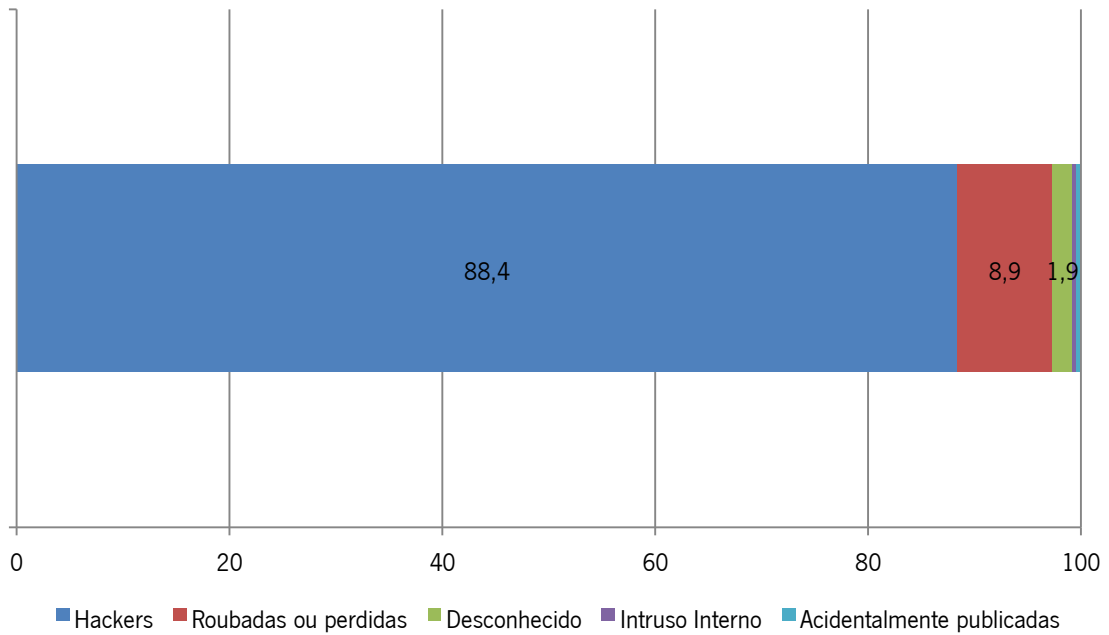
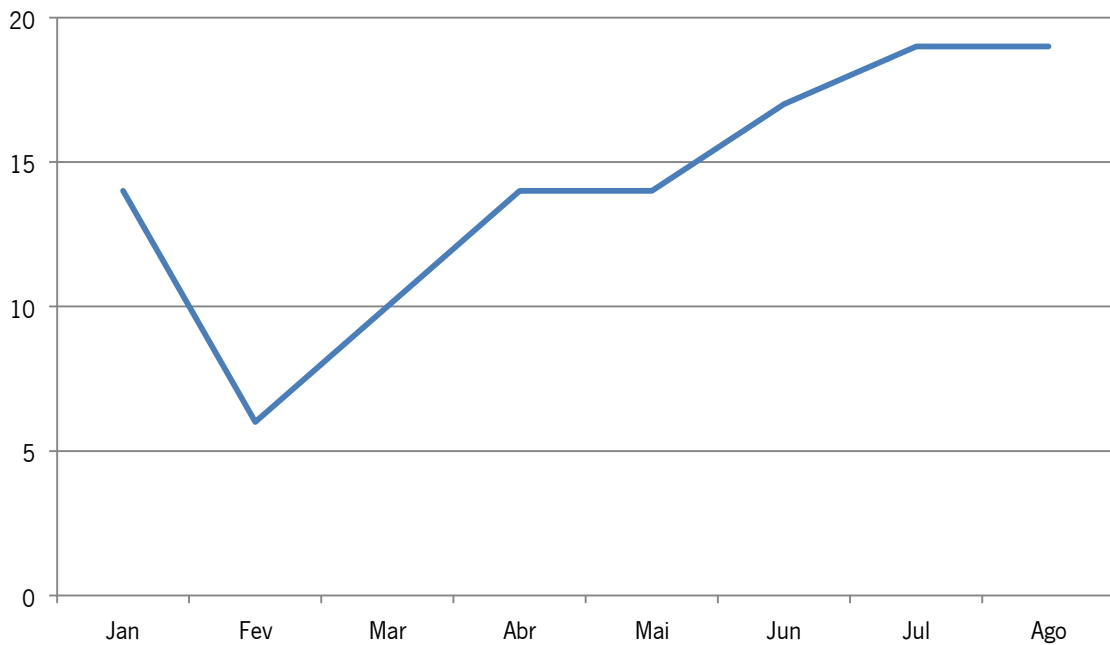


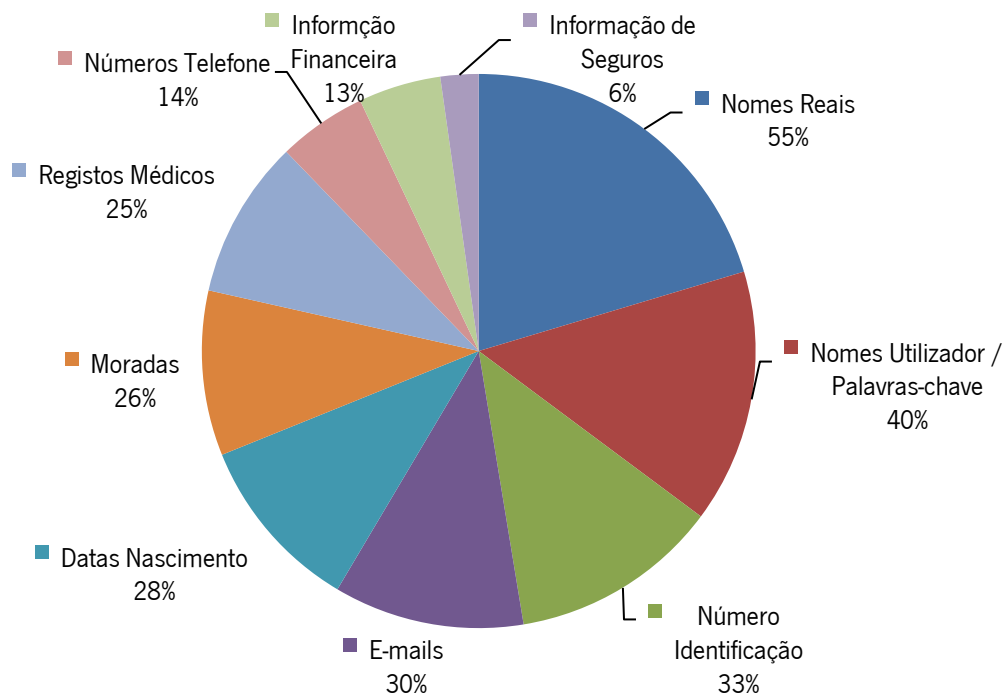
Figura 3 Número de ataques de violação de dados até Agosto 2012 (milhões)



Na figura 3 é apresentado o número de ataques de violações de dados no mesmo período que as figuras anteriores. Comparando com o número de identidades expostas, concluímos que estas seguem caminhos opostos: o número de ataques de violação de dados tem vindo a crescer e

o de roubo de identidades a diminuir. Isto pode significar que os atacantes estão a ficar mais seletivos, apontando para lotes específicos de dados, em vez de os capturar massivamente. A informação que estes estão a capturar pode ser pequena em tamanho, mas é muito mais significativa para fins criminosos.

Figura 4 Tipo de informação capturada em ataques de violação de dados



Como se pode depreender da figura 4, adaptada de (Symantec Intelligence, 2012), vemos à primeira vista que, surpreendentemente, é o verdadeiro nome de um pessoa que é, de longe, a informação mais comum de ser obtida através de um ataque de violação de dados. Superando até a captura da combinação nome de utilizador/palavra-chave, que surge em segundo com 40% dos ataques. Fazendo uma análise mais geral vemos que os quatro tipos de informação mais capturada (nomes reais, nomes de utilizador e password, números de identificação e e-mails) referem-se à tentativa de roubo de identidade ou à captura de dados confidenciais. Isto leva-nos à conclusão que é necessário um reforço dos mecanismos de proteção destes dados sensíveis. Por outro lado, também podemos pensar em formas de essa informação ser inútil depois de capturada, ou seja, fazer com que a informação, quer pessoal quer de autenticação, só possa ser decodificada pela pessoa legítima.

De modo a combater este tipo de problemas têm vindo a ser desenvolvidos métodos de autenticação que visam a garantir a verdadeira identidade dos utilizadores. Esses métodos são apresentados na secção seguinte.

1.1. Métodos de Autenticação

Autenticação é o processo de verificar se as identidades digitais dos computadores e das pessoas físicas são fidedignas. A autenticação de um utilizador é o processo de verificar a sua identidade. Cada ação efetuada deve ser atribuída a uma pessoa específica e isso só é possível de obter com um sistema de autenticação eficiente. Existem várias tecnologias que verificam a identidade de um utilizador antes de lhe conceder o acesso aos recursos do sistema. No entanto, estas tecnologias oferecem diferentes níveis de segurança e nenhuma pode garantir a total segurança do sistema (Karnan, Akila, & Krishnaraj, 2011).

Os sistemas de autenticação são geralmente agrupados em três categorias (O'Gorman, 2003): conhecimento (ex. password), objetos (ex. cartão) e biometrias. Os tipos de autenticação estão descritos de seguida e apresentados de forma resumida na tabela 1 adaptada de (O'Gorman, 2003):

- Baseado em conhecimento (“algo que se sabe”) – é caracterizado pelo sigilo ou obscuridade. Este tipo inclui a palavra-chave memorizada ou outro tipo de informações, que não sendo secretas, podem ser definidas como “segredo para a maior parte das pessoas”. Nome da mãe ou cor preferida podem ser incluídos nessa categoria. O grande problema deste tipo reside na partilha destes dados de autenticação, sendo que estes vão ficando cada vez menos secretos.
- Baseado em objetos (“algo que se tem”) – é caracterizado pela posse física. Chaves físicas são objetos que têm resistido ao longo do tempo. O problema deste tipo de chaves é que, se perdidas, permitem ao seu localizador o acesso ao recurso correspondente. Por isso, é recomendado que chaves digitais combinem sempre outro fator, por exemplo uma palavra-chave, para as proteger do furto ou perda. A vantagem de se ter um objeto físico é que em caso de perda, esta é notada, quase de imediato, e podem-se tomar as respetivas diligências.

- Baseado em identidade (“algo que se é”) – é caracterizado pela singularidade de uma pessoa. Cartas de condução, passaporte, cartão de crédito, diploma universitário, etc., todos pertencem a esta categoria. Também as técnicas biométricas (impressão digital, leitura da íris, assinatura, etc.) se encaixam neste tipo. Para ambos, os documentos de identificação e biometrias, a maior vantagem é a grande dificuldade de cópia ou falsificação. Porém, se uma biometria é comprometida ou se um documento é perdido, eles não são tão facilmente substituíveis como um cartão ou uma palavra-chave.

Tabela 1 Tipos de Autenticação

Métodos de Autenticação				
	Baseados em Conhecimento	Baseados em Objetos	Baseados em ID	
Conhecidos como:	Palavra-chaves, segredos	Cartões de acesso	Biometrias	
Exemplos	Tradicional	Cadeado	Chave metálica	Diploma
	Digital	Palavra-chave	Smart Card	Impressão digital
Principal problema	Menos secreto em cada utilização	Inseguro em caso de perda		Difícil de substituir

Em seguida serão apresentados os grandes representantes de cada categoria: a combinação nome de utilizador/palavras-chave, *tokens* e biometrias.

1.1.1. Combinação Nome de Utilizador/Palavra-chave

Atualmente o método de autenticação dominante é o da combinação nome de utilizador e palavra-chave. Este método pode ser vulnerável já que as credenciais podem ser partilhadas, roubadas ou até esquecidas. Estudos mostram que “adivinhar” uma palavra-chave pode ser mais fácil do que parece e estima-se que 55% dos cibernautas utilizam palavras-chave fracas, que podem ser descobertas em poucas horas (Pozadzides, 2007) (Mahmood, 2010). Como exemplo disso temos a notícia de que em Junho de 2012 cerca de 6,5 milhões de palavras-chave da rede social

LinkedIn foram descobertas por um grupo de *hackers* russos (Brian, 2012) (Kamp, 2012). A tabela 2, adaptada de (ENISA, 2012), demonstra os principais ataques realizados em 2012.

Tabela 2 Principais ataques em 2012

Companhia	Contas afetadas	Data
LinkedIn	6 500 000	06/06/2012
EHarmony	1 500 000	06/07/2012
Formspring	420 000	10/07/2012
Yahoo Voice	400 000	12/07/2012
Android Forums	1 000 000	12/07/2012
NVIDIA	400 000	13/07/2012
Gamigo	8 240 000	24/07/2012

Na tabela 3 vemos o tempo médio, num computador pessoal, para a descoberta de uma palavra-chave. Convém notar que uma palavra-chave com até 6 caracteres o tempo para a sua descoberta é relativamente pequeno. (Pozadzides, 2007) (Mahmood, 2010).

Tabela 3 Tempo de descoberta de uma palavra-chave de acordo com o seu tamanho

Tamanho da Palavra-Chave	Usando Todos os Carateres	Usando Só letras minúsculas
3	0,86 segundos	0,02 segundos
4	1,36 minutos	0,046 segundos
5	2,15 horas	11,9 segundos
6	8,51 dias	5,15 minutos
7	2,21 anos	2,23 horas
8	2,10 séculos	2,42 horas
9	20 milénios	2,07 meses

De modo a tornar esta forma de autenticação mais fiável e segura existem vários manuais de várias práticas. Em (ENISA, 2012) é sugerido:

- A não utilização da mesma palavra-chave em contas diferentes;

- A sua mudança imediata em situações de roubo. No caso da utilização da mesma ou uma variação da palavra-chave noutras contas, estas também deverão ser alteradas;
- A utilização de uma palavra-chave complexa com um mínimo de 8 caracteres que contenham alfanuméricos e caracteres especiais (i.e. caracteres a-z, A-Z, 0-9 e ',&@:?!()\$#/\);
- A mudança regular de palavra-chave (no máximo de 90 em 90 dias);
- A utilização de serviços com autenticação multi-fator.

No entanto, e mesmo com a utilização destas práticas, o problema do esquecimento e da partilha de palavras-chave não fica totalmente resolvido. É necessário um reforço a este sistema de modo a torná-lo mais eficiente, seguro e fiável.

1.1.2. Tokens

Um *token* de identidade, de segurança, de acesso, ou simplesmente *token*, é um dispositivo físico que realiza a autenticação. Estes dispositivos variam desde armazenamentos seguros que contêm uma palavra-chave (ex. um cartão bancário), um comando para abertura de uma garagem, ou um cartão inteligente (*smart card*). Também pode ser um dispositivo ativo que gera palavras-chave de uma só utilização, síncronas (que mudam em sincronismo com um servidor), ou de desafio-resposta (respondendo a um desafio de uma única vez). Um *token*, normalmente, contém uma embalagem resistente a violação e um *hardware* especial que desativa o *token* se este for violado ou se o número de tentativas de autenticação exceder um determinado limite (O'Gorman, 2003).

Devido à sua grande vulnerabilidade ao roubo, um *token* nunca deverá ser utilizado de forma individual. Por isso, na maioria dos casos, é utilizado sempre em combinação com uma palavra-chave (tal como acontece nos cartões bancários). Quando combinado com uma palavra-chave este mecanismo pode oferecer muitas vantagens como (O'Gorman, 2003):

- *Single Sign-on*: Como um *token* pode armazenar ou gerar múltiplas palavras-chave, este pode ser utilizado para aceder a vários sistemas. O problema da memorização de várias palavras-chave pode ser assim resolvido;

- A observação da sua ausência (a perda ou esquecimento de uma palavra-chave não o são). O que permite uma rápida ação contra uma possível perda ou furto;
- Oferecer uma maior segurança à comum palavra-chave. Agora não basta a um impostor descobrir a palavra-chave, para este conseguir aceder ao sistema tem que também ter em sua posse o respetivo *token*.

Olhando agora às desvantagens, este tipo de dispositivos têm um custo elevado e como já referido anteriormente, a sua perda ou roubo têm uma probabilidade elevada. Como exemplo desse caso temos os cartões de crédito que são frequentemente noticiados em fraudes bancárias (Público, 2013). Os certificados digitais, embora não sejam físicos, podem também ser considerados como *tokens*. Porém, sofrem dos mesmos problemas, uma vez que estes podem ser distribuídos por correio eletrónico ou em dispositivos USB.

1.1.3. Biométrica

A autenticação biométrica é um método automático que identifica ou verifica a identidade baseando-se na medida das suas características físicas únicas (face, palma da mão, íris, etc.) ou em características comportamentais (voz, escrita à mão, assinatura, dinâmica da digitação, etc.). Biometrias físicas são traços biológicos que são inatos ou desenvolvidos naturalmente e as biometrias comportamentais são hábitos ou maneirismos que são aprendidos ou adquiridos ao longo do tempo (Karnan, Akila, & Krishnaraj, 2011).

Ambos os sistemas, físicos ou comportamentais, são divididos em duas fases: fase de aprendizagem e fase de autenticação. Durante a fase de aprendizagem são adquiridos os dados biométricos da pessoa, que, posteriormente, são processados e armazenados como perfil de referência para esse indivíduo. Este perfil é usado como um modelo para uso futuro pelo sistema de autenticação. A decisão de autenticação é baseada na correspondência entre a amostra e o perfil de referência previamente armazenado (Karnan, Akila, & Krishnaraj, 2011) (Giot, El-Abed, Hemery, & Rosenberger, 2011).

As biometrias são excelentes candidatas à verificação de identidades pois, ao contrário de cartões de acesso ou palavras-chave, estas não podem ser esquecidas, roubadas ou perdidas, e na ausência de lesões físicas são potencialmente o método mais eficaz para determinar a identidade

de um indivíduo. Adicionalmente, combinando técnicas biométricas com métodos de autenticação tradicionais, estes podem garantir um nível de segurança extraordinário (Monrose & Rubin, 1999).

As biometrias físicas são, atualmente, as mais implementadas, devido à grande variabilidade dos comportamentos humanos. Estes podem variar consecutivamente entre amostras, uma vez que são dependentes do desempenho humano. Porém, as biometrias físicas apresentam, dois grandes inconvenientes: custo e intrusão. A aquisição de equipamentos específicos, como leitores ou câmaras, são necessários para a leitura das características físicas das pessoas. Noutros casos, como por exemplo a leitura da íris, podem ser considerados intrusivos, e por isso não aceites por alguns utilizadores (Ferreira, Santos, & Patrão, 2011).

1.1.4. Combinação de técnicas

Diferentes tipos de técnicas de autenticação podem ser combinadas para aumentar a segurança dos sistemas (ver Tabela 3). Isto é chamado de autenticação multi-fator ou autenticação em vários níveis. Por razões de segurança, a autenticação só é válida caso esta satisfaça todas as técnicas em utilização.

O exemplo mais comum da autenticação multi-fator é a do cartão multibanco. A combinação do cartão com o código numérico – autenticação em dois passos – garante uma maior segurança do que o cartão sozinho, que pode ser roubado, enquanto um cartão protegido com código não pode ser usado sem o conhecimento prévio desse segredo. Este exemplo da utilização de um *token* com uma palavra-chave constitui a maioria dos sistemas multi-fator atuais.

Para evitar o uso da senha, que pode ser esquecida, uma identificação biométrica pode ser uma alternativa para proteger o *token*. Porém esta abordagem implica um custo superior com o equipamento biométrico.

Geralmente a autenticação em três fatores não é utilizada, devido à elevada complexidade do sistema de autenticação, ficando restringida a aplicações de alta segurança (ex. bases militares).

A tabela 4 apresenta as combinações de técnicas de autenticação apresentando as suas vantagens e desvantagens.

Tabela 4 Combinação de técnicas de autenticação

Combinação	Vantagens	Desvantagens	Exemplo
Conhecimento e Objeto	<i>Token</i> perdido/roubado é protegido por um código secreto	É necessário a posse do <i>token</i> e a memorização do código secreto	Cartão bancário protegido por PIN
Objeto e ID	<i>Token</i> perdido/roubado é protegido por uma biometria	É necessária a posse do <i>token</i>	Cartão de identidade com fotografia
Conhecimento e ID	Dois fatores que fornecem segurança	Necessária a memorização de uma palavra-chave	Palavra-chave mais impressão digital para acesso a um computador
Conhecimento, Objeto e ID	Três fatores que fornecem segurança	Memorização de palavra-chave e posse de <i>token</i>	Aplicações militares que requerem cartão de identificação mais palavra-chave

1.1.5. Aplicação

Todas as técnicas acima apresentadas têm vantagens e inconvenientes, porém todas elas podem ser combinadas para maximizar a segurança de um sistema. O relatório (Trustwave, 2012) apresenta algumas técnicas e boas práticas para o desenvolvimento de sistemas de autenticação.

O primeiro passo para a construção de um sistema de autenticação fiável passa pela eliminação de contas genéricas, partilhadas ou desnecessárias. Todos os anos se regista um significativo número de violações de acesso que ocorrem como resultado de atacantes que conseguem obter acesso a estas contas do sistema. Por isso é muito importante a realização de uma análise periódica de todos os utilizadores e a definição de políticas de complexidade de palavras-chave. A adoção de métodos de autenticação ou a utilização de biometrias, também contribuem, de forma importante, para o aumento da robustez e eficiência da autenticação:

- Revisão da Lógica de Acesso – Efetuar análises periódicas a todos os utilizadores e funções (*roles*) irá melhorar a segurança em volta dos níveis de acesso e até pode identificar contas redundantes e que não são necessárias ao sistema;
- Políticas de complexidade de palavras-chave – Definir políticas de complexidade e ensinar os utilizadores sobre as melhores técnicas, como a utilização de frases secretas;
- Autenticação multi-fator – Autenticação em dois passos permite que os utilizadores sejam autenticados pelo que sabem (palavra-chave) e pelo que têm (dispositivo ou certificado). Isto não deve ser só aplicado ao mundo digital, mas também ao mundo físico. A combinação cartão com um código secreto é o maior exemplo deste caso.
- Biometrias – Além de palavras-chave e outros mecanismos de autenticação, as biometrias podem ser necessárias para áreas mais sensíveis ou ambientes mais restritos.

1.2. Sistemas de Detecção de Intrusões

Os sistemas de deteção de intrusões são principalmente baseados em duas abordagens: conhecimento ou comportamento. Neste estudo ir-nos-emos focar na deteção de intrusões baseadas em comportamento. Convém notar que, atualmente, existem poucas soluções que implementem esta abordagem, mesmo que em (Denning, 1987) se reconheça que é um requisito importante para um sistema de deteção de intrusões (IDS).

