

Universidade Federal de Santa Catarina  
Programa de Pós-Graduação em Ciência da Informação



# Text Mining: o uso de dados não-estruturados como insumo para extração de inteligência

II Workshop de Informação, Dados e Tecnologia

UNIVERSIDADE FEDERAL DA PARAÍBA

De 27 a 29 de novembro de 2018

João Pessoa - Paraíba

<http://www.ufpb.br/widat2018>

[f/contatowidat](https://www.facebook.com/contatowidat)



Moisés Lima Dutra  
[moises.dutra@ufsc.br](mailto:moises.dutra@ufsc.br)

# AGENDA

- Preâmbulo
- Text Mining: visão geral
- Etapas da mineração
- Desafios e oportunidades

# Preâmbulo

# Vamos falar de cinema?





08757TNW

Larry Gigli (Ben Affleck) é violento e medíocre e, além disso, tem fama de se meter em grandes confusões. Ricki (Jennifer Lopez) é rude como os mafiosos. Ricki e Larry são unidos por uma tarefa que rapidamente perde o controle. Será que eles resolverão suas diferenças e aceitarão que se sentem atraídos um pelo outro?

Dirigido por Martin Brest (*Perfume de Mulher*), **CONTATO DE RISCO** é uma comédia de ação e muito sensual.



**CONTATO DE RISCO**  
(GIGLI)

**Apresentações Especiais**

- Formato de Tela do Filme: Standard
- Idiomas do Filme: Inglês (5.1 Dolby Digital), Português e Espanhol (ambos Dolby Surround)
- Legendas do Filme: Português, Inglês e Espanhol
- Menus Interativos em Português, Inglês e Espanhol
- Trailers (sem legendas)
- Seleção de Cenas

Os trailers possuem áudio apenas em inglês e não possuem legendas. O idioma do Menu dependerá da tecnologia do aparelho de DVD.

REVOLUTION STUDIOS APRESENTA UMA PRODUÇÃO DE CITY LIGHT FILMS/CASEY SILVER BEN AFFLECK JENNIFER LOPEZ "GIGLI" JUSTIN BARTHA  
 PRODUTORA MICHAEL KAPLAN MÚSICA DE JOHN POWELL EDITORES DE FILM BILLY WEBB JULIE MONROE LEGENDAS GARY FRUTKOFF DIRETOR DE FOTOGRAFIA ROBERT ELSWIT, ASC  
 EXECUTIVO PRODUTORES JOHN HARDY PRODUTORES DE PRODUÇÃO CASEY SILVER MARTIN BREST PRODUTORA EXECUTIVA MARTIN BREST  
 COLUMBIA TRISTAR PICTURES

©2003 Revolution Studios Distribution Company, LLC. Todos os Direitos Reservados. Duração Aproximada: 121min • Comédia

TRILHA SONORA DISPONÍVEL POR VARESE SARABANDA

08757TNW



Layout e Design © 2002 Columbia TriStar Home Entertainment. Todos os Direitos Reservados. ATENÇÃO: Apenas para edição privada. Exibições públicas em por meios eletrônicos e a reprodução deste disco são violações da lei e passíveis de punição. Dolby e o símbolo do Duplo-D são marcas registradas de Dolby Laboratories Licensing Corporation. Produção no Polo Industrial do Manaus e distribuído por Videolar S.A. - Av. Solimões, 909 - Distrito Industrial - Manaus - AM - CNPJ 04.228.767/0004-13 - Indústria Brasileira, sob licença de Columbia TriStar Home Entertainment do Brasil Ltda. - CNPJ 01.342.611/0001-03. VALORES: O prazo de validade do disco DVD é indeterminado desde que observadas as seguintes condições: Armazenar em local seco, livre de poeira, não expor ao sol, não riscar, não dobrar, não enfiar, não manter a uma temperatura superior a 55° C e umidade acima de 60 gr/m3 e segurar o disco sempre pela lateral e pelo furo central.