As técnicas de deteção de intrusões baseadas em comportamento assumem que uma intrusão pode ser detetada através da observação de um desvio de um comportamento normal ou esperado do sistema ou dos utilizadores. O modelo de comportamento normal é extraído a partir da informação recolhida, durante a fase de aprendizagem, através de vários meios. O detetor de intrusões depois compara o modelo com a atividade atual e se um desvio for observado, é despoletado um alarme. Isto significa que qualquer atividade que não corresponda a um comportamento previamente aprendido é considerada intrusivo. Por isso, o sistema de deteção de intrusões deverá detetar todos os ataques mas a sua precisão fica um pouco aquém, já que existem um largo número de falsos alarmes (Debar, 2010).

Este tipo de abordagem tem a vantagem de conseguir detetar tentativas de intrusão e vulnerabilidades novas e imprevistas. Até pode contribuir para a descoberta automática destes

novos ataques. Também são menos dependentes dos mecanismos específicos dos sistemas operativos e ajudam a detetar ataques de “abuso de privilégios” que não envolvem nenhuma vulnerabilidade específica. Em suma, esta abordagem fica resumida numa só frase: “Tudo o que não foi visto anteriormente é uma ameaça” (Debar, 2010).

A elevada taxa de falsos alarmes é geralmente citada como a grande desvantagem das técnicas baseadas em comportamento porque todo o âmbito do comportamento de um sistema de informação pode não ser coberto na fase de aprendizagem. Além disso, os comportamentos variam com o tempo causando a necessidade de treino periódico do perfil de comportamento, resultando, tanto numa indisponibilidade do sistema de deteção de intrusões, como noutros falsos alarmes. O sistema pode também sofrer ataques durante a fase de aprendizagem e por isso conter comportamento errático no seu perfil e que depois não é reconhecido como anormal (Debar, 2010).

Falando do ponto de vista dos intrusos, estes podem ser classificados como externos ou internos. Intrusos externos são aqueles que não conseguem a aceder fisicamente ao sistema alvo. Intrusos internos são os que podem aceder diretamente e fisicamente ao sistema mas não estão autorizados a aceder aos seus dados, programas ou recursos. Este tipo de intrusos inclui:

- Utilizadores mascarados (que operam utilizando as credenciais (nome de utilizador/palavra-chave) de outro indivíduo);
- Utilizadores clandestinos (que quebram o controlo de acesso);
- Utilizadores abusivos (que utilizam recursos que não estão nos seus privilégios).

Os sistemas de deteção de intrusões (IDS) existentes, sendo grande parte deles baseados em conhecimento e direccionados para o tráfego de rede, estão bem preparados para ameaças externas. Mas as ameaças internas também apresentam grande perigo para pessoas e organizações. Por isso é necessário um reforço de sistemas capazes de prevenir e detetar situações de ameaças internas.

No desenvolvimento deste trabalho deparamo-nos com inúmeros desafios que têm de ser identificados. Esses desafios são apresentados na próxima secção.

1.3. Desafios

Os desafios inerentes a este projeto podem ser classificados como de natureza social, económica ou tecnológica.

Do ponto de vista social, é preocupante a forma como as pessoas e algumas organizações desvalorizam a autenticação segura e o problema do roubo de identidade. Quando vemos indivíduos que apontam as suas credenciais e as guardam numa gaveta, ou ainda pior, criam um documento onde têm organizadas todas as suas credencias, isto levanta um enorme desafio de mudança de hábitos e mentalidades. Com a agravante de, com o evoluir do tempo, os serviços essenciais do dia-a-dia (bancos, finanças, segurança social, entre outras) estarem a apostar na sua informatização e disponibilização *online* é necessário instruir e informar a sociedade dos potenciais problemas que o roubo de identidade pode levantar. Simultaneamente, estas preocupações também deveriam ser consideradas pelas empresas. A espionagem industrial é potencialmente um grande problema e basta uma credencial de acesso ao correio eletrónico ir parar a mãos erradas que podemos ter um produto confidencial a ser constantemente vigiado pela concorrência.

Do ponto de vista económico, e tendo em conta o estado financeiro de alguns dos países é necessário a adoção de políticas de segurança que tenham custos reduzidos. Como vimos, a proteção de identidades ainda não é levada a sério e se considerarmos a agravante de os produtos que visam a sua solução sejam bastante caros, estamos a contribuir para que este problema não seja solucionado. No entanto, a proteção e a adoção de sistemas de deteção de intrusões podem contribuir para um decréscimo do número de fraudes existentes neste momento. Como vimos nos seções anteriores, as causas que alimentam o problema do roubo de identidades são maioritariamente económicas. Combatendo este problema estaremos também a resolver e a diminuir a quantidade de crimes eletrónicos existentes neste momento.

Do ponto de vista tecnológico, que não relega preocupações económicas, é preciso tirar o máximo partido do que a tecnologia atual nos oferece. A utilização de equipamento básico e existente em todas as casas é importante quer pelo aspeto económico, quer pelo aspeto da reutilização e eficiência dos sistemas. A adoção de técnicas não intrusivas também é um grande desafio nesta área. A biometria possui técnicas de autenticação que são altamente fiáveis mas o problema da privacidade faz com que estas técnicas (leitura da íris, reconhecimento facial) sejam postas de parte pela população. Por isso é importante a adoção de técnicas que sejam socialmente

aceitáveis e que não exijam grande mudança nos procedimentos normais. Um sistema com uma curva de aprendizagem muito reduzida tem uma maior probabilidade de sucesso e contribui para uma adoção mais fácil e mais generalizada deste tipo de produtos na sociedade.

1.4. Objetivos

O instante de início de sessão num sistema é de grande importância já que, caso o acesso se conceda, são desbloqueados recursos confidenciais. Como já vimos anteriormente, todos os métodos de autenticação tradicionais apresentam problemas. As palavras-chave são facilmente esquecidas, partilhadas ou roubadas, os *tokens* podem ser duplicados e também partilhados, e as biometrias atuais, além de poderem ser intrusivas, pressupõem a aquisição de equipamento caro. É necessário a criação e desenvolvimento de um método de autenticação multi-fator que minimize todos estes problemas.

Outro problema pertinente surge depois de um início de sessão. O período entre início de sessão e fecho de sessão não é monitorizado, isto é, parte-se do princípio que o utilizador autenticado é aquele que usa a sessão durante todo o seu período. Mas isto pode não ser verdade. Quantas vezes existem situações de utilizadores que se esquecem de terminar a sua sessão? Quantas vezes em organizações, os utilizadores, mesmo não estando no computador, os deixam ligados e desbloqueados? E no caso da perda das credenciais? De modo a resolver todas estas vulnerabilidades é necessário a existência de um modelo de autenticação contínua, de modo que o utilizador esteja constantemente a ser vigiado, assegurando assim que as contas ou computadores estão a ser operados por utilizadores fidedignos. A ameaça de intrusos internos e mascarados é cada vez maior, sendo por isso necessário haver ferramentas eficientes, não intrusivas e invisíveis (que não alterem o funcionamento normal dos sistemas) que impeçam este tipo de vulnerabilidades.

Esta dissertação tem como objetivo o estudo de mecanismos que, com base em informação recolhida, consigam construir um perfil comportamental único para cada utilizador, e a construção de protótipos que apresentem e demonstrem a validade deste tipo de abordagem. A informação recolhida pode variar desde o tipo de escrita no teclado, habilidade com o rato, hábitos, cliques, número de páginas abertas, origem do acesso, etc., que depois será sujeita à utilização de

algoritmos comportamentais para autenticar, de forma inequívoca, um utilizador. Assim pretende-se dar resposta à pergunta: É possível caracterizar um utilizador através da forma como ele interage com determinada aplicação ou serviço e com isso obter um determinado grau de confiança de que o utilizador que se identificou/autenticou é verdadeiro?

O produto final será a construção de um protótipo contendo dois sistemas de autenticação: autenticação estática e autenticação contínua. Estes sistemas irão ser incluídos na plataforma IAM (Identity and Access Manager) desenvolvida pela PT Inovação. Esta solução terá duas funções principais: garantir a autenticação eficiente no ato de início de sessão e fornecer proteção contínua contra intrusões até ao ato de fecho de sessão dos utilizadores registados nesta plataforma. O principal foco deste produto será de autenticar utilizadores de forma “invisível”, sem alterar o seu modo de interação e sem a aquisição de qualquer equipamento adicional. Por outras palavras, pretende-se combater intrusões de forma permanente, permitindo aos utilizadores a realização de todas as tarefas sem qualquer mudança de procedimentos.

Em suma, os objetivos deste trabalho são:

- Estudar os comportamentos de utilizadores na sua interação com o computador;
- Definir perfis comportamentais únicos com base nessas interações;
- Monitorizar identidade real do utilizador no ato da autenticação;
- Definição de algoritmos de classificação adequados à deteção de desvios comportamentais;
- Execução de testes em ambientes reais.

Concluindo, nesta dissertação pretende-se implementar um reforço aos atuais sistemas de autenticação e de deteção de intrusões, auxiliando-os pela verificação de perfis comportamentais do proprietário da conta. Este sistema não terá custos extraordinários, já que só serão usados *hardware* e equipamentos básicos, e será completamente invisível para o utilizador, ou seja este irá ser autenticado e verificado de forma completamente silenciosa e não intrusiva.

1.5. Metodologia da Investigação

Este trabalho foi desenvolvido ao longo de um ano letivo, nomeadamente durante o segundo ano do Mestrado em Engenharia Informática. Foi utilizada a metodologia de acção/investigação (*action research* methodology). O tema foi proposto pela empresa PT Inovação, onde foi identificada esta área, ainda pouco desenvolvida, para investigação e produção de soluções. Seguidamente foram construídos um plano de trabalho e um documento contendo os objetivos a atingir nesta investigação.

Após a formulação do plano foi então iniciada a pesquisa e análise de projetos existentes atualmente e projetos já terminados, esta etapa é continuamente renovada de acordo com novas ideias ou informação que vai aparecendo ao longo do trabalho. Na última fase é criado um modelo funcional que possa atingir os resultados esperados. Esta fase contém uma grande porção do projeto pois é a que produz resultados palpáveis e é também uma fase em constante mudança pois normalmente vão chegando continuamente novos dados. Esta metodologia segue o seguinte padrão:

- Especificação do problema e as suas características.
- Atualização constante do estado da arte, renovando os objetivos do projeto.
- Desenvolvimento de um protótipo que corresponda aos resultados definidos no projeto.
- Análise e correção do protótipo tendo em conta os resultados obtidos.
- Difusão e partilha do conhecimento com a comunidade científica.

1.6. Estrutura do Documento

Esta dissertação inclui um levantamento do estado da arte (capítulo 2), uma proposta de solução e implementação (capítulos 3 e 4), uma validação, demonstrando de uma forma prática os seus resultados (capítulo 5) e por fim as conclusões e trabalho futuro (capítulo 6).

O Capítulo 2 contém a descrição das biometrias comportamentais baseadas em interação humano-computador. São explicadas todas as suas qualidades e propriedades. São apresentadas as tecnologias mais importantes desta área, como estas funcionam e como podem resolver o

problema da detecção de intrusões. De seguida são apresentados os principais problemas deste tipo de abordagem e a forma como se podem evitar ou resolver essas lacunas. Concluindo com uma amostra dos projetos e soluções comerciais já existentes nesta área.

No Capítulo 3 é apresentado o sistema de autenticação estática proposto. Começa por ser mostrada de forma simples a arquitetura de funcionamento do sistema. Passando depois para a explicação das opções tomadas, e o funcionamento geral dos diversos componentes do sistema.

O Capítulo 4 segue a ordem de apresentação do capítulo 3, mas aqui é apresentado o sistema de autenticação contínua. Este sistema que por ser radicalmente mais complexo que o primeiro, exige uma maior explicação do seu funcionamento e arquitetura.

No Capítulo 5 são apresentados os algoritmos de classificação usados neste projeto. Estes algoritmos têm a função de classificar um indivíduo como fidedigno ou intruso. É explicada a forma como foram escolhidos estes algoritmos, como são implementados e alguns resultados que estes obtiveram em projetos similares.

No Capítulo 6 são apresentados todos os resultados produzidos pelos sistemas desenvolvidos. Começando por explicar os seus casos de utilização, passa depois por definir os ambientes de teste e conclui com a análise dos resultados produzidos usando as diversas abordagens possíveis e os diversos níveis de “rigidez” dos algoritmos de classificação.

Finalmente no Capítulo 7 apresenta-se a síntese do trabalho feito neste projeto, tal como as importantes contribuições e trabalho gerado. Terminando com as conclusões obtidas e o trabalho futuro a ser realizado.

2. BIOMETRIAS COMPORTAMENTAIS

Neste capítulo serão apresentadas um conjunto de biometrias comportamentais que poderão ser usadas para a construção de sistemas de detecção de intrusões baseados em comportamento. Segundo (Yampolskiy & Govindaraju, 2008) as biometrias comportamentais fornecem um grande número de vantagens em relação às tecnologias biométricas tradicionais. Elas podem ser obtidas de uma maneira não intrusiva ou até sem o conhecimento do utilizador. Para construir coleções de dados de biometrias comportamentais não é necessário nenhum equipamento especial, sendo por isso uma técnica de baixo custo monetário. Enquanto a maior parte das biometrias comportamentais, aplicadas individualmente, não são únicas para fornecer uma identificação humana fiável, elas podem verificar identidades, ou seja autenticar indivíduos, de forma muito acertada.

Um sistema biométrico pode ser dividido em duas aplicações diferentes: autenticação ou identificação. Autenticação fornece uma verificação da identidade dos utilizadores com o objetivo de confirmar ou negar a identidade de um indivíduo, enquanto que a identificação refere-se a estabelecer a identidade de um indivíduo:

- Autenticação: Verificar a identidade de um indivíduo, ao medir o seu comportamento e compará-lo com um perfil de referência, perfil esse já previamente construído durante uma fase de aprendizagem (comparação 1 para 1).
- Identificação: Obter uma grande amostra de dados e identificar o utilizador comparando o seu comportamento com todos os perfis construídos (comparação 1 para muitos).

As biometrias são classificadas segundo características de universalidade, singularidade, continuidade, coletividade, performance, aceitabilidade e evasão. No que diz respeito a biometrias comportamentais, estas são caracterizadas pelas seguintes propriedades:

- Universalidade: As biometrias comportamentais são dependentes das habilidades específicas possuídas por diferentes pessoas, mas tendo em causa a população mundial a sua universalidade é reduzida. Mas desde que as biometrias comportamentais só sejam aplicadas em ambientes restritos e específicos, a sua universalidade pode até ser de 100%.

- Singularidade: Uma vez que existem poucas maneiras de realizar uma qualquer tarefa, a singularidade deste tipo de biometrias é baixa. Estilos de escrita diferentes, estratégias diferentes e preferências variadas são suficientes para autenticar utilizadores mas não para os identificar, a menos que o conjunto de utilizadores seja muito pequeno.

- Continuidade: Estas biometrias apresentam um baixo grau de permanência pois elas medem um tipo de comportamento que muda à medida que a pessoa aprende técnicas mais avançadas e rápidas para cumprir as suas tarefas. Porém, este problema de desvio de conceito é abordado nas investigações de sistemas de intrusão baseados em comportamentos e os sistemas, através da utilização de algoritmos adaptativos, são capazes de se ajustar às mudanças de comportamento dos utilizadores.

- Recolha: A recolha de dados biométricos é relativamente fácil e não-intrusivo para o utilizador. Em alguns casos, o utilizador nem sabe que os dados estão a ser recolhidos. O método de recolha de informação é completamente automatizada e tem um custo muito baixo.

- Performance: A precisão de identificação de indivíduos por maior parte das biometrias comportamentais são particularmente baixas a partir do momento que as bases de dados começam a ficar grandes. Porém a precisão da autenticação de indivíduos é elevada para a maior parte das técnicas.

- Aceitabilidade: A partir do momento em que as biometrias comportamentais podem ser recolhidas sem a participação do utilizador, estas disfrutam de um elevado grau de aceitabilidade, mas podem ser recusadas por razões de ética ou de privacidade.

- Invasão: É relativamente difícil contornar sistemas biométricos comportamentais, pois requerem o conhecimento íntimo do comportamento de uma pessoa, mas uma vez que esse comportamento seja exposto e esteja disponível pode ser replicado. É por isso que é extremamente importante manter os perfis comportamentais dos utilizadores encriptados ou então protegidos de maneira a que não se consiga associar o perfil a um determinado indivíduo.

2.1. Tecnologias

Existem várias categorias de biometrias comportamentais, mas aqui só iremos analisar um tipo: biometrias baseadas em interações humano-computador (*HCI-based biometrics*). A tabela 5, retirada de (Yampolskiy & Govindaraju, 2008), demonstra as propriedades das biometrias HCI (D = dias, M = Minutos, S = segundos).

Tabela 5 Classificação e propriedades das biometrias

Biometria	Interação direta entre humano e computador		Interação indireta entre humano e computador	Habilidade Motora	Puramente Comportamental	Propriedades das biometrias comportamentais			
	Input	Software				Tempo de Aprendizagem	Tempo de Verificação	Identificação	Hardware Necessário
Audit Logs			X			D	D	N	Computador
Biometric Sketch					X	M	S	N	Rato
Call Stack			X			D	H	N	Computador
E-mail Behaviour		x			X	D	M	N	Computador
GUI Interaction			X			D	H	N	Computador
Keystroke Dynamics	x			x		M	S	N	Teclado
Mouse Dynamics	x			x		M	S	N	Rato
Network			X			D	D	N	Computador

Traffic						
Registry		X		D	H	N
Access						Computador
Storage		X		D	D	N
Activity						Computador
System Calls		X		D	H	N
Tapping			x	M	S	N
Web						
Browsing	x		X	D	M	N

Fazendo uma interpretação à tabela 5, vemos que as biometrias estão divididas em dois tipos de interação: direta e indireta.

Interação direta corresponde à interação diária com computadores, onde os seres humanos adotam estratégias diferentes, estilos diferentes e aplicam habilidades e conhecimentos únicos. Essas características, depois de quantificadas, podem ser usadas para a construção de perfis que irão ser usados na verificação de identidades. Este tipo de biometrias pode ser dividido em categorias adicionais. A primeira categoria consiste na interação humana com dispositivos de *input* como teclados – *keystroke dynamics* (dinâmica da digitação) – ou ratos – *mouse dynamics* (dinâmica do rato) – que conseguem registrar ações musculares que são inerentes, distintivas e consistentes. A segunda categoria mede o comportamento humano avançado, tais como conhecimento, estratégia ou habilidade exibida pelo utilizador durante a interação com os diversos *softwares* diferentes.

Interação indireta consiste em eventos que podem ser obtidos ao monitorizar utilizadores através de ações baixo nível dos *softwares*. Estas ações incluem *system calls* (chamadas ao sistema), *audit logs* (registos de auditoria), *registry access* (acessos a registos), *storage activity* (atividade do armazenamento), *call-stack data analysis* (análise de dados da pilha de chamadas) e *call traces* (rastreamento de chamadas). Estes eventos de baixo-nível são proprietários dos sistemas operativos e produzidos involuntariamente pelos utilizadores durante a utilização de diferentes softwares.

Destas tecnologias destacam-se as de *biometric sketch*, *keystroke dynamics*, *mouse dynamics* e *tapping* por serem aquelas que menos tempo de aprendizagem e verificação necessitam. Como iremos ver nas secções seguintes, estas são as técnicas mais implementadas na construção de sistemas de autenticação e de deteção de intrusões.

2.1.1. Audit Logs

A maior parte dos sistemas operativos modernos mantêm alguns registos da atividade do utilizador e das suas interações com os programas. Enquanto esses dados podem ser de algum interesse para a deteção de desvios comportamentais dos utilizadores, sistemas especializados em reforçar a segurança podem fazer um uso muito mais eficiente e poderoso destes dados. Um *audit log* tradicional contém informação como utilização do CPU e I/O, número de conexões para cada local, se um diretório foi acedido, um ficheiro criado, a mudança do ID do utilizador, se o *audit record* foi modificado, quantidade de atividade para o sistema, rede e anfitrião. Com recurso a algumas experiências foi mostrado que recolher dados de auditoria é uma técnica menos intrusiva do que guardar as chamadas ao sistema (*system calls*). Como pode ser criada uma enorme quantidade de dados de auditoria, e que pode sobrecarregar o sistema de deteção de intrusão, tem sido sugerido que recolher uma amostragem aleatória dos dados pode ser a abordagem mais razoável. No entanto, os dados adicionais podem ser importantes para distinguir atividade suspeita e que se distâncie do comportamento normal. Por exemplo, mudanças de estado dos utilizadores, adição de utilizadores novos, mudança de permissões ou a mudança de atribuições aos utilizadores podem ser dados importantes para detetar uma intrusão e assim despoletar um alerta. Como com esta técnica pode ser capturada uma quantidade de informação muito valiosa (ver tabela 6 retirada de (Lunt, 1993)), ela foi alvo de estudo em (Lunt, 1993), (Wespi, Dacier, & Debar, 2000) e (Lee, Stolfo, & Wok, 1999). A tabela 6 (retirada de (Lunt, 1993)) mostra as informações que se podem obter com esta técnica.

Tabela 6 Dados possíveis de obter usando Audit Logs

MEASURE	DESCRIPTION
CPU Usage (ordinal)	CPU usage
I/O Usage (ordinal)	I/O usage
Location of Use (linear categorical)	# of connections from each location
Mailer Usage (linear categorical)	# of times each mailer was used
Editor Usage (linear categorical)	# of times each editor was used
Compiler Usage (linear categorical)	# of times each compiler was used
Shell Usage (linear categorical)	# of times each shell was invoked
Window Command Usage (linear categorical)	# of times each window command was used
Program Usage (linear categorical)	# of times each program was used
System Call Usage (linear categorical)	# of times each system call was used
Directory Usage (linear categorical)	# of times each directory was accessed
Directory Usage (binary categorical)	Whether a directory was accessed
Commands Used (ordinal)	# of different commands invoked
Directories Created (ordinal)	# of directories created
Directories Deleted (ordinal)	# of directories deleted
Directories Read (ordinal)	# of directories read/accessed
Directories Modified (ordinal)	# of directories modified
File Usage (linear categorical)	# of times each file was accessed
File Usage (binary categorical)	Whether a file was accessed
Temp File Usage (ordinal)	# of temporary files accessed
Files Created (ordinal)	# of files created
Files Deleted (ordinal)	# of files deleted
Files Read (ordinal)	# of files read/accessed
Files Modified (ordinal)	# of files modified
User IDs Accessed (linear categorical)	# of times user ID was changed
User IDs Accessed (binary categorical)	Whether another user ID was accessed
System Errors (ordinal)	# of system-related errors
System Errors by Type (linear categorical)	# of times each type of error occurred
Audit Record Activity (linear categorical)	# of audit records for each hour
Hourly Activity (binary categorical)	Whether an audit record was rec'd for each hour
Day of Use (linear categorical)	# of audit records for each day
Day of Use (binary categorical)	Whether audit records were received for each day
Remote Network Activity (ordinal)	Amt. of remote network activity
Network Activity by Type (linear categorical)	Amt. of network activity of each type
Network Activity by Hosts (linear categorical)	Amt. of network activity for each remote host
Local Network Activity (ordinal)	Amt. of network activity within the local system
Local Network Activity by Type (linear categorical)	Amt. of local network activity of each type
Local Network Activity by Hosts (linear categorical)	Amt. of local network activity for each host

2.1.2. Biometric Sketch

Este método de autenticação é proposto em (Al-Zubi, Bromme, & Tonnies, 2003) e (Bromme & Al-Zubi, 2003) e baseia-se no reconhecimento da capacidade de desenho do utilizador e do conhecimento do utilizador sobre o conteúdo de desenhos. O sistema pede ao utilizador que faça um desenho simples (por exemplo três círculos, onde cada utilizador pode desenhá-los como quiser). Como existem um grande número de combinações diferentes para combinar múltiplas formas estruturais simples, esboços de diferentes utilizadores são suficientemente únicos para

fornecer uma autenticação precisa. A abordagem mede o nível de conhecimento do utilizador sobre o esboço, que só está disponível para um utilizador previamente autenticado. Características como a posição do esboço e posição relativa das diferentes primitivas são a base para a criação do perfil do utilizador.

A tabela 7 (adaptada de (Al-Zubi, Bromme, & Tonnie, 2003)) mostra a potencialidade desta técnica, onde podemos ver 10 utilizadores a fazerem quatro tarefas iguais mas de forma completamente distinta.

Tabela 7 Exemplo da utilização de Biometric Sketch

task	description		objects
1	Draw three connected wheels of different sizes		3
2	Draw 3 connected bars one bar is bigger than the others Connect the bars to 3 knots		6
3	Draw 2 connected wheels one wheel is bigger than the other Connect the wheels to a small bar Connect bar to a big base		4
4	draw Task 2 and task 3 connect them with a knot		11

	user 1	user 2	user 3	user 4	user 5	user 6	user 7	user 8	user 9	user 10
task 1										
task 2										
task 3										
task 4										

2.1.3. Call Stack

Um sistema deste tipo, apresentado em (Feng, Kolesnikov, Fogla, Lee, & Gong, 2003), é utilizado para detetar anomalias utilizando a informação da pilha de chamadas. O *program counter* indica o ponto de execução atual de um programa. Uma vez que a cada instrução de um programa corresponde a um único *program counter*, esta informação é útil para a deteção de intrusões. A ideia passa por extrair endereços de retorno da pilha de chamadas e gerar um caminho de execução abstrato entre dois pontos de execução do programa. Este caminho é posteriormente analisado para decidir se é válido, baseando-se no que foi aprendido durante a execução normal do programa. Os endereços de retorno são uma fonte particularmente boa de informação sobre

comportamentos suspeitos. Esta abordagem mostrou a capacidade de detetar alguns ataques, que não seriam detetados por outras abordagens, mas mesmo assim tem uma taxa de falsos positivos similar.

2.1.4. E-Mail Behaviour

O comportamento do envio de mensagens de correio eletrónico não é igual para todos os indivíduos. Algumas pessoas trabalham de noite e enviam grandes quantidades de mensagens para muitos diferentes destinos; outros só verificam a caixa de entrada pela manhã e só trocam correspondência com uma ou duas pessoas. Todas estas peculiaridades podem ser usadas para criar um perfil de comportamento que pode ser como biometria para um indivíduo. Tamanho das mensagens, período do dia em que foi enviado, frequência com que a caixa de entrada é esvaziada e claro os endereços dos destinatários, entre outras variáveis podem ser combinados para criar um vetor de referência para o comportamento de uma pessoa no correio eletrónico (Stolfo, Hershkop, Wang, Nimeskern, & Hu, 2003).

Outra variante (de Vel, Anderson, Corney, & Mohay, 2001) é a utilização de técnicas de identificação para determinar o potencial autor de uma mensagem. Juntamente com as características típicas utilizadas na identificação de textos, também foi usado algumas características únicas das mensagens de correio eletrónico, tais como: uso de uma saudação, assinatura, número de anexos, posição do texto citado no corpo da mensagem, frequência e número total de etiquetas HTML. Ao todo podem ser utilizadas 200 características, mas algumas não são adequadas para o uso em mensagens de correio eletrónico devido ao seu tamanho curto.

Contudo o estudo da utilização do correio eletrónico pode ser considerado intrusivo devido à quantidade de informação pessoal que pode ser adquirida (Yampolskiy & Govindaraju, 2008).

2.1.5. GUI Interaction

Grande parte dos sistemas de deteção de intrusões estão melhor preparados para ameaças externas do que internas. Após a autenticação de um utilizador da sua máquina, este pode deixá-la só e ficar vulnerável a que qualquer intruso a use para seu próprio proveito. Para resolver esse

problema e partindo do facto de que para realizar uma tarefa existem várias hipóteses de a fazer na interface (diferentes atalhos, caminhos, utilizações do rato, etc.) foi criado um sistema baseado no sistema operativo Windows, (Garg, Rahalkar, Upadhyaya, & Kwiat, 2006), para recolher informação da interação do utilizador com a interface gráfica (GUI). Os dados recolhidos permitem a geração de perfis comportamentais dos utilizadores do sistema. Idealmente, os dados recolhidos incluem informações de alto nível sobre interações do utilizador com a interface como: cliques com o botão esquerdo no menu inicial, duplo-clique, fecho de janelas, entre outros. O *software* produzido recolhe também todas as atividades de baixo-nível do utilizador em tempo real, incluindo: processos background do sistema, execuções de comandos, atividade do teclado e cliques do rato. Assim mesmo que o computador tenha autenticado o seu utilizador legítimo, nenhuma outra pessoa o pode utilizar visto que vão haver diferenças comportamentais na utilização da interface do sistema (Garg, Rahalkar, Upadhyaya, & Kwiat, 2006).

2.1.6. Keystroke Dynamics

Os padrões de digitação no teclado são característicos para cada pessoa: alguns são digitadores experientes, outros utilizam a abordagem “*hunt-and-peck*” utilizando somente dois dedos. Estas diferenças fazem com que a verificação da identidade das pessoas baseada nos seus padrões de digitação seja possível. Para a verificação, uma pequena amostra de digitação (como a inserção da palavra-chave) é suficiente, mas para a identificação de utilizadores é necessária uma grande amostra de dados de digitação que depois são comparados com perfis de utilizadores já presentes no sistema. As características desta técnica são baseadas na latência entre teclas seguidas, duração da pressão na tecla, velocidade de digitação, frequência de erros, uso do numpad, etc.. *Keystroke dynamics* é provavelmente a técnica de biometrias HCI mais estudada: (Revett K. , et al., 2007) (Revett K. , et al., 2006) (Revett, de Magalhães, & Santos, 2005) (Monrose & Rubin, 1999) (Joyce & Gupta, 1990) (Haider, Abbas, & Zaidi, 2000) (Ilonen, 2003).

Keystroke dynamics é uma biometria baseada na suposição de que pessoas diferentes teclam de maneiras únicas. Observações de operadores de telégrafo no século 19 revelaram padrões distintos de escrita de mensagens sobre linhas telegráficas, e os seus operadores conseguiam reconhecer-se entre si baseados nos padrões de escrita. Conceptualmente a correspondência mais próxima entre os sistemas de biometria é a de reconhecimento de

assinaturas. Em ambas, no reconhecimento de assinaturas e no *keystroke dynamics*, a pessoa é identificada através de dinâmicas de escrita que são assumidas como únicas.

Os algoritmos de *keystroke dynamics* monitorizam a entrada do teclado milhares de vezes por segundo com o objetivo de identificar utilizadores com base no seu habitual padrão de escrita. Este método biométrico, ao contrário da maior parte dos outros, é quase gratuito - o único *hardware* necessário é o teclado. Medidas possíveis do *keystroke dynamics*:

- Latência entre teclas consecutivas;
- Duração da pressão na tecla (*hold-time*);
- Velocidade geral de digitação;
- Frequência de erros (número de utilizações das teclas *backspace* e *delete*);
- Hábito de uso de teclas adicionais no teclado (ex.: teclado numérico);
- Ordem em que utilizador pressiona as teclas no caso de letras maiúsculas (primeiro solta a tecla *shift* ou a tecla correspondente à letra);
- Força usada na pressão das teclas (só para teclados especiais);

Os mais populares são a latência entre teclas consecutivas e a duração do *keystroke* pois podem ser facilmente obtidas com o *hardware* de um computador normal. Tanto o carregar como o soltar da tecla geram eventos e interrupções no *hardware* que podem ser obtidos. Reunir os dados de *keystroke dynamics* pode ter, às vezes, algumas complicações. Várias teclas podem ser premidas ao mesmo tempo e em escritas rápidas o utilizador pressiona a tecla seguinte sem soltar a anterior. Dependendo do que é medido pode haver tempo negativo ao medir o tempo entre o soltar de uma tecla e a pressão noutra. Também se pode efetuar tarefas mais complexas como querer saber se o utilizador usa as teclas *shift*, *alt* ou outras teclas especiais.

Um dos problemas da monitorização da escrita em teclados é que esta pode ser afetada pelo nível de alerta do utilizador, ele pode estar ensonado ou doente. Acidentes também podem condicionar a escrita, como por exemplo andar com um penso num dedo, ou então escrever só com uma mão enquanto outra segura no café, entre outras. Mudar de teclado ou de computador

também pode levar a um padrão de escrita tremendamente diferente. Todos estes fatores têm que ser tidos em conta na especificação de um sistema de *keystroke dynamics*:

- Fatores de ambiente – altura da cadeira, distância entre teclado e a pessoa, utilização de um teclado novo, etc.
- Perícia do utilizador – ao longo do tempo o utilizador digita de forma mais eficaz e rápida;
- Estado emocional do utilizador – raiva, desespero, felicidade, nervosismo, excitação, etc.;
- Estado físico do utilizador – fadiga, doença, acidentes, escrever só com uma mão, falar ao telefone, etc.;

2.1.7. Mouse Dynamics

Ao monitorizar todas as ações produzidas pelo rato durante a interação do utilizador com a interface, pode ser gerado um perfil único que pode ser utilizado para fins de autenticação. As ações da utilização do rato que podem ser monitorizadas incluem o seu movimento geral, método de arrastar e soltar, apontar e clicar, e o tempo de imobilidade. A partir destas características podem ser extraídas por exemplo a velocidade média, a distância percorrida e a direção do movimento (Ahmed & Traore, 2007) (Pusara & Brodley, 2004). A tabela 8 (retirada de (Shen, Cai, Guan, Sha, & Du, 2009)) representa as métricas que podem ser obtidas num sistema de *mouse dynamics*:

Estudos nesta área (Gamboa & Fred, 2004) revelam que na interação esquemática a variância é muito elevada (tem que se recorrer a algoritmos para retirar ruídos das amostras) mas que as de habilidade já eram mais estáveis e que eram estas que se deviam ter em conta. Resumindo, esta é uma técnica muito semelhante ao *keystroke dynamics*, sendo que são usadas algumas vezes em conjunto (Ahmed & Traore, 2005).

Tabela 8 Características do Mouse Dynamics

Features of Mouse Dynamics

Category of mouse feature	Mouse Dynamics Features
Schematic Features	Mouse action histogram: statistics of occurrences for various mouse action types
	Percentage of silence periods: statistics of idle time of mouse
	distribution of cursor positions on the screen
	distribution of movement distances/directions
Motor-skill Features	Elapsed time of single click: time interval between down and up of left/right/middle button of a click
	Elapsed times of double click: overall time and 3 internal intervals between downs and ups of left/right/middle button of a double click
	Average movement speed compared to directions: average movement speed calculated for different directions
	Average movement speed and acceleration compared to traveled distance: average speed/accelerations calculated for different distance traveled
	Transition time of actions: transition time between consecutive mouse actions

2.1.8. Network Traffic

A detecção de intrusões ao nível da rede é algo diferente dos outros tipos de detecção de intrusão, já que a atividade monitorizada é gerada fora do sistema que está a ser protegido. Com o aumento da popularidade da internet e outras redes, um intruso já nem tem que ter um endereço físico para o alvo que quer penetrar. Isto significa que o fluxo de dados da rede entram por diferentes portas, estão codificados usando diferentes protocolos, precisando de ser processados e revistos. Sistemas de detecção de intrusão baseados em análise do tráfego da rede analisam vários atributos dos pacotes como o tamanho, número de portas, IP's de origem e destino, valores do tempo de vida, cabeçalhos IP/TCP, *checksums* ou *flags*. Durante o período de construção de perfis, o número de pacotes com cada valor é contado e é tido como um comportamento normal (Zhang & Manikopoulos, 2003) (Novikov, Yampolskiy, & Reznik, 2006). Qualquer desvio do perfil pode despoletar um alarme, informando o administrador da rede que está a existir um ataque.

2.1.9. Registry Access

Em (Apap, Honig, Hershkop, Eskin, & Stolfo, 2002) é proposto um novo tipo de abordagem de segurança chamado “*registry anomaly detection (RAD)*” que monitoriza o acesso ao registo do Windows em tempo real e deteta a ação de software malicioso. O registo do Windows guarda informação sobre o *hardware* instalado no sistema, que portas estão a ser utilizadas, perfis de utilizadores, políticas, nomes de utilizador, palavras-chave e configurações de programas. A maioria dos programas acede a uma parte do conjunto de chaves do registo durante a sua execução normal. Do mesmo modo a maior parte dos utilizadores só usa um conjunto dos programas disponíveis no sistema. Isso resulta num elevado nível de regularidade na interação com o registo durante a execução normal do sistema. No entanto, *software* malicioso altera substancialmente a atividade, fazendo com que possa ser detetado. Muitos ataques envolvem programas que arrancam no início da sessão e que são raramente utilizados posteriormente ou que alteram chaves que nunca foram mexidas antes. Se um sistema RAD for treinado sobre dados “limpos”, ou seja em dados gerados durante a execução normal dos programas, esse tipo de operações sobre o registo vão parecer anormais para o sistema e irá despoletar num alerta.

2.1.10. Storage Activity

O estudo (Stanton, Yurcik, & Brumbaugh, 2005) afirma que grande parte das ações de intrusos são visíveis ao nível da interface de armazenamento. Manipulação dos utilitários do sistema (para a adição de *backdoors*), adulteração dos registos de auditoria (*audit logs*) com a finalidade de apagar o rasto, o *reset* de atributos (para esconder mudanças), e adição de conteúdos suspeitos (vírus) mostram-se todos nas mudanças ao nível da camada de armazenamento do sistema. Um sistema de segurança baseado em armazenamento analisa todas as solicitações recebidas pelo servidor de armazenamento e pode emitir alertas sobre atividades suspeitas. Adicionalmente, também pode abrandar o acesso de um suspeito ou isola-lo através de técnicas específicas. Segurança baseada em armazenamento tem a vantagem de ser independente do sistema operativo e consegue continuar ativa após uma invasão, ao contrário de outros sistemas que podem ser desativados pelo intruso.

2.1.11. System Calls

As chamadas ao sistema são um método usado por um programa para requerer um serviço ao sistema operativo. Estas chamadas usam uma instrução especial que faz com que o processador transfira o controlo para um segmento de código mais privilegiado. A deteção de intrusões pode ser alcançada comparando as chamadas do sistema de uma aplicação com um modelo de chamadas normais, aprendidas durante a execução normal do sistema. O pressuposto é que, desde que o intruso não pode fazer chamadas ao sistema aleatoriamente, é pouco provável que ele possa atingir os seus objetivos maliciosos (Bhatkar, Chaturvedi, & Sekar, 2006) (Giffin, Jha, & Miller, 2004) (Kosoresow & Hofmeyr, 1997).

2.1.12. Tapping

(Henderson, Papakostas, White, & Hartel, 2001) e (Henderson, White, Veldhuis, Hartel, & Shump, 2002) estudam a ideia de reconhecimento do toque, baseado no pressuposto que é possível reconhecer alguém que está a bater numa porta. O estudo concentra-se nas propriedades das ondas dos impulsos que resultam do toque de um sensor de polímero presente numa superfície. Os impulsos de pressão são depois processados para retirar características como: altura e duração do impulso, e duração dos intervalos entre impulsos.

Evoluindo esta técnica vemos que pode ser utilizada em *smartphones*, *tablets* ou em *touchpads* da maioria dos computadores portáteis atuais.

2.1.13. Web Browsing

No que diz respeito a técnicas de navegação Web existem alguns estudos para a construção de sistemas de aprendizagem dos interesses dos utilizadores. Estes sistemas de acordo com o número de ações na página, velocidade do “*scroll down*” e número de cliques em *links* definem o nível de interesse dessa página para o utilizador. A partir destes interesses consegue-se traçar um perfil de utilizador e sugerir páginas de interesse para ele (Goecks & Shavlik, 2000) (Liang & Lai, 2002).

Outro estudo, (Fu & Shih, 2002), refere que através dos dados de utilização da Web se consiga entender os interesses, comportamentos e preferências dos utilizadores. Aqui é guardado dois tipos de atividades: atividades remotas e atividades locais. Nas atividades remotas todos os pedidos feitos aos servidores Web são guardados, estes incluem pesquisas, dados de formulários, *cookies* e endereços das páginas visitadas. No que diz respeito às atividades locais, ações como salvar ou imprimir uma página, fazer Retroceder, Avançar, Atualizar ou Parar, adicionar um favorito, maximizar, minimizar ou fechar a janela são guardados numa base dados. Depois através de técnicas de *data mining* são construídos perfis de utilização para cada utilizador. Num contexto de segurança isto pode ser útil num sistema de deteção de intrusões pois se o comportamento de utilização do navegador Web for diferente do utilizador legítimo pode-se lançar um alerta e prevenir a intrusão.

2.2. Problemas da Análise Comportamental

Como já foi dito anteriormente o processo de analisar e construir perfis comportamentais para posterior autenticação ou identificação de utilizadores levanta alguns problemas como a elevada taxa de falsos positivos, a possibilidade de haver intrusões durante a fase de aprendizagem do algoritmo e, se calhar o mais importante, a grande variância que os comportamentos humanos apresentam.

Os comportamentos humanos variam de forma frequente e gradual ao longo do tempo, o que nos leva a um problema na construção de perfis comportamentais que permitam a verificação e a identificação de utilizadores. A adaptação rápida do sistema às constantes variações de comportamentos é um aspeto muito importante e que se tem de resolver de forma eficaz. A grande dificuldade reside em distinguir mudanças verdadeiras ou ruído. O algoritmo de aprendizagem ideal será aquele que combina robustez ao ruído e sensibilidade às mudanças (Tsymbol, 2004). Existem já alguns algoritmos adaptativos capazes de lidar com esta mudança de comportamento (*concept drift*):

- Em (Schlimmer & Granger, 1986) é apresentado um sistema de aprendizagem incremental, chamado STAGGER, que deteta a mudança de comportamentos dinamicamente. É usada uma representação em grafo onde os nodos representam

atributos booleanos e as conexões pesadas através de um algoritmo bayesiano que associam os nodos dos atributos para um nó do conceito. Depois o sistema aprende e deteta as mudanças de comportamento ao adicionar novos nós de atributo ou ao ajustar o peso das conexões.

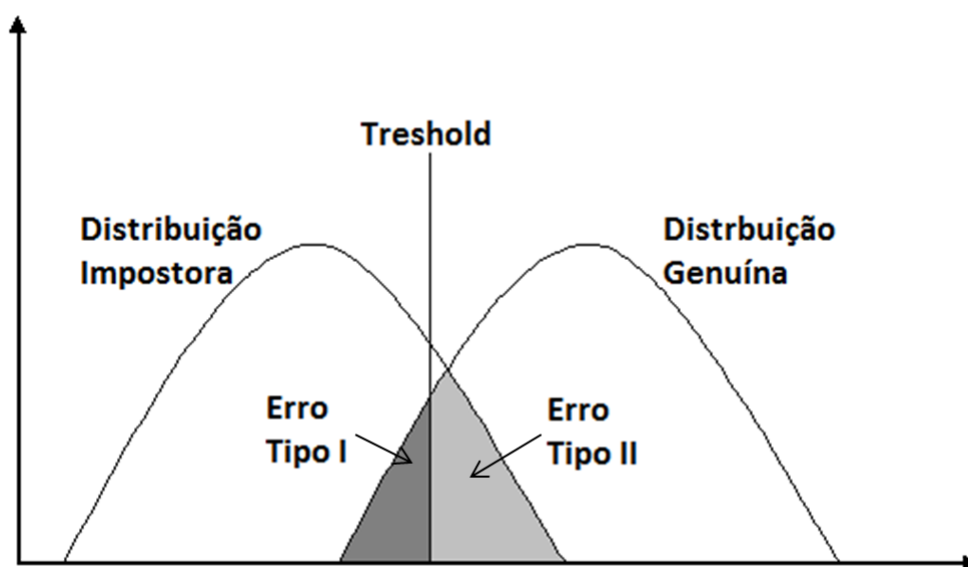
- O sistema FLORA (Widmer & Kubat, 1996) usa uma técnica de esquecimento com um tamanho de janela adaptável. Os exemplos são esquecidos quando são mais velhos que um *threshold*. O tamanho da janela e a taxa de esquecimento são supervisionados e ajustados dinamicamente pela heurística que monitoriza o processo de aprendizagem.
- Outra técnica e semelhante à anterior (Koychev & Schwab, 2000) afirma que o esquecimento natural é um processo gradual, por isso as últimas observações devem ser mais importantes que as mais antigas e a importância de uma observação deve reduzir durante o tempo. Para isso foi construída uma função temporal que é capaz de produzir pesos para cada observação de acordo com a sua idade. Depois o algoritmo treina o sistema utilizando estas observações pesadas.

Em relação à elevada taxa de falsos positivos podemos ver na secção seguinte que já existem algumas técnicas que conseguem ter resultados razoáveis e com baixas taxas de falsos positivos.

2.3. Análise Comparativa

Nesta secção serão apresentados resultados e testes realizados de implementações com as biometrias acima apresentadas. A tabela 9 lista as biometrias com as suas taxas de detenção, taxas de falsos positivos e para as biometrias diretas dois tipos de erros, apresentados na figura 5.

Figura 5 Tipos de erro



- Tipo I – Probabilidade de um utilizador válido ser rejeitado (falso negativo);
- Tipo II – Probabilidade de um utilizador inválido ser aceite (falso positivo);

Tabela 9 Comparação de resultados entre biometrias

Biometria	Fonte	Taxa de detecção	Taxa falsos positivos	Erro tipo I	Erro tipo II
Audit Logs	(Lee, Stolfo, & Wok, 1999)	93%	8%	-	-
Biometric Sketch	(Bromme & Al-Zubi, 2003)	-	-	7.2%	7.2%
E-mail Behaviour	(de Vel, Anderson, Corney, & Mohay, 2001)	90.5 %	-	-	-
GUI Interaction	(Garg, Rahalkar, Upadhyaya, & Kwiat, 2006)	-	-	3.85%	3.85%
	(Monrose & Rubin, 1999)	-	-	7.9%	7.9%
	(Revett K. , et al., 2006)	-	-	4%	4%
Keystroke Dynamics	(Joyce & Gupta, 1990)	-	-	16.36%	0.25%
	(Bergadano, Gunetti, & Picardi, 2002)	-	-	4%	0.01%

	(Shen, Cai, Guan, Sha, & Du, 2009)	-	-	3%	0.55%
Mouse dynamics	(Pusara & Brodley, 2004)	-	.	1.75%	0.43%
Keystroke Dynamics e Mouse Dynamics	(Ahmed & Traore, 2005)	-	1.31%	-	-
	(Novikov, Yampolskiy, & Reznik, 2006)	93.2%	0.8%	-	-
Network Traffic	(Zhang & Manikopoulos, 2003)	96.2%	0.04%	-	-
Registry Access	(Apap, Honig, Hershkop, Eskin, & Stolfo, 2002)	86.9%	3.8%	-	-
Storage Activity	(Stanton, Yurcik, & Brumbaugh, 2005)	97%	4%	-	-
Tapping	(Henderson, Papakostas, White, & Hartel, 2001)	-	-	2.3%	2.3%

Analisando a tabela 9 podemos concluir que já existem sistemas de autenticação e de deteção de intrusões com resultados aceitáveis. Os valores da taxa de deteção rondam os 90% e os da taxa de falsos positivos raramente ultrapassam os 4%. No que diz respeito a biometrias diretas estas também têm boas perspetivas mas os resultados ainda são um pouco flutuantes pois ao construir um sistema mais rígido a taxa de falsos negativos (erro tipo I) aumenta mas o de falsos positivos diminui (erro de tipo II) e vice-versa. A chave está em regular o *threshold* do algoritmo para que este possa garantir o resultado mais aceitável.