COLUMBIA PICTURES

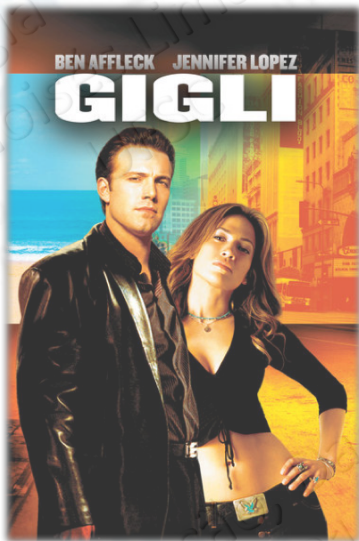
# BEN AFFLECK JENNIFER LOPEZ

# CONTATO DE RISCO

(GIGLI)



<https://coversblog.wordpress.com/2009/01/09/contato-de-risco/>



- Filme lançado em 2003
- Título Original: Gigli
- Título no Brasil: Contato de Risco
- Gênero: Comédia Romântica
- Elenco:
  - **Ben Affleck**
  - **Jennifer Lopez**
  - **Christopher Walken**
  - **Al Pacino**
  - **entre outros**

O que há de  
tão importante  
sobre este filme?

Orçamento:

75 milhões de dólares

Receita:

7,2 milhões de dólares



# "PRÊMIOS" E NOMEAÇÕES

- Ganhou seis prêmios no **Framboesa de Ouro**
  - **Pior Filme**
  - **Pior Diretor**
  - **Pior Ator (Ben Affleck)**
  - **Pior Atriz (Jennifer Lopez)**
  - **Pior Roteiro**
  - **Pior Par (Ben Affleck e Jennifer Lopez)**



# "PRÊMIOS" E NOMEAÇÕES

- Mais três nomeações nas categorias
  - **Pior ator coadjuvante (Al Pacino)**
  - **Pior ator coadjuvante (Christopher Walken)**
  - **Pior atriz coadjuvante (Lanie Kazan)**
- Pra finalizar, ganhou ainda o troféu especial:
  - **Framboesa de Ouro de Pior Comédia dos 25 anos do prêmio**



# COMO MITIGAR ESTE PROBLEMA?

- É possível prever o sucesso de um filme?
- Como saber se um roteiro é adequado?
- Os orçamentos podem ser baseados nas possibilidades de sucesso que determinada obra possui?



# Epagogix

# EPAGOGIX

- Companhia britânica, fundada em 2003
- Utiliza redes neurais e softwares analíticos
- Procura prever se filmes ou programas de TV irão oferecer uma boa possibilidade de retorno do investimento, ainda na etapa anterior à pré-produção
- Também avalia o **grau de sucesso** dos **roteiros** e das **tramas** criadas pelos escritores

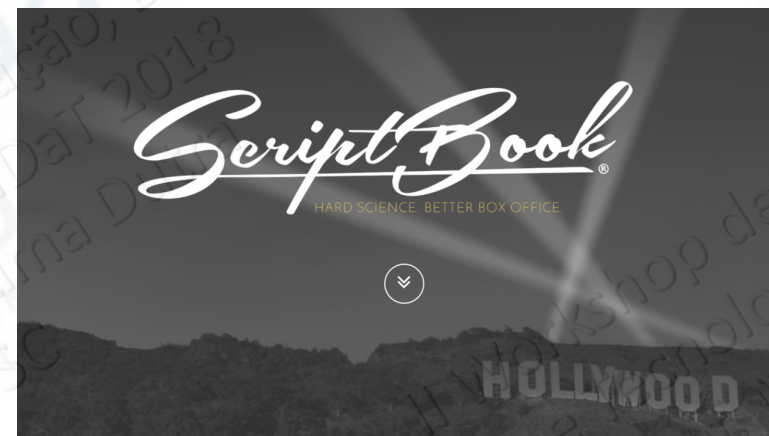
# EPAGOGIX E WARNER BROS.