Na seguinte secção apresentam-se algumas abordagens e aplicações contendo as metodologias mencionadas.

3. Projetos Relacionados

A área da autenticação e da deteção de intrusões com base em biometrias comportamentais é já um alvo de investigações há alguns anos. Porém, nunca resultaram em produtos públicos devido aos problemas apresentados na secção 2.2. Atualmente, e com o aparecimento de técnicas cada vez mais inovadoras, já existem, embora em pouco número, produtos comerciais nesta área. Convém notar que as soluções existentes no mercado são todas baseadas nas tecnologias *keystroke dynamics*, *mouse dynamics*, *GUI Interaction* e *Biometric Sketch*. A razão para a adoção destas tecnologias está no facto de serem independentes do sistema operativo, terem tempos baixos de aprendizagem e verificação e podem ser utilizadas em ambientes Web. Em seguida serão apresentados algumas dos projetos mais relevantes.

3.1. Keytrac

A ferramenta Keytrac (Keytrac, 2013), desenvolvida pela empresa alemã TM3 Software, consiste na autenticação de utilizadores com base nos seus padrões de escrita (*keystroke dynamics*). Um utilizador possui o tradicional par de credenciais (nome de utilizador e palavra-chave), e no momento de introdução dessas credenciais o seu comportamento de escrita está a ser monitorizado. Depois toda a informação é enviada para os servidores da Keytrac que tratam da classificação desta amostra. Para concluir o processo, esses servidores respondem ao cliente com a pontuação da amostra e conseqüente autenticação ou rejeição.

Com a grande vantagem de suportar grande parte das linguagens de programação e de ser fácil de integrar, esta ferramenta tem, ainda, algumas lacunas. O facto de toda a informação ter de ser processada nos servidores da Keytrac, constitui um grande problema para clientes que não tenham ou não queiram estar ligados à Internet. Uma simples falha na rede ou quebra da ligação inutiliza todo o sistema de autenticação fazendo com que os utilizadores não consigam aceder aos recursos que precisam.

Outro problema consiste na inexistência de uma fase de aprendizagem, o que leva com que o sistema, numa fase inicial, não adquira todos os comportamentos distintos de um só utilizador.

Ou seja, numa fase inicial da utilização desta ferramenta a taxa de rejeição de utilizadores válidos pode ser elevada.

3.2. BioSig-ID

O *software* BioSig-ID (Biometric Signature ID, 2013), criado pela empresa norte-americana Biometric Signature ID, substitui os tradicionais portais de início de sessão por uma grelha de desenho. Com a utilização das tecnologias *mouse dynamics* e *biometric sketch* o utilizador é autenticado com base na sua capacidade de desenhar quatro letras ou números. Numa fase de aprendizagem o utilizador desenha o seu código três vezes (ver Figura 6). Na fase de autenticação, este é autorizado a aceder ao sistema caso o seu código for igual e a sua forma de o desenhar coincidir com o perfil construído na fase de aprendizagem.

Figura 6 Utilização do BioSig-ID



Este sistema parte da vantagem que existem várias formas diferentes de desenhar uma figura, o que torna a “assinatura” do utilizador única, fazendo com que seja muito difícil o acesso de intrusos através deste método de autenticação. No entanto, a grande variância que os utilizadores apresentam na interação com o rato constitui a grande desvantagem deste sistema. Outro fator a ter em conta é a habilidade do utilizador a desenhar formas no computador. Esta habilidade pode fazer com que a fase de aprendizagem do sistema não seja suficiente para cobrir todos os tipos de comportamento do utilizador.

3.3. CVMetrics

Desenvolvido pelos norte-americanos Intensity Analytics, este *software* tem duas componentes de autenticação: uma estática e uma contínua. Utilizando a tecnologia *keystroke dynamics* o sistema oferece proteção na fase de início de sessão, e posteriormente funciona de forma contínua ao monitorizar o estado do teclado constantemente.

O sistema de autenticação estática é semelhante ao apresentado pela Keytrac, embora aqui exista uma fase de aprendizagem. Nesta fase de aprendizagem o utilizador insere a sua palavra-chave 10 vezes. De notar que só a palavra-chave é monitorizada, por isso tem que ter uma dimensão mínima de 15 caracteres. Uma vez ultrapassada a fase de aprendizagem o utilizador é autenticado pela comparação da amostra com o seu perfil de escrita da palavra. Embora tenha apresentado resultados apelativos, este sistema peca pela não monitorização da escrita do nome do utilizador forçando a que este utilize uma palavra-chave muito extensa. Essa dimensão pode ser um grave problema, pois aumenta as probabilidades de um utilizador esquecer, fazendo com que a autenticação não seja possível.

Em relação ao sistema contínuo, este exige que na fase de aprendizagem o utilizador insira 10 vezes uma frase ou que escreva um texto livre com um mínimo de 1000 caracteres. Para finalizar esta fase é necessário que outra pessoa introduza o mesmo texto de modo ao sistema poder ter um ponto de comparação para futuras classificações. Esta abordagem apresenta resultados muito positivos em relação às concorrentes, mas tem como principal inconveniente a necessidade da existência de um “impostor” durante a fase de aprendizagem.

3.4. TypeWATCH

Este produto, desenvolvido pela empresa portuguesa Watchful Software, é um detetor de intrusões baseado em comportamento. Utiliza também a tecnologia *keystroke dynamics* e um utilizador é continuamente identificado através da forma como interage com o teclado. A aplicação corre em segundo plano, sendo invisível ao utilizador, e acompanha as suas mudanças de comportamento.

Quando um comportamento errático é detetado, a aplicação despoleta uma série de alertas e ações (como o bloqueio do computador, notificação ao administrador, etc.) para assim minimizar os potenciais danos infringidos por um intruso. Embora não sejam conhecidos os resultados e a eficiência da solução, esta apresenta uma grande lacuna ao só monitorizar o estado do teclado. Desta maneira, todo o tipo de intrusões que não use como fonte de interação o teclado passa despercebida ao sistema, ficando este vulnerável à maioria dos ataques de roubo de identidade.

3.5. Síntese

As quatro soluções em cima apresentadas são muito similares e utilizam maioritariamente a técnica *keystroke dynamics*. Com a grande vantagem de um sistema deste tipo ser relativamente fácil de desenvolver e manter, também apresenta bons resultados no que diz respeito à autenticação e deteção de intrusões.

Porém, nem todas as soluções cobrem o problema de forma total e/ou apresentam arquiteturas que podem apresentar alguns problemas aos utilizadores. O fato de ser necessário uma constante ligação ao servidor do fabricante ou da necessidade de haver uma pessoa externa na fase de aprendizagem dos algoritmos constituem algumas desvantagens da utilização destas técnicas.

Nesta dissertação será construído um protótipo de sistemas de autenticação entre indivíduos e sistemas informáticos, utilizando biometrias comportamentais que sejam independentes de sistemas operativos e que utilizem recursos e *hardware* já existentes num sistema informático tradicional. Assim optou-se pela utilização das técnicas de *keystroke dynamics* e *mouse dynamics* pois unicamente necessitam do rato e teclado, equipamentos onde é feita toda a interação entre um indivíduo e um sistema. Outra especificidade passa pela necessidade de os sistemas serem independentes, sem a necessidade de ligações ao exterior. Pela consulta da tabela 5 também podemos constatar que estas técnicas são as que precisam de menos tempo de aprendizagem/verificação, o que é importante para a deteção de ataques curtos e esporádicos. Na tabela 9 vemos que a eficiência e fiabilidade de sistemas deste tipo são atingem já um patamar superior em relação a todas as outras técnicas aqui abordadas.

4. ALGORITMOS DE CLASSIFICAÇÃO

A escolha do algoritmo de classificação adequado é um dos fatores mais importantes num sistema de deteção de intrusões. Este algoritmo é responsável por classificar as sessões em relação ao perfil do utilizador. A classificação resultante da execução deste algoritmo determina se o autor da sessão é ou não fidedigno. Em (Killourhy & Maxion, 2009) é efetuado um estudo comparativo utilizando 14 algoritmos distintos. Estes algoritmos são utilizados num problema de *keystroke dynamics* semelhante ao desenvolvido neste trabalho. Na tabela 10, retirada de (Killourhy & Maxion, 2009), estão apresentados os resultados obtidos.

Tabela 10 Comparação de algoritmos de classificação

Detector	equal-error rate
1 Manhattan (scaled)	0.096 (0.069)
2 Nearest Neighbor (Mahalanobis)	0.100 (0.064)
3 Outlier Count (z-score)	0.102 (0.077)
4 SVM (one-class)	0.102 (0.065)
5 Mahalanobis	0.110 (0.065)
6 Mahalanobis (normed)	0.110 (0.065)
7 Manhattan (filter)	0.136 (0.083)
8 Manhattan	0.153 (0.092)
9 Neural Network (auto-assoc)	0.161 (0.080)
10 Euclidean	0.171 (0.095)
11 Euclidean (normed)	0.215 (0.119)
12 Fuzzy Logic	0.221 (0.105)
13 <i>k</i> Means	0.372 (0.139)
14 Neural Network (standard)	0.828 (0.148)

Os algoritmos foram classificados segundo a *equal-error rate* ou taxa de erro igual que significa uma taxa igual para erros do tipo I e tipo II (ver secção 2.3). Dentro de parêntesis são apresentados os desvios padrão dos resultados. Com base nestes resultados e na complexidade dos algoritmos foram implementados neste trabalho os algoritmos *Outlier Count* e *Nearest Neighbor (Mahalanobis)*.

4.1. Outlier Count

Este algoritmo baseia-se na contagem dos parâmetros de sessão que quando comparados com a matriz de perfil do utilizador são considerados atípicos ou distantes. Quanto maior o número de parâmetros atípicos, ou seja não válidos, menor a probabilidade da sessão a avaliar seja do

utilizador fidedigno. Na prática o algoritmo funciona de forma contrária, é contado o número de parâmetros válidos e se esse número ultrapassar um limiar (*threshold*) então estamos diante de uma sessão válida. Os parâmetros são avaliados segundo a função 1.

$$O(x_i) = \begin{cases} 1 & \text{se } (\mu_i - \sigma_i) \leq x_i \leq (\mu_i + \sigma_i) \\ 0 & \text{nos outros casos} \end{cases} \quad (1)$$

Em que μ_i e σ_i correspondem à média e desvio padrão da matriz de perfil para o parâmetro i respetivamente, e x_i corresponde ao valor da sessão para o parâmetro i . n diz respeito ao número de parâmetros presentes na sessão. A classificação final é calculada nas funções 2 e 3.

$$\text{Classificação}_{outlier} = \frac{\sum_{i=1}^n O(x_i)}{n} \quad (2)$$

$$\text{Se } \text{Classificação}_{outlier} > \text{threshold}_{outlier} \rightarrow \text{Sessão Válida} \quad (3)$$

4.2. Nearest Neighbor (Mahalanobis)

A distância de Mahalanobis (McLachlan, 1999) (Maesschalck, Jouan-Rimbaud, & Massart, 2000) é já um algoritmo clássico de deteção de anomalias. Este pode ser visto como uma extensão da distância de Euclides com a vantagem de contar com a correlação entre os diversos parâmetros. Em termos práticos é calculada a matriz de covariância da matriz de perfil de utilizador (função 4).

$$C = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^t \quad (4)$$

Onde n corresponde ao número de sessões presente na matriz, X_i o valor da sessão para o parâmetro i e \bar{X} a média da matriz.

Depois de calculada a matriz de covariância é calculada a distância da sessão a avaliar para cada uma das sessões presentes na matriz de perfil. O resultado final será a menor distância

calculada, daí o nome *Nearest Neighbor* (vizinho mais próximo). Este cálculo é traduzido nas funções 5,6 e 7.

$$D(x, X) = \sqrt{(x - X)^T C^{-1} (x - X)} \quad (5)$$

Onde x é a sessão a classificar e X a sessão da matriz de perfil. No final a sessão é considerada fidedigna se a distância para a sessão de perfil mais próxima não ultrapassar um limiar (*threshold*) previamente definido.

$$Classificação_{mahalanobis} = \min(D(x, X)) \quad (6)$$

$$Se Classificação_{mahalanobis} < threshold_{mahalanobis} \rightarrow Sessão Válida \quad (7)$$

4.3. Algoritmo Combinado

O algoritmo combinado é uma junção dos dois algoritmos apresentados anteriormente. Com esta combinação esperam-se melhores resultados do que utilizando os algoritmos tradicionais individualmente. A sua classificação é realizada segundo a função 8.

$$Se \frac{Classificação_{mahalanobis}}{Classificação_{outlier}} < \frac{threshold_{mahalanobis}}{threshold_{outlier}} \rightarrow Sessão Válida \quad (8)$$

4.4. Síntese

Neste capítulo foram apresentados 3 algoritmos de classificação capazes de avaliar desvios comportamentais e diferenças entre perfis de utilizadores.

Estes algoritmos foram escolhidos com base no estudo feito por (Killourhy & Maxion, 2009). Neste estudo os algoritmos são avaliados em relação à sua eficácia, complexidade e eficiência. Com base nessas avaliações foram escolhidos os algoritmos *Outlier Count* e *Mahalanobis*.

Por fim é apresentado um algoritmo que combina os dois apresentados anteriormente. Tendo em conta que os algoritmos *Outlier Count* e *Mahalanobis* têm bons comportamentos quando

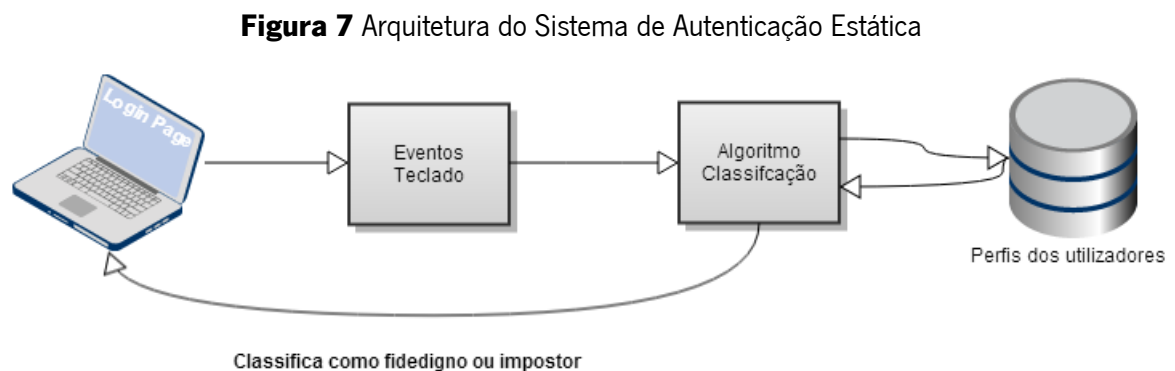
aplicados a situações diferentes, espera-se que este algoritmo combinado seja mais completo, obtendo melhores resultados em todas as situações.

Estes 3 algoritmos serão utilizados pelos 2 sistemas explicados nos próximos capítulos. No final irá ser feita uma análise comparativa, de forma a perceber qual o algoritmo mais indicado para este projeto.

5. SISTEMA DE AUTENTICAÇÃO ESTÁTICA

Nesta secção será apresentado o sistema de autenticação estática, desenvolvido com a tecnologia *keystroke dynamics*. Este é um sistema de autenticação tradicional, onde é inserido o nome de utilizador e a respetiva palavra-chave. A grande diferença reside que o utilizador só é autenticado caso as suas credenciais estejam corretas e o seu padrão de escrita seja igual ao do seu perfil, anteriormente construído. Através da utilização de *keystroke dynamics*, pretende-se verificar através dos padrões de escrita se a pessoa que se está a autenticar é realmente a proprietária da conta que está a tentar aceder ao sistema.

5.1. Arquitetura do Sistema



O sistema, apresentado na figura 7, é composto por três componentes:

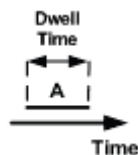
- Interface Web: página que contém um *plugin javascript* onde são capturados todos os eventos de premir e largar uma tecla do teclado;
- Servidor: Recurso que classifica o utilizador e fornece a autenticação de acordo com essa classificação. Este servidor comunica com a interface web utilizando um serviço REST;
- Base de dados: Contém todos os perfis e sessões do sistema.

5.2. Aplicação

O sistema de autenticação, como já dito anteriormente, é baseado numa página de início de sessão tradicional, ou seja, um formulário onde o utilizador insere o seu nome de utilizador e palavra-chave. Porém, neste caso, os eventos do teclado (tecla premida e tecla largada) são monitorizados para depois através de algumas métricas construir o perfil de comportamento do utilizador. Na fase de recolha de dados só os eventos de teclas que produzem texto é que são capturados. Todos os outros (teclas de funções, teclas especiais, etc.) são ignorados. Depois de recolhidos os dados, estes são processados seguindo as seguintes métricas:

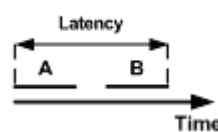
- Tempo de pressão (*dwell time*): tempo que o utilizador preme uma determinada tecla (figura 8);

Figura 8 Tempo de Pressão



- Latência entre duas teclas consecutivas (dígrafo): tempo entre premir uma tecla e largar a tecla seguinte (figura 9);

Figura 9 Latência



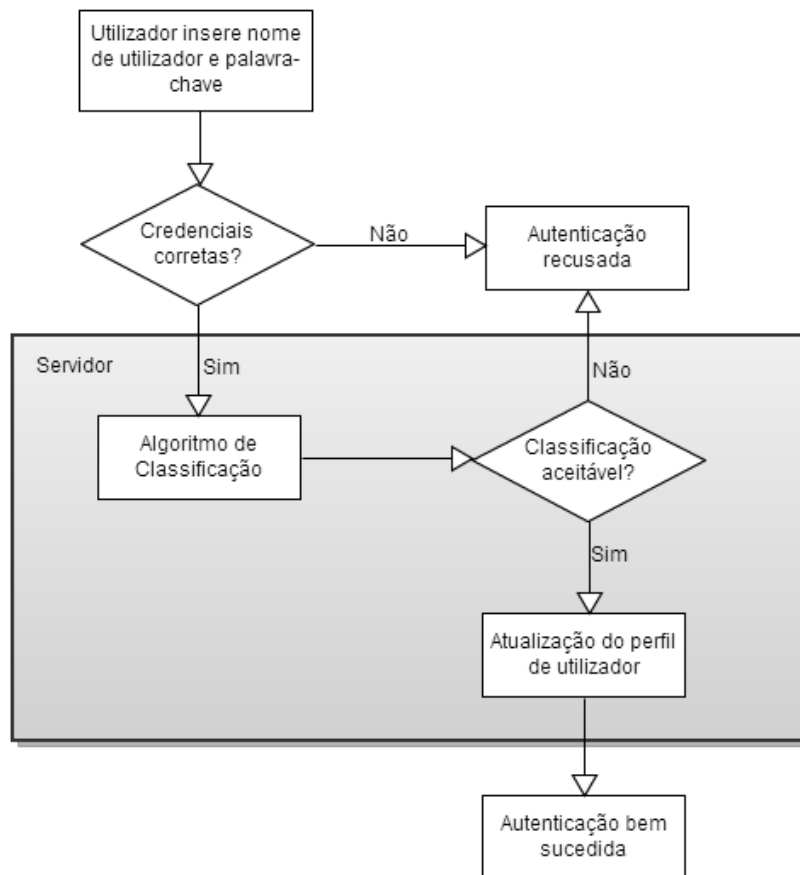
- Velocidade total de escrita do nome de utilizador e da palavra-chave (latência entre a primeira e a última tecla);
- Número de erros (uso da tecla *backspace* ou *delete*);
- Método de transição entre caixas de texto (tecla *Tab* ou utilização do rato);
- Método para escrita de letras maiúsculas (tecla *shift* ou *caps lock*);
- Utilização ou não de *copy & paste*;
- Dispositivo de inserção de texto (*touchscreen* ou teclado).

O perfil do utilizador é construído através das quinze primeiras sessões. A partir daí o algoritmo entra em modo de deteção e o utilizador só consegue a autenticação caso a sessão respetiva tenha uma classificação que garanta a sua identidade.

De modo a que o perfil acompanhe a evolução do comportamento do utilizador é utilizado um algoritmo adaptativo de esquecimento, em que sempre que houver uma autenticação correta, esta contribuirá para o perfil e a autenticação mais antiga será esquecida. Por outras palavras o perfil do utilizador é construído utilizando as últimas quinze sessões em que a autenticação foi bem-sucedida.

A figura 10 apresenta o modo de funcionamento deste sistema.

Figura 10 Funcionamento Sistema Autenticação Estática



5.3. Síntese

Com a utilização deste sistema pretende-se garantir que a pessoa que insere as credenciais é de facto o utilizador correspondente. Assim, o problema do roubo e partilha de credenciais é minimizado já que as pessoas são autenticadas pela forma como escrevem.

A utilização de *keystroke dynamics* numa área sensível como esta vem-se revelando muito útil partindo do princípio que o comportamento de escrita é único, não pode ser partilhado de pessoa em pessoa, e é muito fiável, não pode ser esquecido ou perdido. Outra das grandes vantagens desta tecnologia é a sua ausência de custos e de qualquer mudança de procedimentos. Um teclado e uma tradicional página de início de sessão são os únicos requisitos necessários para o funcionamento de um sistema deste tipo.

A escolha de um sistema cliente/servidor foi óbvia para este caso. Quer pelos exemplos da literatura quer pela sua eficiência no que diz respeito a tempo de computação. O fato do cliente só capturar os eventos do teclado e de todo o processamento e classificação da informação ser feita pelo servidor, faz com que o utilizador não tenha qualquer custo, além do tempo da espera pelos resultados ser bastante reduzido. A utilização de *javascript* para a recolha dos eventos é justificada pelo seu desempenho e por ser compatível com a grande maioria dos navegadores Web existentes.

Embora a aplicação funcione em ambientes *touchscreen*, devido à grande variedade de produtos que existem com esta tecnologia fazem com que os resultados não sejam fiáveis para todos os dispositivos que utilizem o toque como método de introdução de texto.

6. SISTEMA DE AUTENTICAÇÃO CONTÍNUA

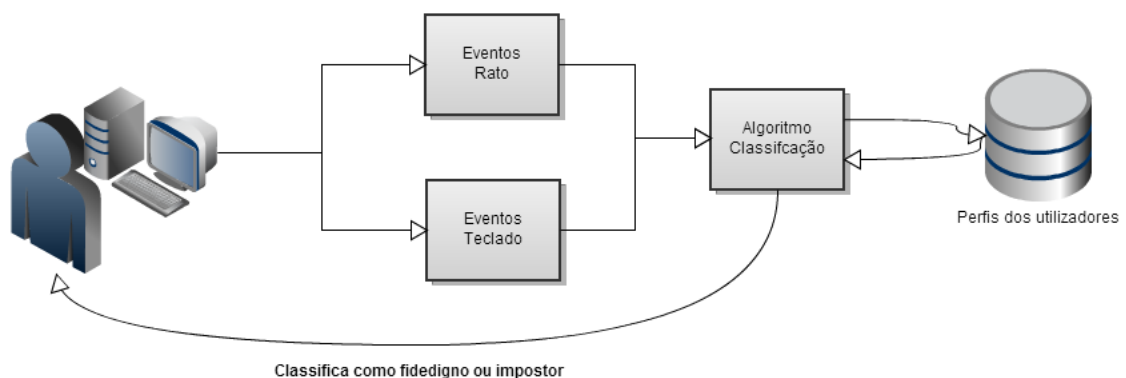
Este sistema foi proposto para combater possíveis intrusões a computadores com sessões de autenticação ativas, ou seja só entrará em funcionamento na fase posterior ao início de sessão. O seu principal objetivo é, com base na interação do utilizador com o rato e teclado, construir um perfil de comportamento que depois será usado para detetar se a pessoa que está a usar um determinado computador é fidedigna ou um intruso.

O perfil de comportamento de cada utilizador é construído com base na sua interação com dispositivos como o teclado e rato. Com o desenvolvimento de algumas métricas, apresentadas nas próximas secções, torna-se possível quantificar e classificar o comportamento e posteriormente concluir se o sujeito que está a operar o computador é de fato aquele que se autenticou.

6.1. Arquitetura do Sistema

Esta solução resume-se a uma aplicação nativa desenvolvida para ambiente Windows (.NET Framework 4.5) que monitoriza constantemente o estado do teclado e rato e partir daí verifica a identidade do utilizador.

Figura 11 Arquitetura do Sistema de Identificação



A figura 11 apresenta uma visão geral à aplicação. Esta recolhe os eventos produzidos pelo rato e teclado no sistema operativo. Devido à enorme quantidade de dados agora recolhidos, estes são submetidos a um módulo de limpeza que filtra e identifica informação necessária. No fim, e analogamente ao sistema anterior, os dados são introduzidos nos algoritmos de classificação, que identificam se o autor dos eventos foi, ou não, o utilizador autorizado.

A aplicação está dividida em dois módulos distintos: uma que trata da monitorização do teclado e outra da monitorização do rato. Embora estas sejam distintas elas podem cooperar de modo a fundir os seus resultados para haver uma maior certeza na verificação da identidade do utilizador.

6.2. Monitorização do teclado

O teclado é monitorizado de forma semelhante à apresentada no sistema anterior (o sistema de recolha e filtragem de eventos é igual) mas o seu processamento é diferente de modo a poder suportar e classificar texto contínuo e livre. Desta forma as métricas utilizadas são:

- Tempo de pressão (herdado do sistema anterior);
- Latência de uma palavra: são medidas todas as latências das combinações das n letras da palavra (n-grafo).

A tabela 11 apresenta um exemplo do cálculo da latência para a palavra “TESE”. De modo a calcular os diversos grafos da palavra são capturadas as *timestamp* dos eventos de pressão e libertação de cada tecla, e a partir daí são calculadas as latências (em milissegundos) das letras da palavra.

Tabela 11 Exemplo de N-Grafo

Tecla	Keystrokes		Dígrafo		Tri-grafo		Tetra-grafo	
	Premida	Largada	Grafo	Latência	Grafo	Latência	Grafo	Latência
T	798340	798409	T + E	187	T + E + S	358	T + E + S + E	510
E	798403	798527	E + S	295	E + S + E	447		
S	798605	798698	S + E	245				
E	708746	708850						

A fase de aprendizagem do algoritmo corresponde às quinze primeiras sessões do utilizador no sistema. Depois é aplicado o algoritmo de esquecimento apresentado no sistema anterior. Cada sessão é constituída por 250 caracteres.

6.3. Monitorização do rato

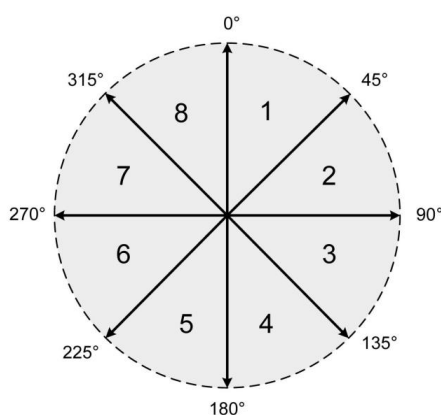
A monitorização do rato é feita através do registo de todos os eventos de movimento e cliques. A aplicação é baseada na descrita em (Ahmed & Traore, 2007) e monitoriza o estado do rato de 100 em 100 milissegundos capturando as seguintes ações:

- *Silence*: ausência total de movimento;
- *Mouse Movement*: movimento do rato;
- *Point and Click*: movimento seguido de um clique ou duplo clique;
- *Drag and Drop*: Pressão de um botão seguido de movimento seguido de largada do botão.

Devido às diversas aplicações presentes no computador, e tendo a cada uma delas associado um tipo de comportamento do rato diferente, as ações são guardadas consoante as aplicações onde foram efetuadas. Assim em vez de compararmos o comportamento geral, podemos classificar o comportamento do utilizador para cada uma das aplicações que usa obtendo assim uma informação mais precisa sobre a sua identidade.

Além das ações acima referidas também é registado o tempo que o utilizador mantém cada botão premido e tempo de duplo clique.

Figura 12 Direção do movimento



O movimento do rato é dividido em 8 direções, seguindo o esquema da figura 12. Onde qualquer movimento entre os 0 e os 45 graus corresponde à direção 1 e assim sucessivamente.

Como exemplo num ambiente Windows, temos que o movimento na direção 4 pode significar a consulta do relógio e na direção 6 a navegação no menu iniciar.

De modo a reduzir a enorme quantidade de dados recolhidos, foi usado a fórmula de amplitude interquartil (NIST/SEMATECH, 2012) que deteta e elimina os valores atípicos (*outliers*) presentes nos dados. Estes *outliers* podem ser erros de leitura, erros de interpretação ou comportamentos dispersos que não traduzem qualquer significado específico. Após a utilização deste algoritmo verificou-se uma eliminação de cerca de 25% dos dados o que aumentou a eficiência e a precisão do sistema.

Depois de recolhidos os dados da sessão e para uma melhor interpretação dos valores obtidos são construídos gráficos que servem como métrica de comparação. Cada sessão do rato é composta por 34 pontos retirados das 5 métricas aqui classificadas:

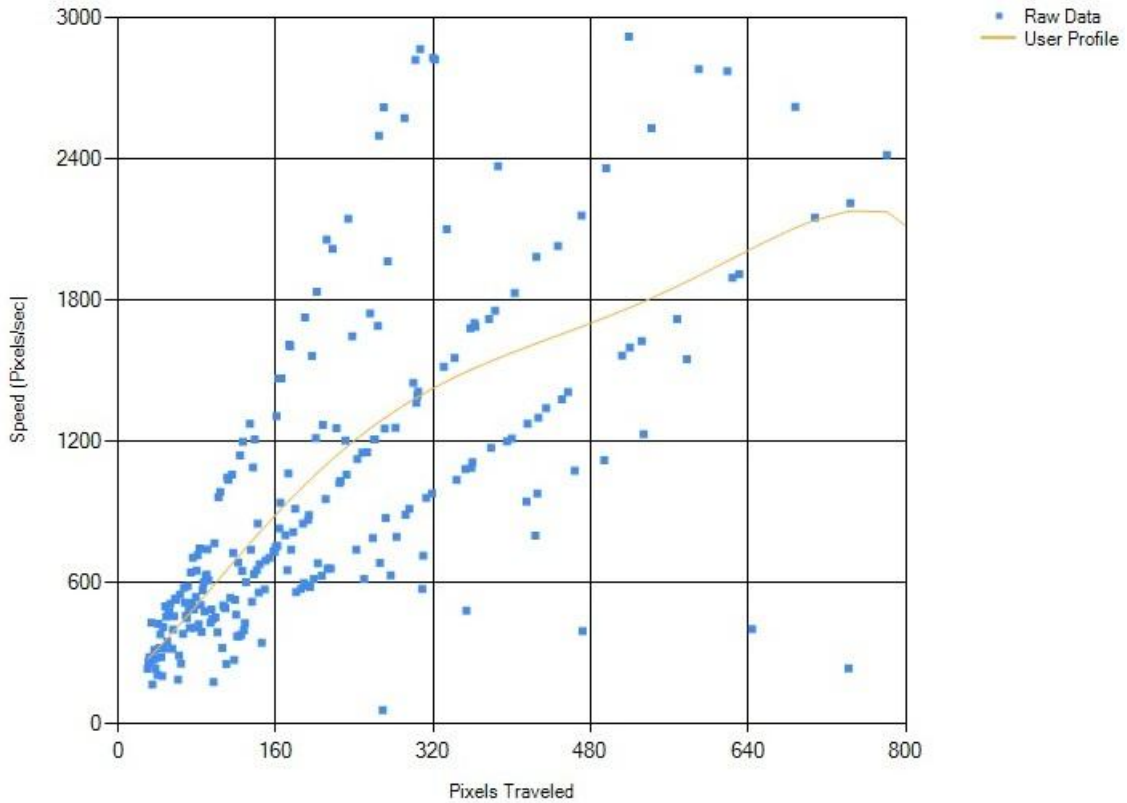
- 12 pontos da velocidade de movimento em comparação com a distância;
- 8 pontos da percentagem de movimento por direção;
- 8 pontos da velocidade média por direção;
- 3 pontos da percentagem do tipo de ação;
- 3 pontos da velocidade média por tipo de ação.

Cada sessão contém 1000 ações do rato. As métricas são explicadas de seguida.

Velocidade de movimento em comparação com a distância

Com a medição da velocidade para cada ação do rato proferida pelo utilizador, é possível traçar o seu perfil de interação. Assim conseguimos saber se o utilizador é rápido a percorrer grandes distâncias, ou se é mais cuidadoso quando tem que fazer pouco movimento. A partir destas medições conseguimos extrair um perfil de utilização normal do rato, que pode ser utilizada para potenciais deteções de intrusão.

Figura 13 Distância do movimento e a sua Velocidade



Na figura 13 os pontos a azul correspondem aos dados brutos recolhidos durante a sessão. Como estes pontos são dispersos e de difícil interpretação, estes são aproximados através de um curva polinomial (curva amarela) utilizando um algoritmo de regressão polinomial de grau 15 (NIST/SEMATECH, 2012). Para a construção da matriz do perfil comportamental do utilizador são considerados 12 pontos sobre a curva. Estes pontos (x_i) são calculados através das fórmulas 9, 10, 11e 12 apresentadas de seguida. Convém notar que $f(x)$ corresponde à curva polinomial de grau 15.

$$x_0 = x_{min} \quad (9)$$

$$y_0 = f(x_0) \quad (10)$$

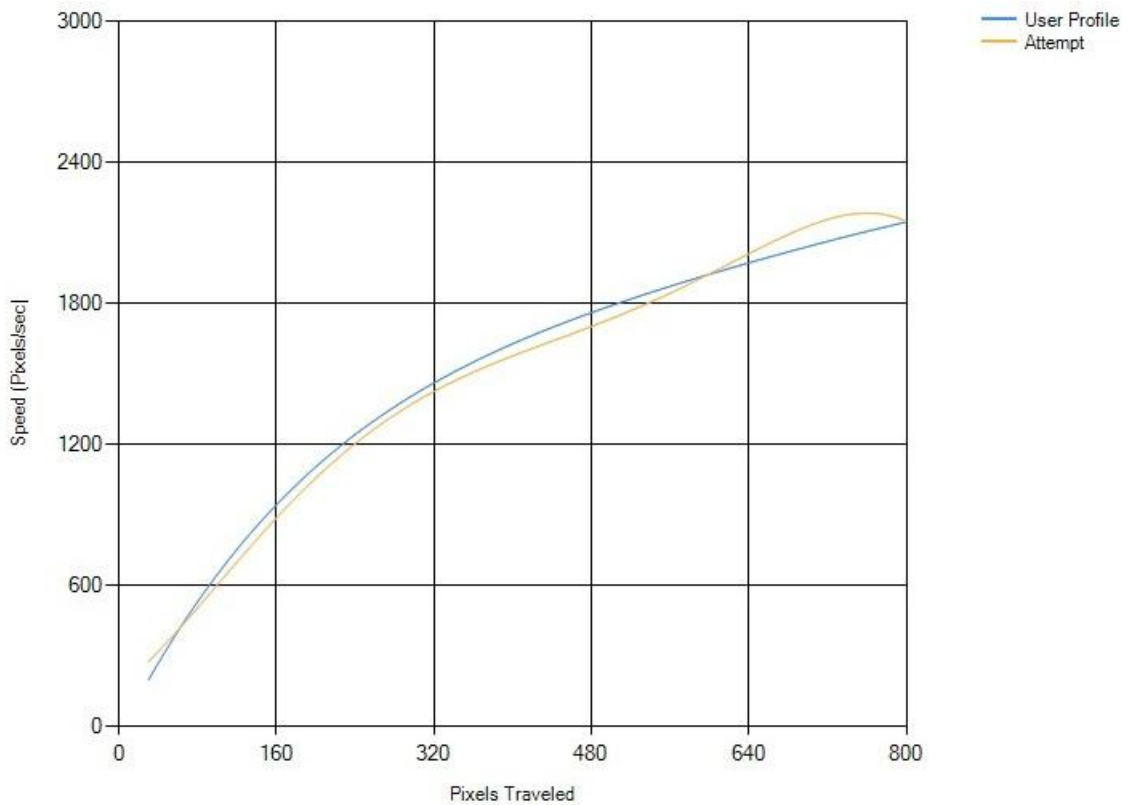
$$\text{Para } 1 \leq i \leq 11: x_i = \frac{(x_{max} - x_{min})}{11} + x_i - 1 \quad (11)$$

$$y_i = f(x_i) \quad (12)$$

A partir desta aproximação é produzida uma curva que traduz as características do utilizador. As intrusões podem assim ser detetadas através da comparação da curva de perfil do

utilizador com a curva da utilização atual. Como exemplo, temos a comparação do perfil de um utilizador com uma sessão fidedigna (figura 14) e a comparação de um perfil com a sessão de um intruso (figura 15).

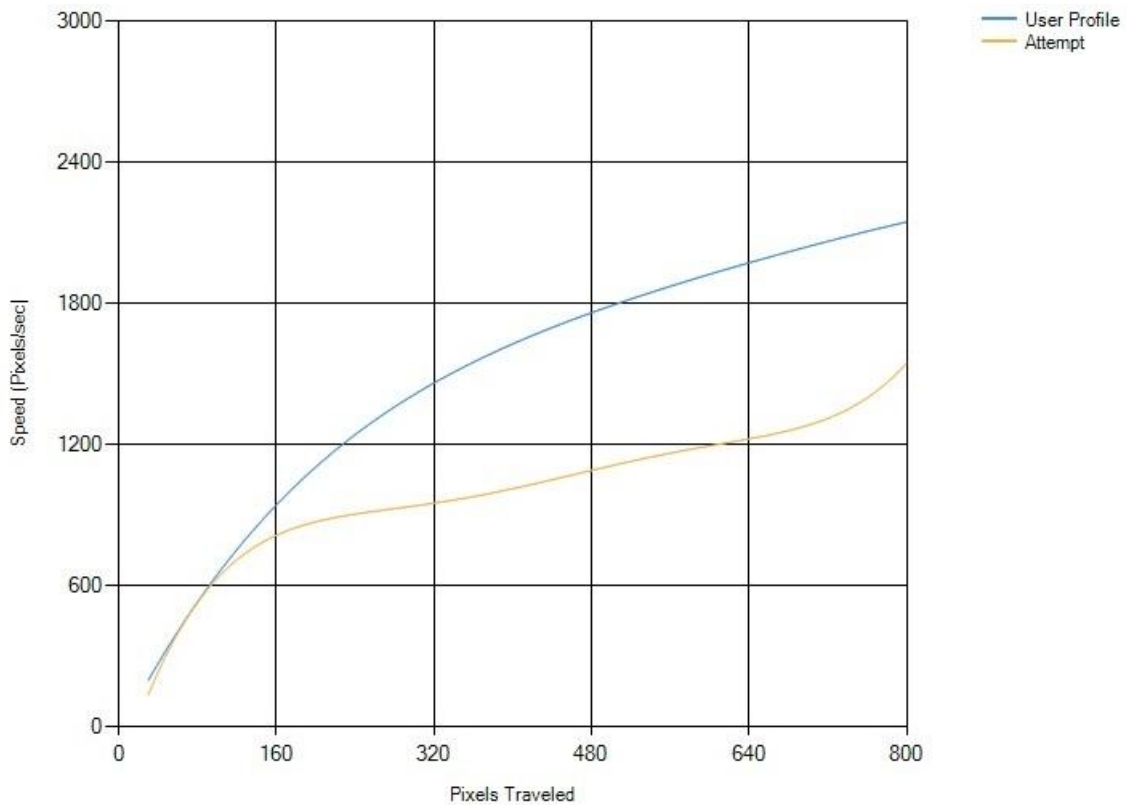
Figura 14 Comparação entre sessões do mesmo utilizador



Na figura 14 vemos que a curva que representa o perfil de utilizador (curva azul) é muito semelhante à sessão atual (linha amarela). Com esta proximidade conseguimos deduzir que a sessão atual foi produzida pelo utilizador fidedigno.

Por outro lado, na figura 15 vemos que a curva azul, representante do perfil de utilizador, está algo distante da curva amarela, representante da sessão atual. Consequentemente, e devido à grande diferença entre as curvas, a sessão atual é marcada com sendo da autoria de um intruso e é despoletado um alerta.

Figura 15 Comparação entre sessões de utilizadores diferentes



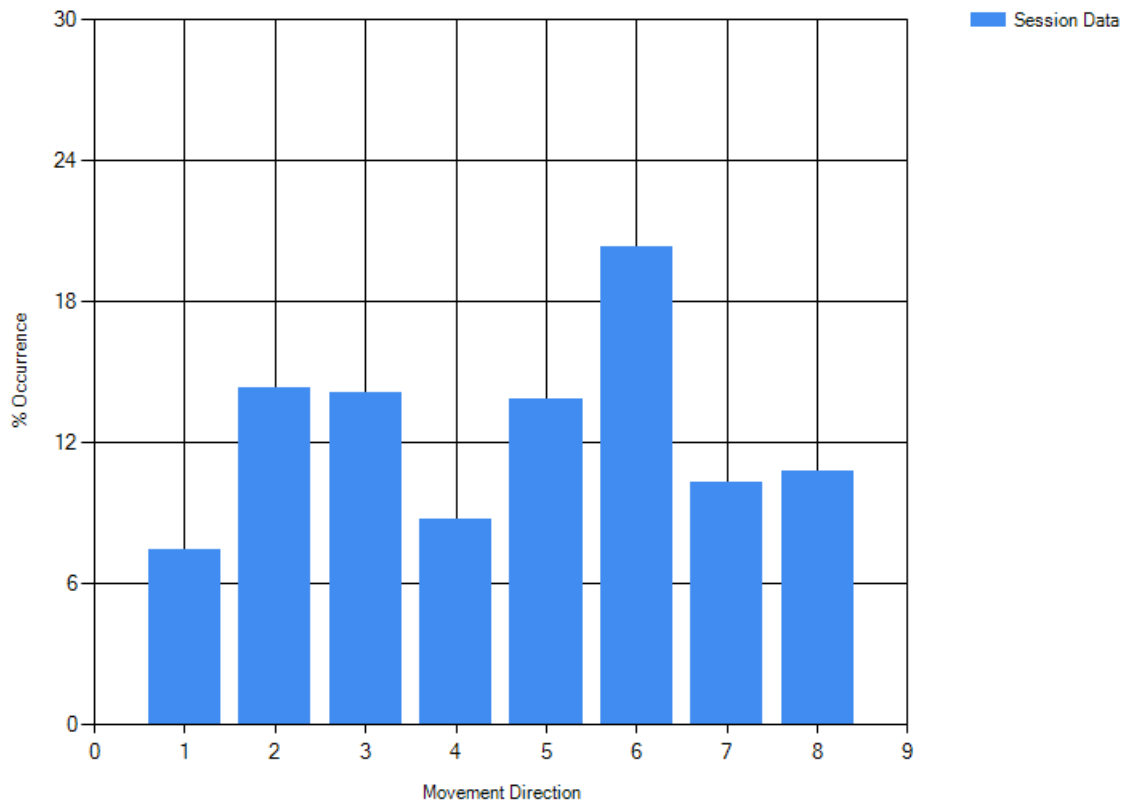
Porcentagem de movimento por direção

Nesta métrica (figura 16) é quantificada a percentagem de movimento para cada uma das direções indicadas na figura 12.

Os dados recolhidos podem ser de grande importância já que cada aplicação tem uma interface diferente. O local onde os botões e barras de ferramentas são colocados leva a ações e comportamentos distintos em termos de movimentos, velocidade e cliques.

Ao obter a direção do movimento, somos capazes de perceber a interação do utilizador com cada aplicação. Devido às diversas tarefas levadas a cabo pelos utilizadores, estes têm de operar com diferentes aplicações e contextos. Por exemplo, escrever um relatório é diferente de programação ou codificação. Além disso, o uso do teclado é também distinta para cada aplicação utilizada.