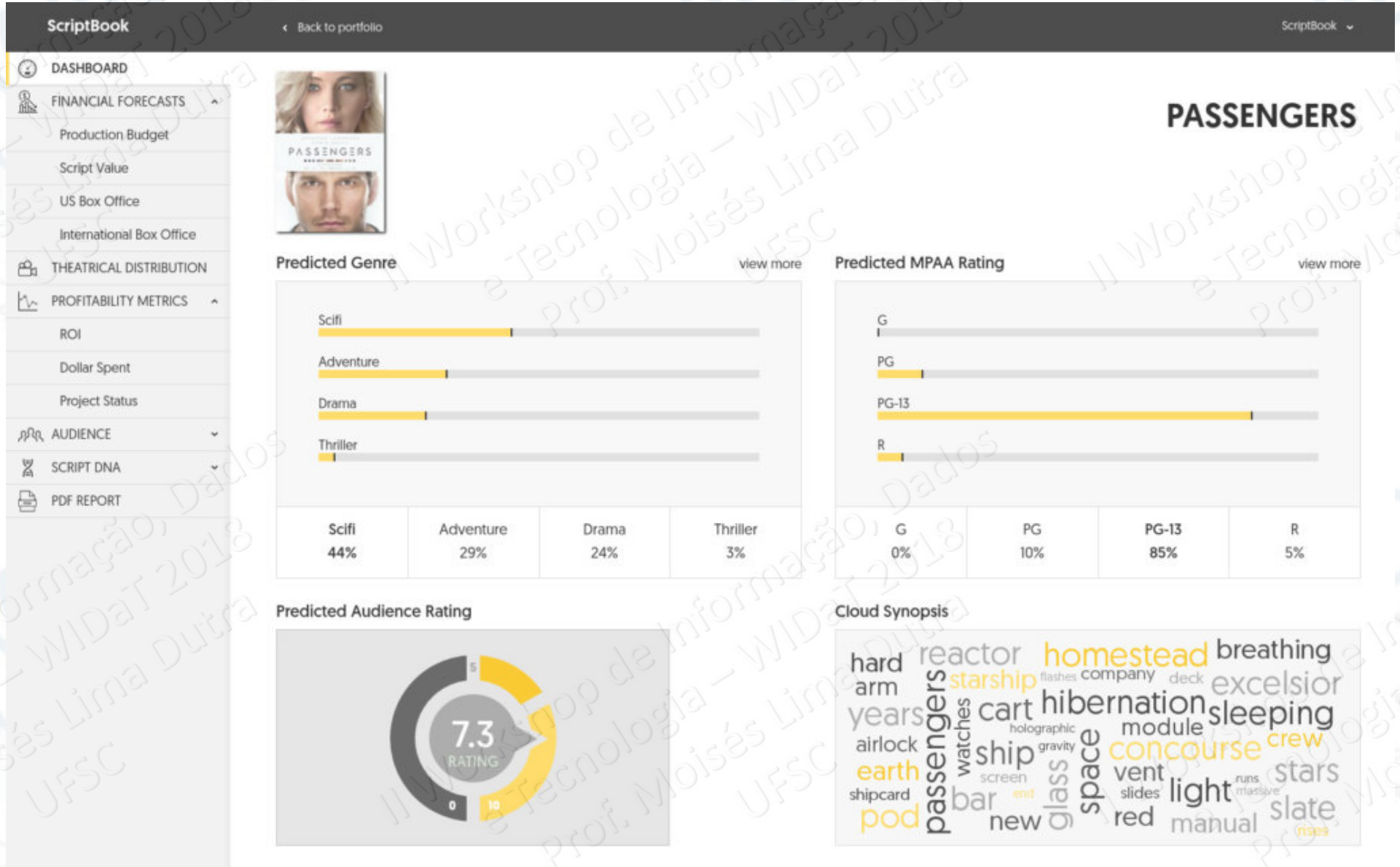
- Parceria com a Warner Brothers Europa
- 16 programas piloto de televisão foram analisados, na busca de se avaliar se eles fariam sucesso nos Estados Unidos
- Para 13 programas, a taxa de audiência calculada foi correta, dentro de uma margem de erro de **2 pontos** percentuais
- Dentre estes, 6 programas tiveram suas taxas calculadas numa margem de erro de apenas **0,06**

# ANÁLISE DE ROTEIROS

- Surgimento da área chamada **Screenwriting Script Analysis**
- Hoje, existem diversas opções de análise de scripts, dentre as quais se destaca o software **ScriptBook**, utilizado a partir de 2016
- O **ScriptBook** cobra 100 dólares para analisar qualquer roteiro enviado pelo seu site