Figura 16 Histograma das direções de movimento



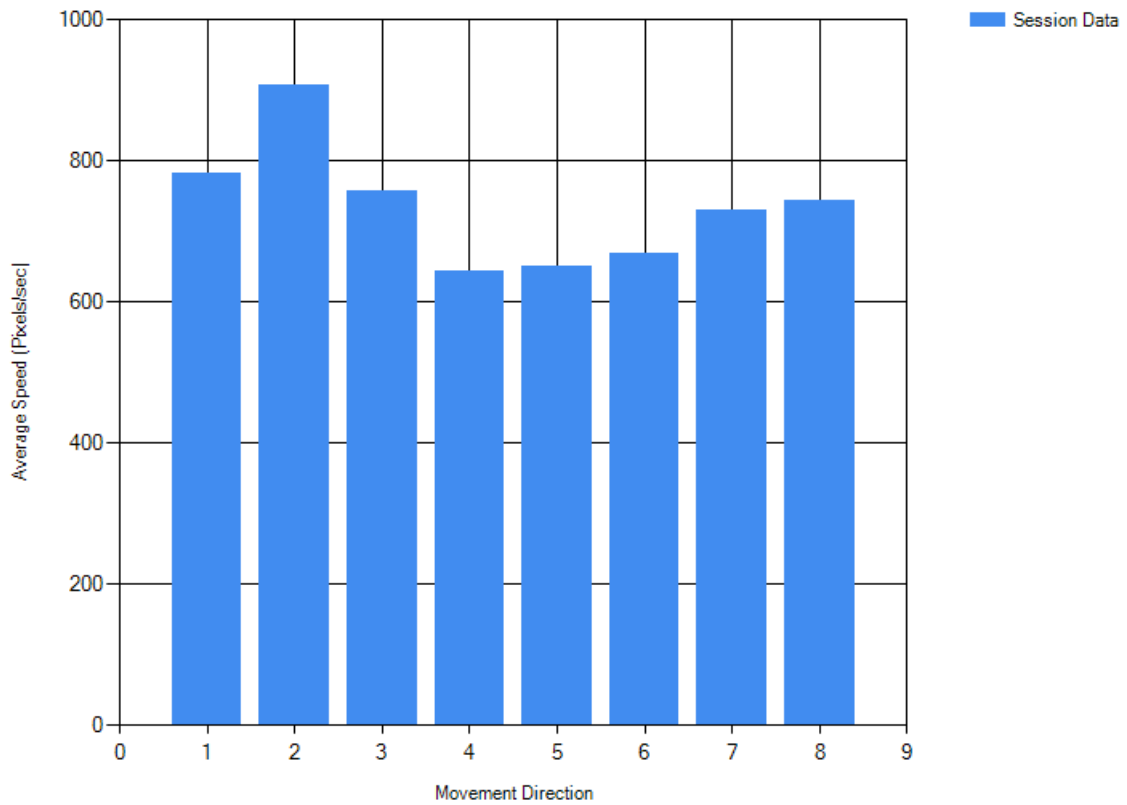
Um exemplo claro na figura 16 é a incidência na direção 6, o que pode indicar uma utilização frequente do menu Iniciar presente nos sistemas Windows.

Velocidade média por direção de movimento

Nesta fase é medida a velocidade média, em pixels/segundo, para cada direção diferente. Olhando para a figura 17 vemos que o utilizador tem velocidades constantes em quase todas as direções embora seja mais rápido na direção 2. A partir destes valores podemos extrair algum conhecimento para construir o perfil comportamental. Tendo em conta que na direção 2 se encontram os botões de minimizar, maximizar e fechar janelas, vemos que ações conhecidas e de alta utilização levam a velocidades mais rápidas por parte deste utilizador.

Com esta métrica conseguimos obter dados mais específicos do utilizador, obtendo informações sobre os seus gestos comuns e interações.

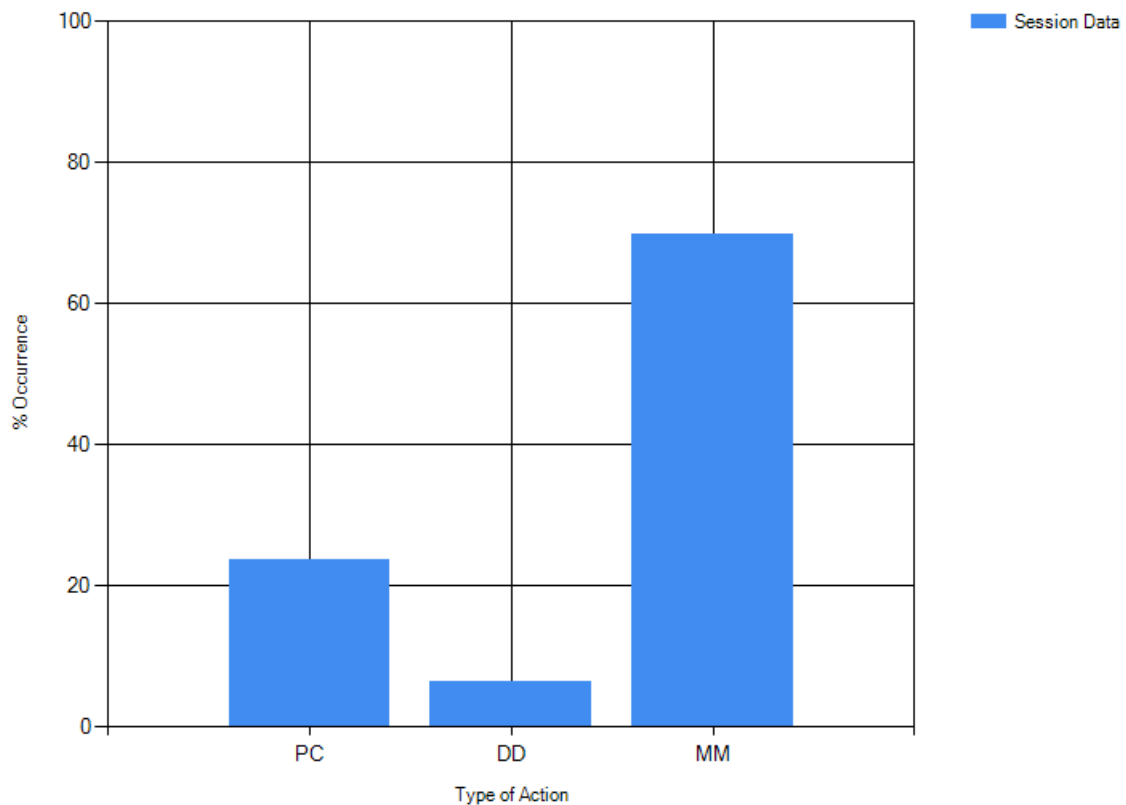
Figura 17 Velocidade média do movimento por direção



Percentagem de tipos de ação

Como já foi dito anteriormente, as ações são cruciais para a autenticação de um utilizador. Adicionalmente, as decisões e ações do utilizador podem ser obtidas para complementarem o seu perfil comportamental.

Figura 18 Histograma dos tipos de ação



A figura 18 mostra as diferentes tipos de ações que o utilizador efetua durante a sua utilização do computador. Conseguimos retirar que o movimento (MM) é a ação mais frequentemente efetuada pelo utilizador, sendo procedida de apontar e clicar (PC) e arrastar e soltar (DD).

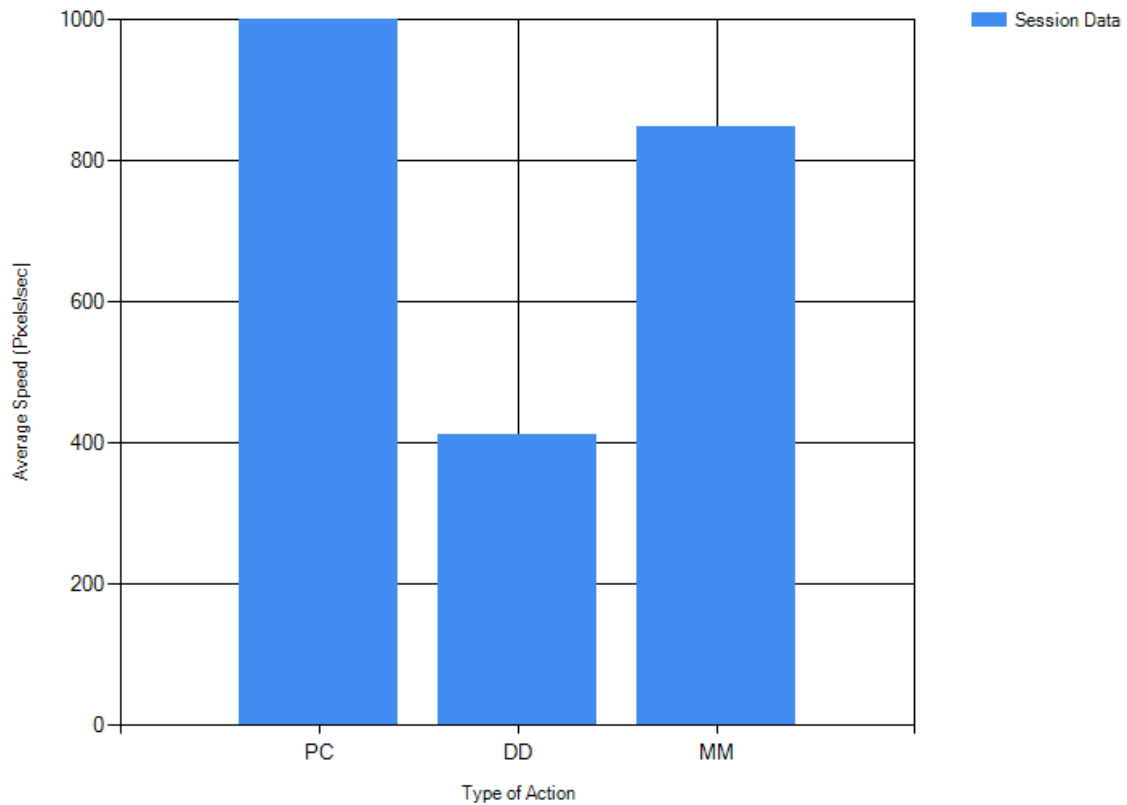
Estes valores são esperados como sendo normais do utilizador. Porém esta métrica usada isoladamente não pode ser considerada fiável. Embora se consiga traçar um perfil global do utilizador, existirão vários indivíduos com o mesmo perfil de comportamento. Para isso foi adicionado uma outra métrica: velocidade.

Velocidade média por tipo de ação

A velocidade é uma característica que relaciona uma decisão a uma resposta muscular. A forma como uma pessoa move a mão ou braço é intrínseca ao seu corpo e personalidade, fazendo parte dos seus traços básicos. Esta informação pode ser transformada em conhecimento, contribuindo para o perfil comportamental de cada utilizador.

Enquanto que os dados sobre as ações do utilizador podem ser um pouco genéricos, quando estas são relacionadas com a velocidade transformam-se em informação muito útil no que diz respeito à interação do utilizador com o rato.

Figura 19 Velocidade média por tipo de ação



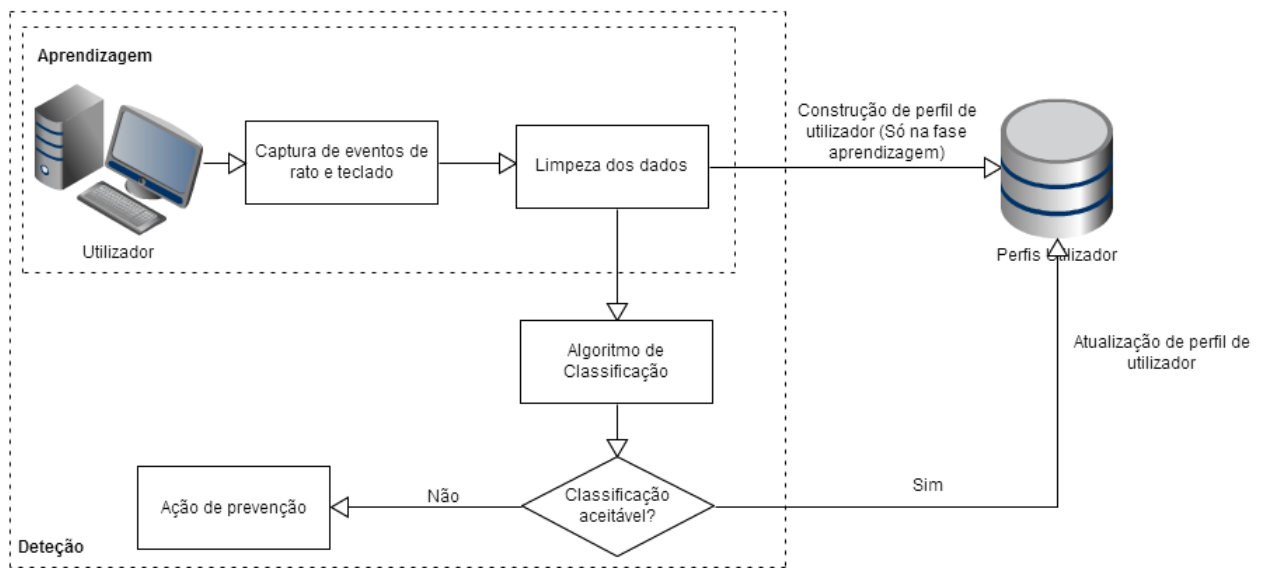
Na figura 19 é apresentada a velocidade média, em pixels/segundo, em relação às ações. Relativamente a este utilizador podemos retirar que é mais rápido na ação de apontar e clicar (PC) do que em ações livres como o movimento (MM). Com isso podemos assumir que o utilizador é mais eficaz quando tem uma ideia clara da tarefa que vai realizar.

6.4. Aplicação

Na fase de aprendizagem a aplicação recolhe os eventos do teclado e rato e constrói o perfil comportamental do utilizador. Este perfil é composto com 12 sessões de rato e teclado.

Posteriormente, e concluída a aprendizagem, o algoritmo entra em modo de deteção e todas as sessões serão classificadas com o intuito de se perceber se uma intrusão está a ocorrer. A figura 20 ilustra o funcionamento geral deste sistema.

Figura 20 Funcionamento Sistema Autenticação Contínua



6.5. Síntese

O problema das intrusões internas, através de utilizadores mascarados, pode ser resolvida pela adoção deste sistema. Com ele podemos garantir que só as pessoas fidedignas podem utilizar os computadores, prevenindo assim problemas como o esquecimento de terminar sessão, a perda de computador e o roubo de credenciais.

Com a utilização de *keystroke dynamics* e *mouse dynamics*, monitorizando permanentemente as ações do rato e do teclado, podemos assegurar e adicionar segurança desde o período de início de sessão até ao fecho da mesma. A aplicação em segundo plano e a utilização de equipamentos já existentes fazem com que os utilizadores nem notem a existência da aplicação.

Este sistema vem combater a falta de sistemas de deteção de intrusões e de deteção de ameaças internas. Os atuais detetores de intrusões não são capazes de detetar ataques novos, já que são baseados em regras previamente definidas, e isso torna este sistema ainda mais interessante.

A grande complexidade da análise do comportamento do rato faz com que se tenham que construir sessões grandes e que podem exceder os 10 minutos de utilização deste equipamento.

Este espaço temporal é um pouco elevado para uma ferramenta desta área sendo que são necessários mais avanços nesta tecnologia de modo a aumentar a sua taxa de deteção.

7. CASOS DE ESTUDO

Ao longo do desenvolvimento deste trabalho verificou-se a necessidade de verificar a acuidade e a possibilidade de implementação num ambiente real. Devido às condicionantes que se entropõem, como a captura de dados sensíveis e a partilha de informação, foram desenvolvidos dois casos de estudo implementados num ambiente controlado. Estes casos de estudo contaram com participantes que estavam cientes da monitorização, contudo, a sua colaboração incidiu na perspetiva de um trabalho normal.

No decorrer destes casos de estudo participaram 10 pessoas, todas elas docentes ou alunos do Departamento de Informática da Universidade do Minho. A sua duração total foi de 1 semana onde os utilizadores foram interagindo com os sistemas de autenticação aqui desenvolvidos. Como resultado desta atividade temos uma grande quantidade de dados de interação que foram submetidos aos algoritmos de classificação para assim medir a eficiência e taxa de deteção dos protótipos.

Este tipo de testes são de elevada importância uma vez que na eventualidade de estes protótipos sejam integrados no produto IAM da PT Inovação, a existência de testes em ambientes reais ajudam a perceber se o projeto está pronto para a utilização genérica do público ou se por outro lado ainda existem lacunas a resolver.

7.1. Caso de Estudo – Sistema de Autenticação Estática

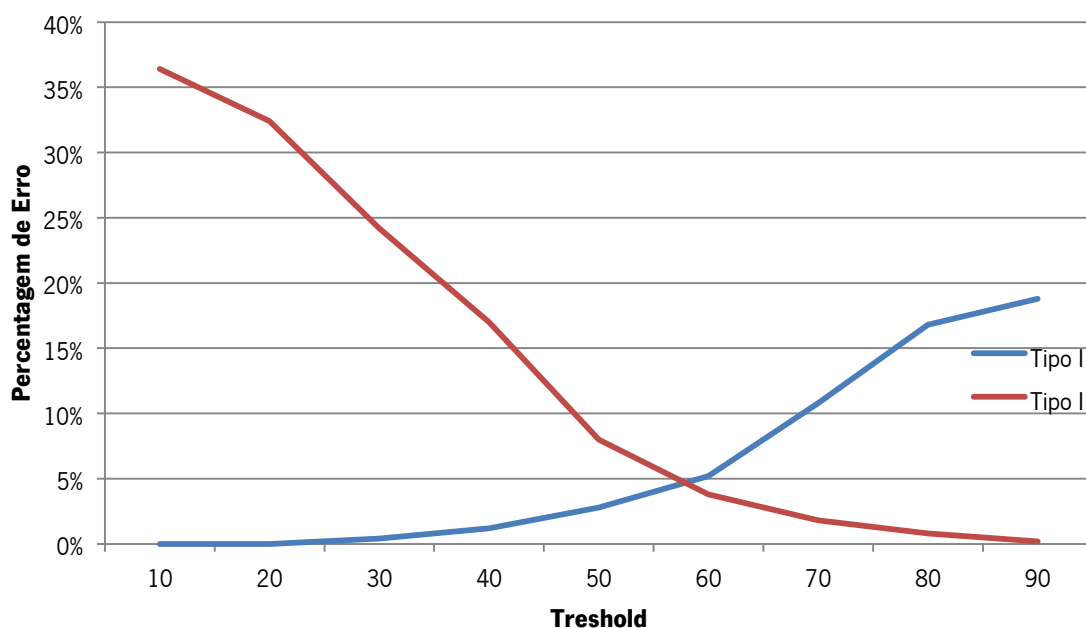
O primeiro caso de estudo diz respeito ao sistema de autenticação estático. Foi criada uma página Web com um formulário de início de sessão simples e algumas contas de utilizadores para teste. Cada utilizador escolheu um nome e palavra-chave livremente e tentou iniciar sessão no sistema, usando as suas credenciais, 30 vezes. No fim foi calculada a taxa de falsos negativos (erro tipo I) ou seja a quantidade de tentativas que um utilizador válido não conseguiu aceder à sua conta. Depois a totalidade das credenciais foram partilhadas e todos os utilizadores tentaram fazer início de sessão 2 vezes para cada uma das contas que não eram suas. Com este teste calculou-se a percentagem de falsos positivos (erros tipo II) ou seja a quantidade de vezes em que uma conta foi acedida por um intruso. No total foram realizados 500 tentativas de início de sessão (300 legítimas e 200 de tentativa de intrusão) por 10 pessoas num período de uma semana.

7.1.1. Resultados Obtidos

Nos próximos gráficos (figuras 21, 22 e 23) estão representados os resultados dos três algoritmos em relação à percentagem de erros de tipo I e tipo II (ver figura x). Para efeitos de comparação foi utilizada a Taxa de Erros Igual (TEI), que é o ponto no gráfico onde as curvas se cruzam. Esta métrica é o ponto onde as taxas de erro de tipo I e II são iguais.

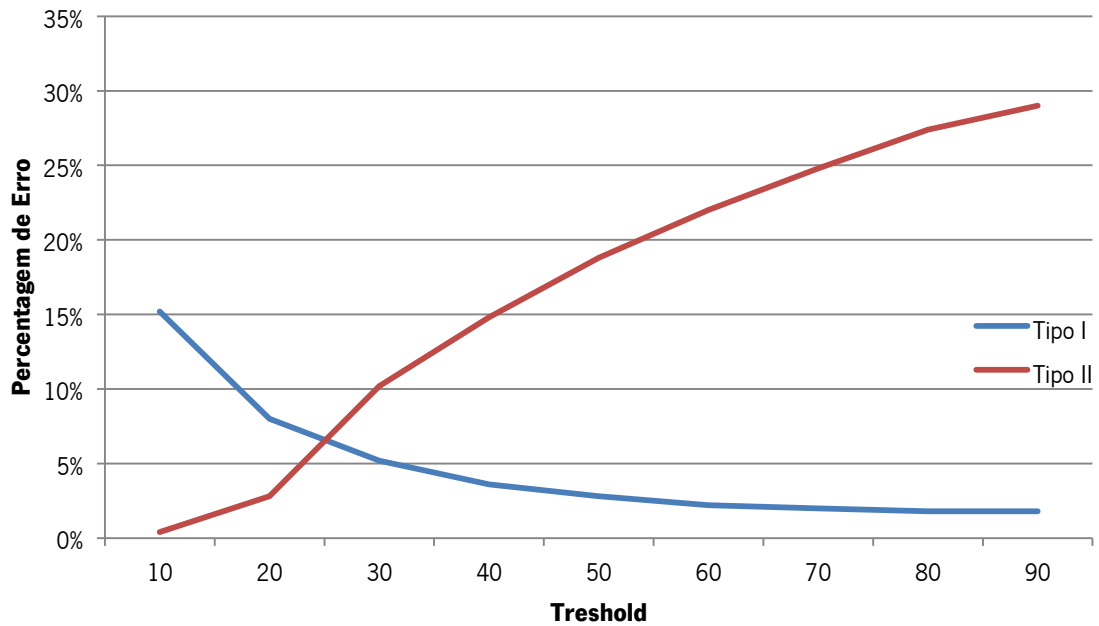
Para o algoritmo *Outlier Count*, figura 21, à medida que o *threshold* aumenta a percentagem de erros de tipo I aumenta, acontecendo o inverso com a percentagem de erros de tipo II. As curvas encontram-se num ponto em que a percentagem de erro é de 4,9%.

Figura 21 Percentagem de erro do algoritmo Outlier Count



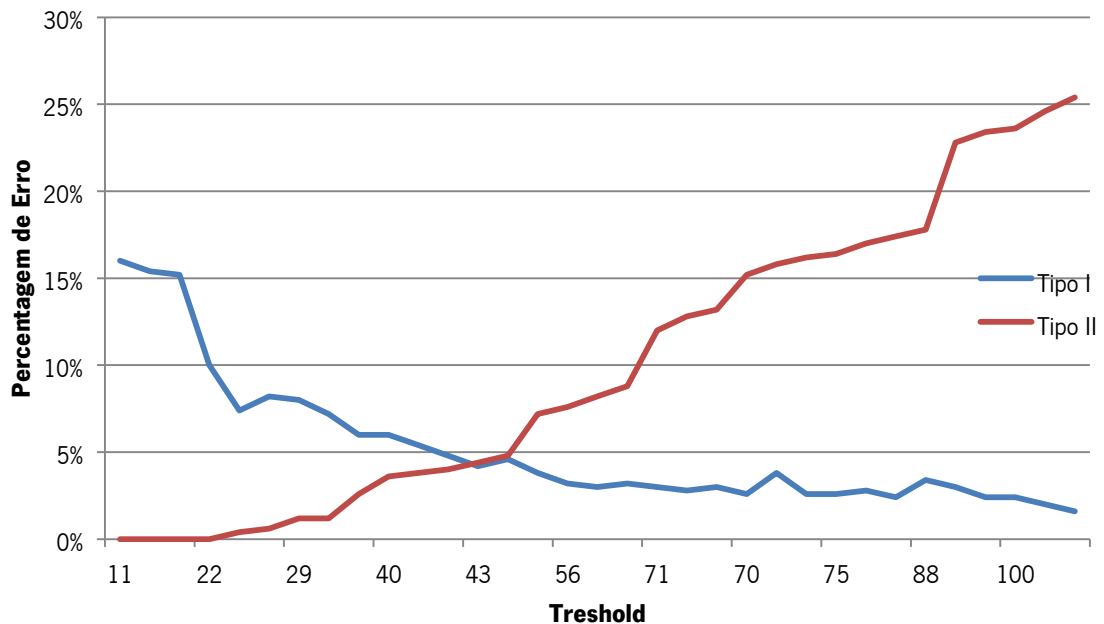
Em relação ao algoritmo da distância de *Mahalanobis*, figura 22, o seu comportamento é contrário ao Outlier Count. Aqui à medida que o *threshold* aumenta, a percentagem de erro tipo I diminui e a de tipo II cresce. A taxa de erro igual para este algoritmo é de 6,5%.