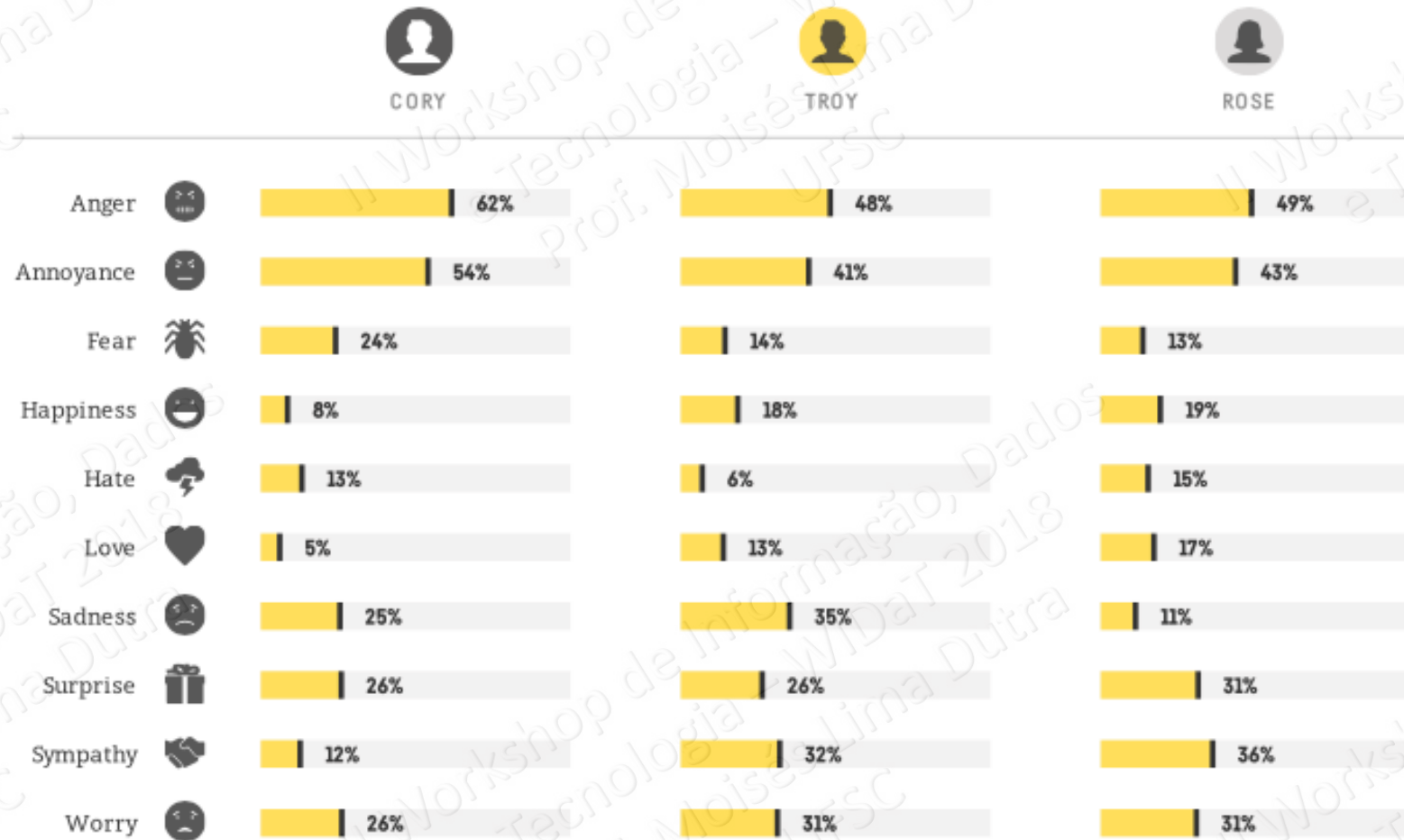
# SCRIPTBOOK DASHBOARD





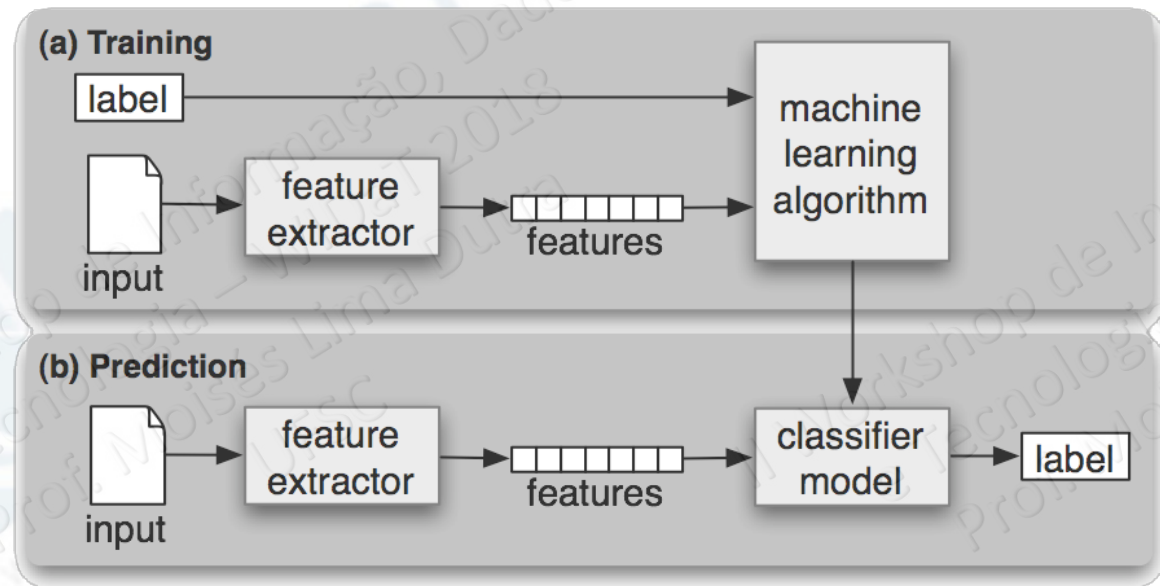
# SCRIPTBOOK DASHBOARD

## Character Sentiment

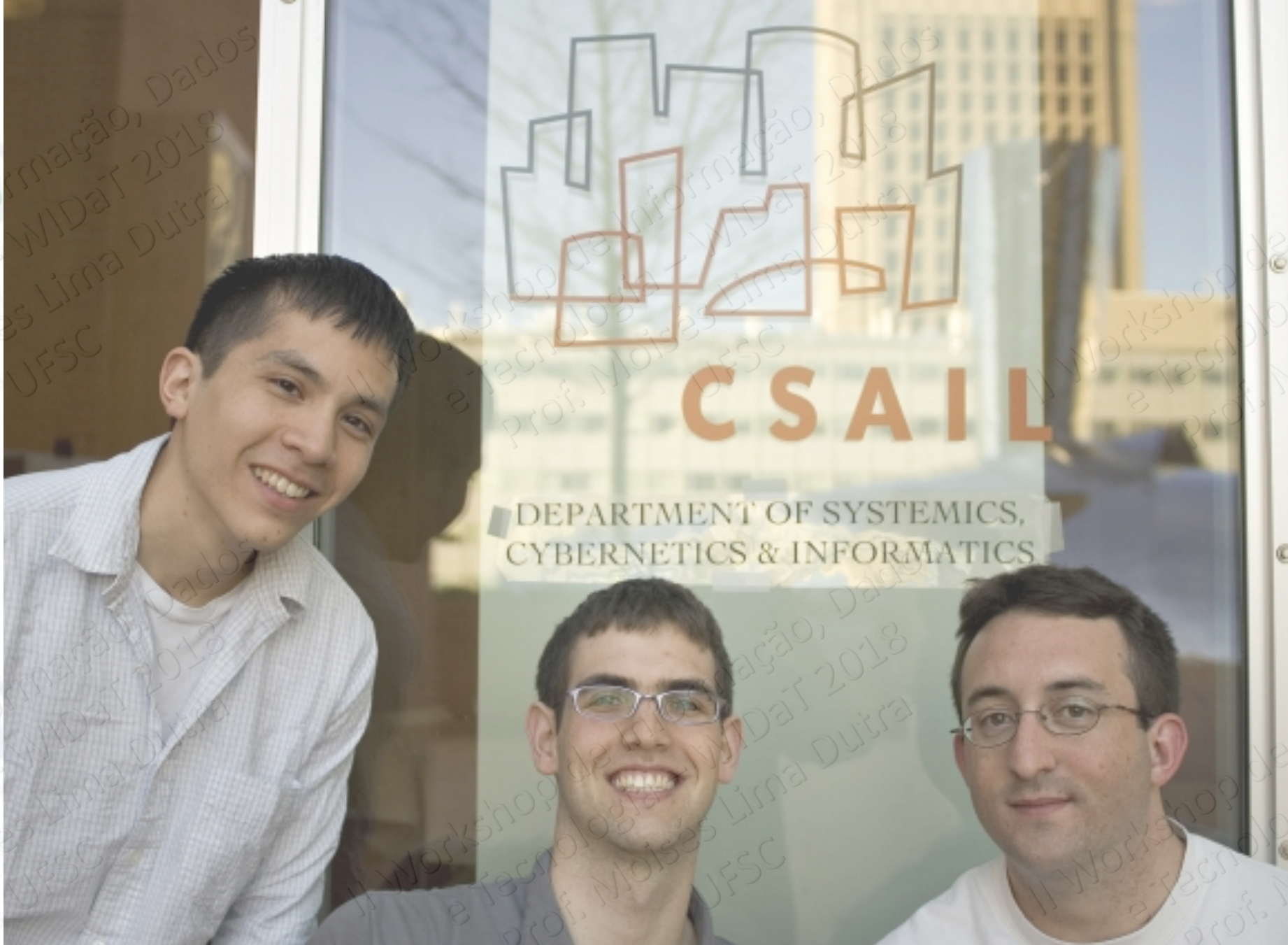


# SCRIPTBOOK

- Procura analisar a possibilidade de sucesso do filme
- Faz uma análise de sentimento de cada personagem
- Calcula quanto cada filme pode gerar de lucro nos Estados Unidos
- Possui um **dataset** composto por **todos os filmes exibidos em cinema nos EUA a partir de 1970**



Vamos falar agora  
de publicação  
científica?



# Router: A Methodology for the Typical Unification of Access Points and Redundancy

Jeremy Stribling, Daniel Aguayo and Maxwell Krohn

## ABSTRACT

Many physicists would agree that, had it not been for congestion control, the evaluation of web browsers might never have occurred. In fact, few hackers worldwide would disagree with the essential unification of voice-over-IP and public-private key pair. In order to solve this riddle, we confirm that SMPs can be made stochastic, cacheable, and interposable.

## I. INTRODUCTION

Many scholars would agree that, had it not been for active networks, the simulation of Lamport clocks might never have occurred. The notion that end-users synchronize with the investigation of Markov models is rarely outdated. A theoretical grand challenge in theory is the important unification of virtual machines and real-time theory. To what extent can

The rest of this paper is organized as follows. For starters, we motivate the need for fiber-optic cables. We place our work in context with the prior work in this area. To address this obstacle, we disprove that even though the much-touted autonomous algorithm for the construction of digital-to-analog converters by Jones [10] is NP-complete, object-oriented languages can be made signed, decentralized, and signed. Along these same lines, to accomplish this mission, we concentrate our efforts on showing that the famous ubiquitous algorithm for the exploration of robots by Sato et al. runs in  $\Omega((n + \log n))$  time [22]. In the end, we conclude.

## II. ARCHITECTURE

Our research is principled. Consider the early methodology by Martin and Smith; our model is similar, but will actually overcome this grand challenge. Despite the fact that such

# Router: A Methodology for the Typical Unification of Access Points and Redundancy

Jeremy Stribling, Daniel Aguayo, and Maxyed Kishn

## ABSTRACT

Many physicists would agree that, had it not been for congestion control, the evaluation of web browsers might never have occurred. In fact, few hackers worldwide would disagree with the essential unification of voice-over-IP and public-private key pair. In order to solve this riddle, we confirm that SMPs can be made stochastic, cacheable, and interposable.