Figura 22 Percentagem de erro do algoritmo Mahalanobis



Como esperado, o algoritmo que combina os dois anteriores foi o que obteve melhores resultados. Como apresentado na figura 23 a sua taxa de erro igual foi de 4,2%.

Figura 23 Percentagem de erro do algoritmo Combinado



Na tabela 12, temos uma comparação entre os três algoritmos. Existem três fatores de comparação:

- Percentagem de Erro de Tipo I a que corresponde uma percentagem de Erro de Tipo II nula. Ou seja, o cenário em que nenhum impostor consegue quebrar contas do sistema;
- Percentagem de Erro de Tipo II a que corresponde uma percentagem de Erro de Tipo I nula. Ou seja, a situação em que um utilizador válido consegue sempre a autenticação usando a sua conta;
- Taxa de Erro Igual. Situação onde os Erros de Tipo I e II são iguais.

Tabela 12 Comparação de algoritmos

Algoritmo	Erro Tipo I para Tipo II = 0%	Taxa de Erro Igual	Erro Tipo II para Tipo I = 0%
Outlier Count	19%	4,9%	32,4%
Mahalanobis	16,6%	6,5%	30,6%
Combinado	10%	4,2%	27%

7.1.2. Análise dos Resultados

Analisando com mais detalhe a tabela 12, vemos que o algoritmo combinado apresenta os melhores resultados para os três fatores de comparação.

Comparativamente aos algoritmos *Outlier Count* e *Mahalanobis*, é de salientar que no algoritmo de *Mahalanobis* apresenta melhores resultados no Erro Tipo I para Tipo II quando este é 0 (primeira coluna da Tabela 12) e para Erro Tipo II para Tipo I quando este é 0 (terceira coluna da Tabela 12) do que o algoritmo *Outlier Count*. Contrariamente, para a Taxa de Erro Igual o algoritmo *Outlier Count* apresenta melhores resultados que o algoritmo de *Mahalanobis*.

Tendo em conta estas duas características, verificou-se que uma junção dos dois algoritmos obteria resultados fiáveis para todos os três fatores de comparação, logo, foi introduzido um novo algoritmo *Combinado*.

Isto significa que em alguns casos o algoritmo mais indicado será um e noutros o outro. Por isso a combinação dos dois algoritmos funciona de forma superior e com resultados bem mais positivos.

Comparando com os resultados apresentados na tabela 9, na área de *keystroke dynamics* só uma solução, (Bergadano, Gunetti, & Picardi, 2002), apresenta resultados melhores do que os aqui apresentados. O projeto (Revett K. , et al., 2006) apresenta resultados similares, com uma taxa de erro igual de 4%.

De notar que todos os testes realizados foram com nomes de utilizador e palavras-chave escolhidos livremente pelos testadores. Em alguns casos o comprimento dessas palavras era de 5 caracteres e noutros de 10, o que leva a uma grande variância de resultados. A maioria dos erros de tipo II ocorreram para palavras de dimensão pequena, o que é expetável devido à menor complexidade da sua escrita. Foi também notado que para todos os algoritmos os melhores resultados foram obtidos com palavras de 8 ou 9 caracteres de dimensão. Caso esta dimensão fosse uma restrição e uma regra na criação de contas de teste, a percentagem de erros do sistema seria menor.

7.2. Caso de Estudo – Sistema de Autenticação Contínuo

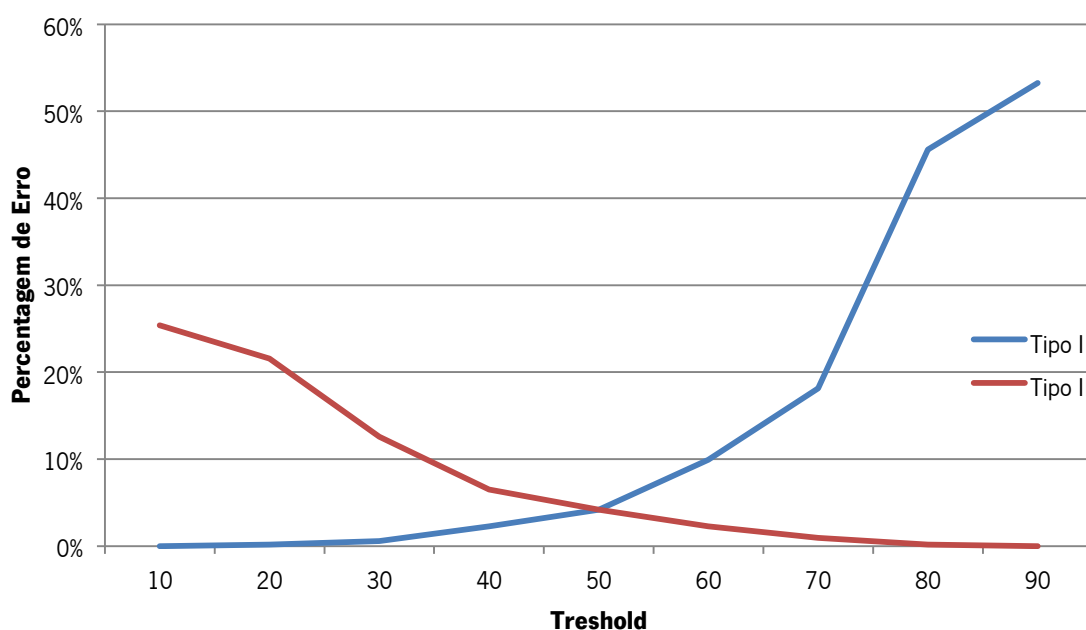
Neste cenário de utilização foi instalado nativamente o sistema de identificação contínua em 6 computadores, sendo estes computadores operados só por uma pessoa. O sistema irá executar durante uma semana com a pessoa fidedigna para se calcular a taxa de erro de tipo I. Depois cada pessoa irá usar a conta de outros durante um dia para assim conseguirmos obter a taxa de erro de tipo II. No total foram realizadas 525 sessões legítimas (390 de teclado e 135 de rato) e 205 sessões de tentativa de intrusão (135 de teclado e 70 de rato). O número de sessões de teclado é significativamente maior em relação às de rato já que o tempo necessário, em utilização contínua, para a construção de uma sessão de teclado ronda os 2 minutos contra os 10 minutos de uma sessão de rato. Quanto mais baixa a taxa de ambos os erros, melhor a eficiência do sistema.

7.2.1. Resultados Obtidos

Tal como na secção 6.1.2 os resultados aqui apresentados terão como referência a tara de erro igual. Como este sistema monitoriza o teclado e o rato de forma distinta e separada, os seus resultados também irão ser analisados dessa forma.

Começando com os resultados da monitorização do teclado, temos as Figuras 24, 25 e 26 que classificam os utilizadores segundo os algoritmos *outlier count*, *mahalanobis* e combinado, respetivamente.

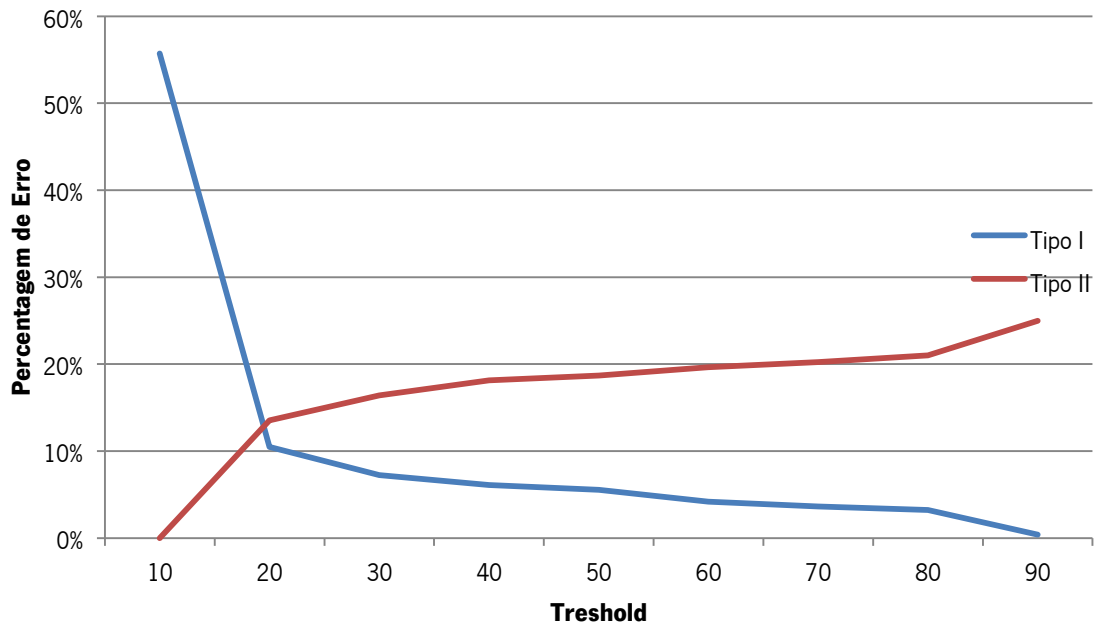
Figura 24 Percentagem de erro do algoritmo Outlier Count para o teclado



Como podemos perceber da figura 24, a monitorização do teclado apresenta uma taxa de erro igual de 4,2%. Também podemos ver que quando a taxa de erro de tipo I se aproxima de 0 a taxa de erro de tipo II ronda os 25%. Na outra extremidade os resultados são um pouco piores já que o erro de tipo I é de 53% quando o erro de tipo II alcança os 0%.

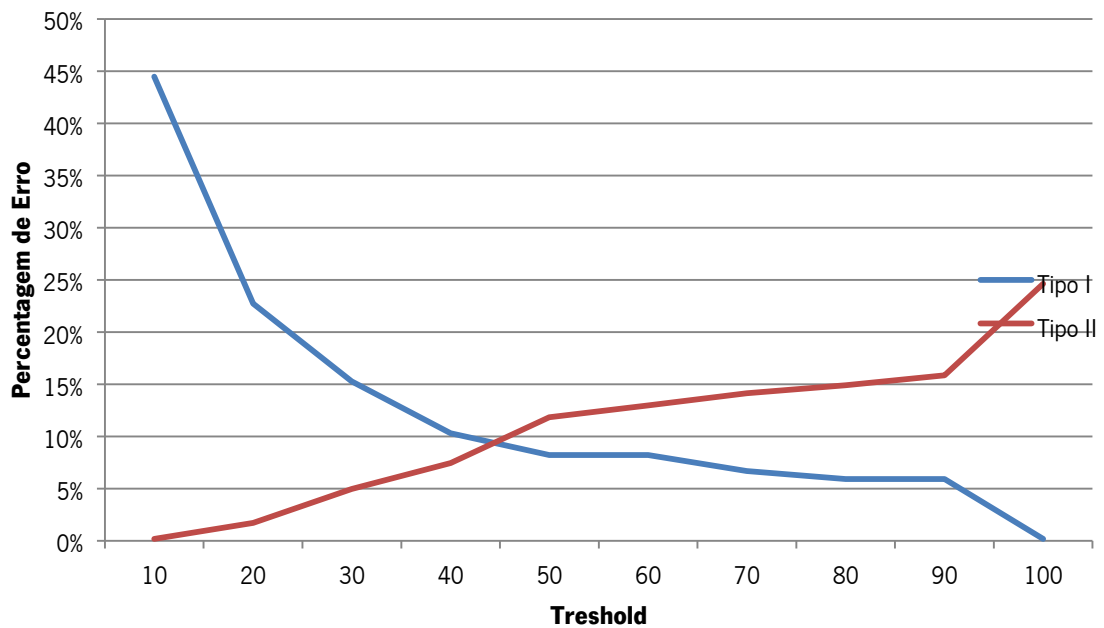
Na figura 25 temos a classificação do algoritmo de Mahalanobis para o mesmo problema. Este algoritmo apresenta resultados um pouco abaixo dos alcançados pelo algoritmo Outlier Count. Com uma taxa de erro igual de 12,6% e 55,7% de erros de tipo I para erros de tipo II nulos e 25% de erros de tipo II na situação oposta, este algoritmo apresenta resultados um pouco aquém do esperado.

Figura 25 Percentagem de erro do algoritmo Mahalanobis para o teclado



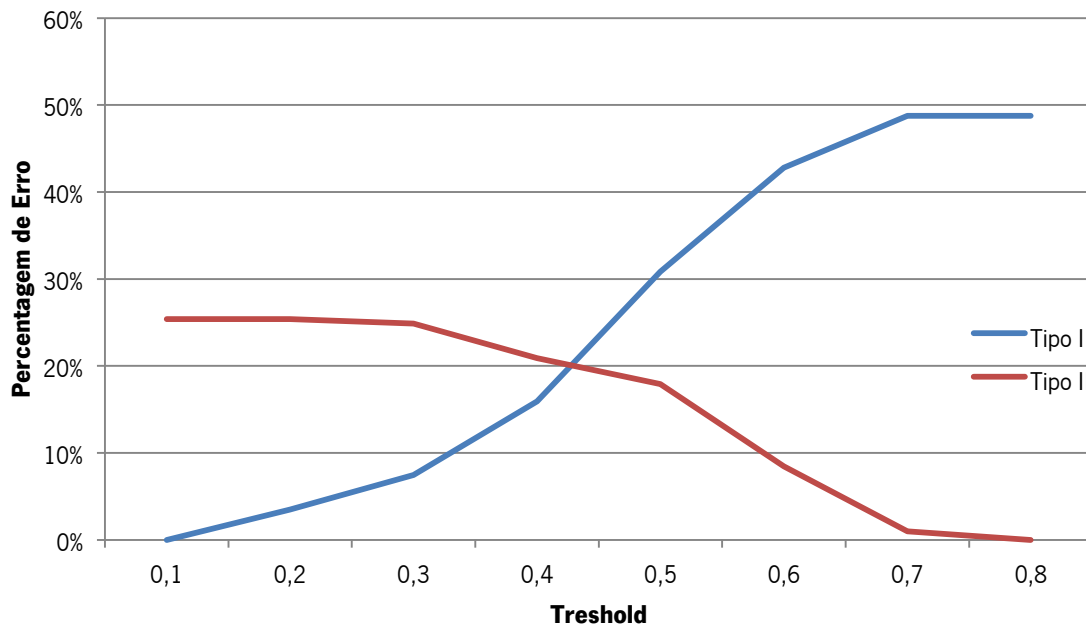
Finalizando o estudo da classificação da monitorização do teclado, o algoritmo combinado, figura 26, apresenta o melhor resultado global, sendo por isso escolhido na maior parte das situações.

Figura 26 Percentagem de erro do algoritmo Combinado para o teclado



Em relação à monitorização do rato, e seguindo a mesma linha do teclado, os resultados estão presentes nas figuras 27, 28 e 29.

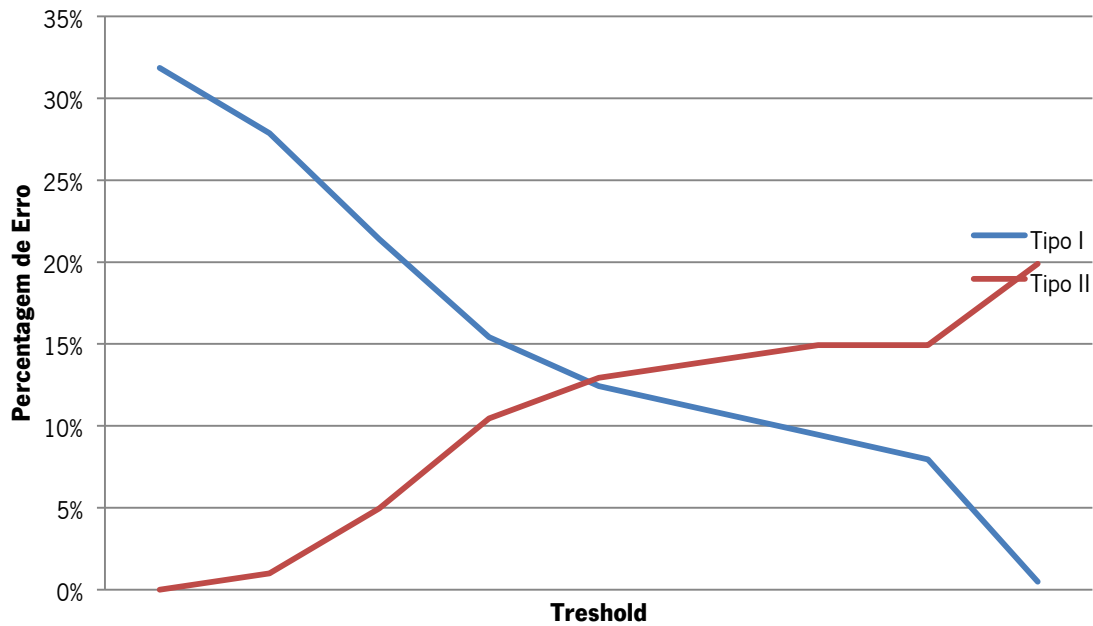
Figura 27 Percentagem de erro do algoritmo Outlier para o rato



Como podemos perceber da figura 27, a monitorização do rato apresenta uma taxa de erro igual de 20,3%. Também podemos ver que quando a taxa de erro de tipo I se aproxima de 0% a taxa de erro de tipo II ronda os 25,3%. Na outra extremidade os resultados são um pouco piores já que o erro de tipo I é de 48,7% quando o erro de tipo II alcança os 0%. Tendo como referência a taxa de erro igual podemos concluir que este algoritmo não é indicado para a problemática da monitorização do rato.

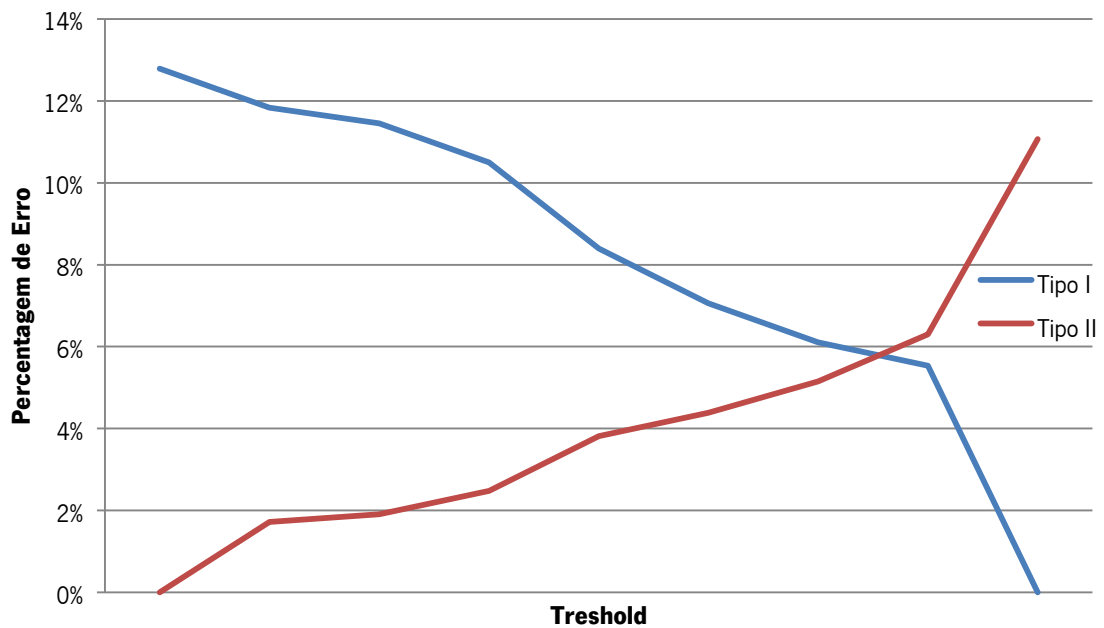
Na figura 28 temos a classificação do algoritmo de Mahalanobis para o mesmo problema. Este algoritmo apresenta resultados um pouco superiores dos alcançados pelo algoritmo Outlier Count. Com uma taxa de erro igual de 12,5% e 31% de erros de tipo I para erros de tipo II nulos e 20% de erros de tipo II na situação oposta, este algoritmo apresenta resultados um pouco melhor que o anterior mas mesmo assim uma taxa de erro igual de 12,5% é ainda um pouco acima do pretendido.

Figura 28 Percentagem de erro do algoritmo Mahalanobis para o rato



Como esperado o algoritmo combinado apresenta as melhores classificações para a monitorização do rato. Através da figura 29 vemos que com uma taxa de erro igual de 5,9% este algoritmo apresenta resultados mais encorajadores.

Figura 29 Percentagem de erro do algoritmo Combinado para o rato



De forma a compreender e analisar melhor os resultados, todos os dados relevantes estão apresentados na tabela 13. Tal como na tabela 12, aqui temos a taxa de erro para erros de tipo II nulos, a taxa de erro igual e a taxa de erro para erros de tipo I nulos para cada um dos algoritmos e cada um dos dispositivos.

Tabela 13 Comparação de algoritmos para os diferentes dispositivos

Algoritmo	Dispositivo	Erro Tipo I para	Taxa de Erro	Erro Tipo II para
		Tipo II = 0%	Igual	Tipo I = 0%
Outlier Count	Teclado	53%	4,2%	25%
	Rato	48,7%	20,3%	25,3%
Mahalanobis	Teclado	55,7%	12,6%	25%
	Rato	31%	12,5%	20%
Combinado	Teclado	44%	9%	24,6%
	Rato	12,4%	5,9%	11%

7.2.2. Análise dos Resultados

Os resultados do sistema de autenticação contínua, devido à sua robustez e complexidade, eram exetáveis que fossem significativamente piores aos resultados obtidos pelo outro sistema aqui apresentado. Porém, estes resultados são muito similares, algo que não era esperado.

O algoritmo combinado é aquele que apresenta melhores resultados globais (ver tabela 13). Apresenta os melhores resultados em 5 dos 6 fatores de comparação. Só perde para o algoritmo *outlier count* a taxa de erro igual no que diz respeito à monitorização do teclado. Com um resultado de 4,2%, esta taxa é igual ao do sistema estático, apesar da maior complexidade deste tipo de monitorização já que a análise de texto livre envolve outro esforço.