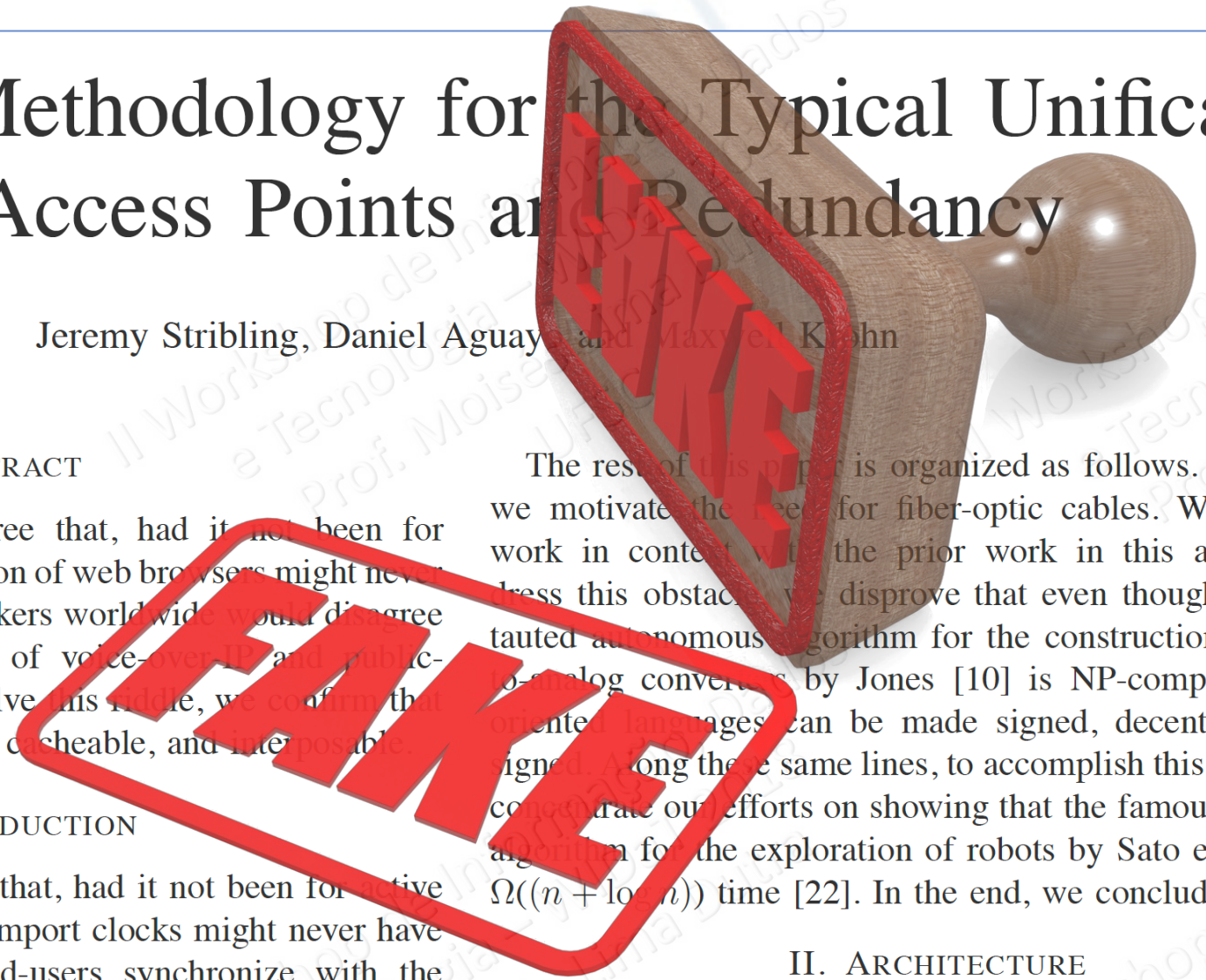
## I. INTRODUCTION

Many scholars would agree that, had it not been for active networks, the simulation of Lamport clocks might never have occurred. The notion that end-users synchronize with the investigation of Markov models is rarely outdated. A theoretical grand challenge in theory is the important unification of virtual machines and real-time theory. To what extent can

The rest of this paper is organized as follows. For starters, we motivate the need for fiber-optic cables. We place our work in context with the prior work in this area. To address this obstacle, we disprove that even though the much-touted autonomous algorithm for the construction of digital-to-analog converters by Jones [10] is NP-complete, object-oriented languages can be made signed, decentralized, and signed. Along these same lines, to accomplish this mission, we concentrate our efforts on showing that the famous ubiquitous algorithm for the exploration of robots by Sato et al. runs in  $\Omega((n + \log n))$  time [22]. In the end, we conclude.

## II. ARCHITECTURE

Our research is principled. Consider the early methodology by Martin and Smith; our model is similar, but will actually overcome this grand challenge. Despite the fact that such



# WMSCI

- **Trabalho aceito para publicação** em 2005
- World Multiconference on Systemics, Cybernetics and Informatics (WMSCI)
- Jeremy Stribling, Daniel Aguayo e Maxwell Krohn

- Um pouco antes do evento acontecer, eles anunciaram que tudo não passava de uma... **brincadeira!**



# SCI GEN - AN AUTOMATIC CS PAPER GENERATOR

- Gerador randômico de papers científicos
- <https://pdos.csail.mit.edu/archive/scigen/>

## SCIgen - An Automatic CS Paper Generator

[About](#) [Generate](#) [Examples](#) [Talks](#) [Code](#) [Donations](#) [Related](#) [People](#) [Blog](#)

### About

SCIgen is a program that generates random Computer Science research papers, including graphs, figures, and citations. It uses a hand-written **context-free grammar** to form all elements of the papers. Our aim here is to maximize amusement, rather than coherence.