No que diz respeito à monitorização do rato, esta apresenta uma taxa de erro igual de 5,9%, que é algo superior às indicadas nas publicações (Ahmed & Traore, 2005), (Pusara & Brodley, 2004) e (Shen, Cai, Guan, Sha, & Du, 2009). Contudo, este pode ser considerado um resultado positivo, devido à maior simplicidade dos algoritmos de classificação aqui adotados em comparação com os escolhidos pela literatura.

8. CONCLUSÕES E TRABALHO FUTURO

8.1. Síntese Do Trabalho

Este trabalho foi desenvolvido com o objetivo de ser integrado na suite de soluções que o projeto IAM (Identity and Access Management) da PT Inovação apresenta. Segundo as suas especificações era necessário a conceção de um sistema de autenticação fiável e que não acrescente custos com equipamento extraordinário.

Para isso, foi proposta uma solução baseada em perfis comportamentais dos utilizadores, perfis esses que são criados a partir da interação normal entre o utilizador com o rato e teclado. Partindo do princípio que os comportamentos e interações são individuais e únicos para cada pessoa, foram criados dois sistemas de autenticação com o objetivo de verificar a identidade do utilizador que está a interagir com o sistema.

Primeiro, foi desenvolvido o sistema de autenticação estática, que não é mais do que um reforço adicional de segurança para os atuais portais de início de sessão. Com a verificação da dinâmica de digitação do nome de utilizador e respetiva palavra-chave, conseguimos, com uma taxa de erro de 4,2%, garantir que a pessoa que se está a autenticar é de fato o utilizador fidedigno.

Aproveitando o desenvolvimento do módulo de interpretação de eventos do teclado construído no sistema acima falado, foi proposto a construção de um sistema de autenticação contínuo, que verifica, constantemente, a identidade da pessoa que está a operar o computador. Para isso foi adicionado um módulo de captura e interpretação dos eventos do rato com o objetivo de conseguir verificar a identidade de um utilizador de forma mais completa e com o maior grau de certeza possível. Embora este sistema de monitorização do rato tenha tido resultados um pouco acima dos apresentados por publicações relacionadas, o seu resultado final, tendo em vista a falta de existência de sistemas deste género, pode ser considerado positivo.

Em termos de trabalho realizado pode-se salientar:

- Estudo e escolha das biometrias comportamentais baseadas em interação humano-computador mais indicadas para este problema;
- Criação de perfis comportamentais únicos para cada utilizador;

- Escolha e implementação de 3 algoritmos de classificação capazes de identificar desvios de comportamento;
- Desenvolvimento e implementação de um sistema de autenticação estático que autentica o utilizador pela forma de digitação do seu nome/palavra-chave;
- Desenvolvimento e implementação de um sistema de autenticação contínua que verifica a identidade do utilizador com base no seu perfil comportamental;
- Realização de testes num ambiente laboratorial;

Embora tenha sido desenhado para ser integrado no projeto IAM, devido à sua arquitetura, este projeto pode ser implementado em vários cenários e plataformas. Também podem ser construídos outros sistemas, com as mais diversas funcionalidades, utilizando os módulos de captura e interpretação de eventos do teclado e rato aqui desenvolvidos. Concluindo, apresenta-se um projeto multifacetado e que pode servir de alavanca a outros projetos futuros.

8.2. Trabalho Relevante Realizado

Este trabalho, como já dito anteriormente, nasceu no âmbito de complementar o produto IAM da PT Inovação. O projeto está a ser desenvolvido nas instalações da PT Inovação com a colaboração do Departamento de Informática da Universidade do Minho.

Como resultado deste projeto foram submetidos os seguintes artigos:

- Henrique Martins, Ângelo Costa, Ricardo Azevedo, Paulo Novais, Henrique Santos: “An Dynamic Intrusion Detection System Based on Profiles”: in Biometrics Journal, 2013, IBS, ISSN: 1541-0420;
- Henrique Martins, Ricardo Azevedo, Paulo Novais, Ângelo Costa: “Biometric and Behavioural Identification of Users in Intrusion Scenarios”. INForum 2013 – Simpósio de Informática, 2013, Évora, Portugal;
- Henrique Martins, Ricardo Azevedo: “Identificação biométrica e comportamental de utilizadores em cenários de intrusão”. Revista Saber & Fazer Telecomunicações nº 11, 2013, PT Inovação.

O sistema de autenticação contínua foi também apresentado e experimentado na iniciativa Open Day em celebração do 14º Aniversário da PT Inovação.

8.3. Trabalho Futuro

Apesar de o trabalho proposto estar realizado, existem algumas melhorias a implementar. Também podem ser desenvolvidos sistemas similares que utilizem as técnicas de *keystroke dynamics* e *mouse dynamics* para a verificação de identidades. Estas sugestões são apresentadas de seguida.

Em relação às melhorias que podem ser implementadas nos dois sistemas, elas são:

- Garantir o suporte a erros ortográficos no sistema de autenticação estática;
- Implementar suporte a dispositivos *touchscreen*. Devido ao seu crescimento, uma expansão dos sistemas aqui propostos para um ambiente de toque é crucial. Existem já algumas propostas para esta temática: (Luca, Hang, Brudy, Lindner, & Hussmann, 2012), (Sae-Bae, Ahmed, Isbister, & Memon, 2012), (Chang, Tsai, & Lin, 2012) e (Saevanee & Pattarasinee, 2008);
- Melhorar os resultados obtidos pelo sistema de monitorização do rato. A utilização de algoritmos de classificação mais complexos e poderosos, por exemplo *Support Vector Machines* (Bennett & Campbell, 2000), levaria a uma melhor avaliação e classificação dos dados.

Falando agora de sistemas que poderiam usar as mesmas técnicas, poderiam ser desenvolvidos algoritmos que garantissem a autenticidade de autores de textos. Isto teria aplicabilidade em aplicações como:

- Clientes de correio eletrónico (Outlook, Gmail, etc.) que poderiam verificar a identidade do autor das mensagens antes do seu envio. O envio só seria efetuado caso o autor da mensagem e o proprietário da conta de correio eletrónico fossem o mesmo;
- Facebook, ou outras redes sociais, de modo a evitar o denominado *“facejacking”*. De modo semelhante ao dos clientes de correio eletrónico, sempre que uma

publicação fosse escrita e antes de a tornar pública, efetuar uma verificação de identidade seria bastante útil;

- Processadores de texto em geral, para a verificação de autoria e possível detecção de plágios.

REFERÊNCIAS

- Ahmed, A. A., & Traore, I. (2005). Anomaly Intrusion Detection based on Biometrics. In *Proceedings from the Sixth Annual IEEE SMC Workshop on Information Assurance* (pp. 452-453).
- Ahmed, A. A., & Traore, I. (2005). Detecting Computer Intrusions Using Behavioral Biometrics. In *Third Annual Conference on Privacy, Security and Trust*.
- Ahmed, A. A., & Traore, I. (2007). A New Biometric Technology Based on Mouse Dynamics. In *IEEE Transactions on Dependable and Secure Computing* (pp. 165-179).
- Altman, A. (23 de Outubro de 2003). *Review of BioPassword 4.5*. Obtido de <http://www.lfca.net/Biometritech%2022502%20review.pdf>
- Al-Zubi, S., Bromme, A., & Tonnes, K. (2003). Using an Active Shape Structural Model for Biometric Sketch Recognition. In *Pattern Recognition* (pp. 187-195). Springer Berlin Heidelberg.
- Apap, F., Honig, A., Hershkop, S., Eskin, E., & Stolfo, S. (2002). Detecting malicious software by monitoring anomalous windows registry accesses. In *In Proceedings of the Fifth International Symposium on Recent Advances in Intrusion Detection* (pp. 16-28).
- Bennett, K., & Campbell, C. (2000). Support vector machines: Hype or hallelujah? . *SIGKDD Explorations 2*, 1-13.
- Bergadano, F., Gunetti, D., & Picardi, C. (2002). User authentication through keystroke dynamics. In *ACM Transactions on Information and System Security (TISSEC)* (pp. 367-397). ACM New York.
- Bhatkar, S., Chaturvedi, A., & Sekar, R. (2006). Dataflow Anomaly Detection. In *Proceedings of the 2006 IEEE Symposium on Security and Privacy* (pp. 48-62). IEEE Computer Society Washington.
- Biometric Signature ID. (2013). *BioSig-ID™*. Obtido de <http://www.biosig-id.com/products/biosig-id%E2%84%A2/>
- BioPassword*. (s.d.). Obtido de BioPassword: <http://www.biopassword.com/>
- Brian, M. (6 de Junho de 2012). *Bad day for LinkedIn: 6.5 million hashed passwords reportedly leaked*. Obtido de TNW - The Next Web: <http://thenextweb.com/socialmedia/2012/06/06/bad-day-for-linkedin-6-5-million-hashed-passwords-reportedly-leaked-change-yours-now/>
- Bromme, A., & Al-Zubi, S. (2003). Multifactor Biometric Sketch Authentication. In *IN PROCEEDINGS OF THE FIRST CONFERENCE ON BIOMETRICS AND ELECTRONIC SIGNATURES OF THE GI WORKING GROUP BIOSIG* (pp. 81-90).

- Chang, T.-Y., Tsai, C.-J., & Lin, J.-H. (2012). A graphical-based password keystroke dynamic authentication system for touch screen handheld mobile devices. *The Journal of Systems and Software* 85, 1157-1165.
- Cho, S., Han, C., Han, D. H., & Kim, H.-I. (2000). Web-Based Keystroke Dynamics Identity Verification Using Neural-Network. In *Journal of Organizational Computing and Electronic Commerce* (pp. 295-307).
- CLUSIT. (2012). *Italian Information Security Association 2012 Report*.
- de Magalhães, S. T., Revett, K., & Santos, H. M. (2005). Password Secured Sites – Stepping Forward With Keystroke Dynamics. In *International Conference on Next Generation Web Services Practices*. IEEE Computer Society.
- de Vel, O., Anderson, A., Corney, M., & Mohay, G. (2001). Mining E-mail Content for Author Identification Forensics. In *Sigmod Record* (pp. 55-64).
- Debar, H. (19 de Maio de 2010). *Intrusion Detection FAQ: What is behavior-based intrusion detection?* Obtido de SANS: http://www.sans.org/security-resources/idfaq/behavior_based.php
- Delac, K., & Grgic, M. (2004). A survey of biometric recognition methods. In *46th International Symposium Electronics in Marine* (pp. 184-193).
- Denning, D. E. (1987). An Intrusion-Detection Model. In *IEEE Transactions on Software Engineering - Special issue on computer security and privacy* (pp. 222-232). IEEE Press Piscataway.
- ENISA. (2012). *Password security: a joint effort between end-users and service providers*.
- Federal Office for Information Security. (2011). *The IT Security Situation in Germany in 2011*.
- Feng, H. H., Kolesnikov, O. M., Fogla, P., Lee, W., & Gong, W. (2003). Anomaly Detection Using Call Stack Information. In *In Proceedings of the 2003 IEEE Symposium on Security and Privacy* (pp. 62-75).
- Ferreira, J., Santos, H., & Patrão, B. (2011). Intrusion detection through keystroke dynamics. In *The Proceedings of the 10th European Conference on Information Warfare and Security* (pp. 81-90). Tallin.
- Fu, Y., & Shih, M.-Y. (2002). A Framework for Personal Web Usage Mining. In *In Intl Conf. on Internet Computing* (pp. 595-600).
- Gamboa, H., & Fred, A. (2003). An Identity Authentication System Based On Human Computer Interaction Behaviour. In *In Proceedings of the 3rd International Workshop on Pattern Recognition in Information Systems*.

- Gamboa, H., & Fred, A. (2004). A Behavioural Biometric System Based on Human Computer Interaction. In *Proceedings of SPIE*.
- Garg, A., Rahalkar, R., Upadhyaya, S., & Kwiat, K. (2006). Profiling Users in GUI Based Systems for Masquerade Detection. In *Proceedings of the 2006 IEEE Workshop on Information Assurance* (pp. 48-54).
- Giffin, J. T., Jha, S., & Miller, B. P. (2004). Efficient Context-Sensitive Intrusion Detection. In *Network and Distributed Systems Security Symposium*.
- Giot, R., El-Abed, M., Hemery, B., & Rosenberger, C. (2011). Unconstrained keystroke dynamics authentication with shared secret. *Computers & Security* 30, 427-445.
- Goecks, J., & Shavlik, J. (2000). Learning Users' Interests by Unobtrusively Observing Their Normal Behaviour. In *Proceedings of the 5th international conference on Intelligent user interfaces* (pp. 129-132). ACM New York.
- Haider, S., Abbas, A., & Zaidi, A. K. (2000). A Multi-Technique Approach for User Identification through Keystroke Dynamics. In *IEEE International Conference on Systems* (pp. 1336-1341).
- Henderson, N. J., Papakostas, T. V., White, N. M., & Hartel, P. H. (2001). Polymer Thick-Film Sensors: Possibilities for Smartcard Biometrics. In *Sensors and their applications XI* (pp. 83-89).
- Henderson, N., White, N., Veldhuis, R., Hartel, P., & Shump, K. (2002). Sensing Pressure For Authentication. In *Proceedings of 3rd IEEE Benelux Signal Processing Symposium*.
- Ilonen, J. (2003). Keystroke dynamics. In *Advanced Topics in Information Processing–Lecture*.
- Intensity Analytics Corporation. (2011). *CVMetrics™ in Summary*. Obtido de CVMetrics: <http://www.intensityanalytics.com/media.aspx>
- Joyce, R., & Gupta, G. (1990). Identity Authentication Based on Keystroke Dynamics. In *Communications of the ACM* (pp. 168-176).
- Kamp, P.-H. (7 de Junho de 2012). *LinkedIn Password Leak: Salt Their Hide*. Obtido de ACM - Queue: <http://queue.acm.org/detail.cfm?id=2254400>
- Karnan, M., Akila, M., & Krishnaraj, N. (2011). Biometric personal authentication using keystroke dynamics: A review. *Applied Soft Computing Volume 11*, 1565-1573.
- Kayacik, H. G., Zincir-Heywood, A. N., & Heywood, M. I. (2012). Intrusion Detection Systems. In *Signal Processing*.
- Keytrac. (2013). *Protects your application's access data*. Obtido de <https://www.keytrac.net/>

- Killourhy, K. S., & Maxion, R. A. (2009). Comparing Anomaly-Detection Algorithms for Keystroke Dynamics. In *Proceedings of DSN*, 125-134.
- Kosoresow, A. P., & Hofmeyr, S. A. (1997). Intrusion Detection via System Call Traces. In *IEEE Software* (pp. 35-42). IEEE Computer Society Press Los Alamitos.
- Koychev, I., & Schwab, I. (2000). Adaptation to Drifting User's Interests. In *In Proceedings of ECML2000 Workshop: Machine Learning in New Information Age* (pp. 39-46).
- Lee, W., Stolfo, S. J., & Wok, K. W. (1999). A Data Mining Framework for Building Intrusion Detection Models. In *In IEEE Symposium on Security and Privacy* (pp. 120-132).
- Liang, T.-P., & Lai, H.-J. (2002). Discovering User Interests from Web Browsing Behavior: An Application to Internet News Services. In *Proceedings of the 35th Annual Hawaii International Conference on System Sciences* (pp. 2718-2727). IEEE Computer Society Washington.
- Luca, A. D., Hang, A., Brudy, F., Lindner, C., & Hussmann, H. (2012). Touch me once and I know it's you! Implicit Authentication based on Touch Screen Patterns. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 987-996.
- Lunt, T. (1993). A survey of intrusion detection techniques. In *Computers & Security* (pp. 406-418). Elsevier.
- Maesschalck, R. D., Jouan-Rimbaud, J., & Massart, D. L. (2000). The Mahalanobis Distance. *Chemometrics and Intelligent Laboratory Systems Vol. 50*, 1-18.
- Mahmood, N. (31 de Março de 2010). *Here's How Easily A Hacker Can Crack Your Weak Passwords*. Obtido de The Tech Journal: <http://thetechjournal.com/electronics/computer/security-computer-electronics/heres-how-easily-a-hacker-can-crack-your-weak-passwords.shtml>
- Marinos, L., & Sfakianakis, A. (2012). *ENISA Threat Landscape - Responding to the Evolving Threat Environment*.
- McLachlan, G. J. (1999). Mahalanobis Distance. *Resonance Vol. 4*, 20-26.
- Monrose, F., & Rubin, A. D. (1999). Keystroke Dynamics as a Biometric for Authentication. In *Future Generation Computer Systems* (pp. 351-359). Elsevier.
- Moskovitch, R., Feher, C., Messermann, A., Kirschnick, N., Mustafic, T., Camtepe, A., . . . Elovici, Y. (2009). Identity Theft, Computers and Behavioral Biometrics. In *IEEE International Conference on Intelligence and Security Informatics* (pp. 155-160).
- NIST/SEMATECH. (2012). *e-Handbook of Statistical Methods*. <http://www.itl.nist.gov/div898/handbook/>.

- Novikov, D., Yampolskiy, R. V., & Reznik, L. (2006). Artificial intelligence approaches for intrusion detection. In *Systems, Applications and Technology Conference* (pp. 1-8). IEEE Long Island.
- Novikov, D., Yampolskiy, R., & Reznik, L. (2006). Anomaly Detection Based Intrusion Detection. In *Third International Conference on Information Technology: New Generations* (pp. 420-425). IEEE Computer Society Washington.
- O'Gorman, L. (2003). Comparing Passwords, Tokens, and Biometrics for User Authentication. *Proceedings of the IEEE, Vol. 91, No. 12*, 2019-2040.
- Pozadzides, J. (26 de Março de 2007). *How I'd Hack Your Weak Passwords*. Obtido de One Man's Blog: <http://onemansblog.com/2007/03/26/how-id-hack-your-weak-passwords/>
- Público. (09 de 05 de 2013). *Hackers roubam 45 milhões de dólares em 27 países*. Obtido de Jornal Público: <http://www.publico.pt/mundo/noticia/hackers-roubam-45-milhoes-de-dolares-em-27-paises-1593943>
- Pusara, M., & Brodley, C. E. (2004). User re-authentication via mouse movements. In *Proceedings of the 2004 ACM workshop on Visualization and data mining for computer security* (pp. 1-8). ACM New York.
- Revet, K., de Magalhães, S. T., & Santos, H. (2005). DataMining a Keystroke Dynamics Based Biometrics Database Using Rough Sets. In *Workshop on Extraction of Knowledge from Databases and Warehouses: proceedings*. Covilhã: IEEE.
- Revet, K., Gorunescu, F., Gorunescu, M., Ene, M., de Magalhães, S. T., & Santos, H. M. (2006). Authenticating computer access based on keystroke dynamics using a probabilistic neural network. In *2nd Annual International Conference on Global e-Security*. Docklands.
- Revet, K., Gorunescu, F., Gorunescu, M., Ene, M., de Magalhães, S. T., & Santos, H. M. (2007). A machine learning approach to keystroke dynamics based user authentication. In *International Journal of Electronic Security and Digital Forensics* (pp. 55-70). Inderscience Publishers.
- Roesch, M. (1999). Snort - Lightweight Intrusion Detection for Networks. In *Proceedings of LISA '99: 13th Systems Administration Conference*.
- Sae-Bae, N., Ahmed, K., Isbister, K., & Memon, N. (2012). Biometric-Rich Gestures: A Novel Approach to Authentication on Multi-touch Devices. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 977-986.
- Saevanee, H., & Pattarasinee, B. (2008). User Authentication using Combination of Behavioral Biometrics over the Touchpad acting like Touch screen of Mobile Device. *International Conference on Computer and Electrical Engineering*, 82-86.

- Schlimmer, J. C., & Granger, R. H. (1986). Incremental Learning from Noisy Data. In *Machine Learning* (pp. 317-354). Kluwer Academic Publishers Hingham.
- Schuckers, S. A. (2002). Spoofing and Anti-Spoofing Measures. In *Information Security Technical Report* (pp. 56-62).
- Shanmugapriya, D., & Padmavathi, G. (2009). A Survey of Biometric keystroke Dynamics: Approaches, Security and Challenges. In *International Journal of Computer Science and Information Security* (pp. 115-119).
- Shen, C., Cai, Z., Guan, X., Sha, H., & Du, J. (2009). Feature Analysis of Mouse Dynamics in Identity Authentication and Monitoring. In *IEEE International Conference on Communications* (pp. 1-5).
- Software, W. (2013). *TypeWATCH*. Obtido de <http://www.watchfulsoftware.com/en/products/typewatch/overview>
- Stanton, P. T., Yurcik, W., & Brumbaugh, L. (2005). FABS: File and Block Surveillance System for Determining Anomalous Disk Accesses. In *Proceedings of the 2005 IEEE Workshop on Information Assurance and Security* (pp. 207-214).
- Stolfo, S. J., Hershkop, S., Wang, K., Nimeskern, O., & Hu, C.-W. (2003). A Behavior-Based Approach to Securing Email Systems. In *Mathematical Methods, Models and Architectures for Computer Networks Security* (pp. 57-81). Springer Berlin Heidelberg.
- Symantec Intelligence. (2012). *Symantec Intelligence Report: August 2012*.
- Symantec Intelligence. (2012). *Symantec Intelligence Report: November 2012*.
- Trustwave. (2012). *2012 Global Security Report*.
- Tsymbal, A. (2004). *The problem of concept drift: definitions and related work*. Dublin: Department of Computer Science, Trinity College Dublin.
- Wespi, A., Dacier, M., & Debar, H. (2000). Intrusion Detection Using Variable-Length Audit Trail Patterns. In *Recent Advances in Intrusion Detection* (pp. 110-129). Springer Berlin Heidelberg.
- Widmer, G., & Kubat, M. (1996). Learning in the presence of concept drift and hidden contexts. In *Machine Learning* (pp. 69-101). Kluwer Academic Publishers Hingham.
- Xiang, S., Nie, F., & Zhang, C. (2008). Learning a Mahalanobis distance metric for data clustering and classification. *Pattern Recognition Vol. 41*, 3600-3612.
- Yampolskiy, R. V. (2007). Human Computer Interaction Based Intrusion Detection. In *Proceedings of the International Conference on Information Technology* (pp. 837-842). IEEE Computer Society Washington.

- Yampolskiy, R. V., & Govindaraju, V. (2008). Behavioural biometrics: a survey and classification. In *International Journal of Biometrics* (pp. 81-113). Geneva: Inderscience Publishers.
- Yoohwan, K., Jo, J.-Y., & Suh, K. K. (2006). Baseline Profile Stability for Network Anomaly Detection. In *ITNG '06 Proceedings of the Third International Conference on Information Technology: New Generations* (pp. 720-725). IEEE Computer Society Washington.
- Zhang, Z., & Manikopoulos, C. (2003). Investigation of Neural Network Classification of Computer Network Attacks. In *International Conference on Information Technology: Research and Education*. IEEE.
- Zilberman, A. G. (1998). Security method and apparatus employing authentication by keystroke dynamics.