One useful purpose for such a program is to auto-generate submissions to conferences that you suspect might have very low submission standards. A prime example, which you may recognize from spam in your inbox, is SCI/IIIS and its dozens of co-located conferences (check out the very broad conference description on the [WMSCI 2005](#) website). There's also a list of [known bogus conferences](#). Using SCIgen to generate submissions for conferences like this gives us pleasure to no end. In fact, one of our papers was accepted to SCI 2005! See [Examples](#) for more details.

We went to WMSCI 2005. Check out the [talks and video](#). You can find more details in our [blog](#).

Also, check out our 10th anniversary celebration project: [SCIpher!](#)

### Generate a Random Paper

Want to generate a random CS paper of your own? Type in some optional author names below, and click "Generate".

Author 1:   
Author 2:   
Author 3:   
Author 4:   
Author 5:

SCIgen currently supports Latin-1 characters, but not the full Unicode character set.



# SCIDTECT

- Entre 2008 e 2013, a IEEE (Institute of Electrical and Electronics Engineers) detectou mais de 120 papers que foram gerados com o SCigen
- Detector de artigos gerados com o SCigen: <http://scidetect.forge.imag.fr/>
- Ainda em 2005, a IEEE retirou o patrocínio do WMSCI



SciDetect



UNIVERSITÉ Grenoble Alpes SPRINGER NATURE

SciDetect is a collaboration between Springer-Verlag GmbH and Université Grenoble Alpes.

SciDetect is an open source software program for evaluating academic papers. The software discovers text that has been generated with the SCigen computer program and other fake-paper generators like Mathgen and Physgen.

SciDetect scans Extensible Markup Language (XML) and Adobe Portable Document Format (PDF) files and compares them against a corpus of fake scientific papers. Using intertextual distance, SciDetect can discover any automatically generated text materials and indicate whether an entire document or its parts are genuine or not. The software relies on sensitivity thresholds to report suspicious activity.

SciDetect is publicly released under the [GNU General Public License \(GPL\), Version 3.0](#).

**Announcements**

[\[2015-03-23\]](#) Springer and Université Joseph Fourier release SciDetect to discover fake scientific papers.

**Download**

To download SciDetect, get a copy of the [Git project](#) for the software. For example: `git clone https://forge.imag.fr/anonymous/git/scidetect/scidetect.git`

# OUTROS EXEMPLOS DE GERADORES DE TEXTO

- O Fabuloso Gerador de Lero-Lero (<http://lerolero.miguelborges.com/>)
- Mathgen (<http://thatsmathematics.com/mathgen/>)
- Lero Lero (<https://www.lerolero.com/>)
- Gerador de Artigo Pós-Modernista (<http://www.elsewhere.org/pomo/>)
- etc.

O que une estes  
dois cenários  
(cinema e publicação  
científica)?

# TEXT MINING



# TEXT MINING

- Mineração de Textos
- Mineração de Dados Textuais
- Baseia-se na **busca por padrões** em **textos digitais em linguagem natural**

# TEXT MINING: MOTIVAÇÕES

- O texto digital em linguagem natural é uma das mais importantes fontes de informação que existe
- International Data Corporation: estimou que mais de **90% dos dados** existentes hoje são **dados textuais**

# TEXT MINING: MOTIVAÇÕES

- Os dados textuais estão majoritariamente na forma não-estruturada, que não é apropriada para ser analisada por um software
- Ainda há um longo caminho a percorrer até a Inteligência Artificial ser capaz de compreender por completo um texto em linguagem natural

Text Mining

ou

Data Mining?



# Profusão de Dados

Informação Digital  
Transações Eletrônicas  
Nuvem Computacional  
Big Data  
Computação Pervasiva  
IoT  
Machine Learning  
etc.

# DATA MINING

**"Mostre-me dados, que eu lhe apresentarei padrões!"**

- Pré-requisitos:
  - Datasets bem estruturados e organizados
  - Formatos de dados bem definidos
  - Processo prévio intensivo de preparação dos dados
  - Ou os dados são fornecidos como se espera, ou então será necessário trabalhá-los

# DADOS *VERSUS* TEXTO

# DATA MINING

Dados  
Não-  
Estruturados

Dados  
Semi-  
Estruturados

Dados  
Estruturados

# TEXT MINING

CIÊNCIA DE DADOS

DATA MINING

Dados  
Não-  
Estruturados

Dados  
Semi-  
Estruturados

Dados  
Estruturados

TEXT MINING

DATA SCIENCE

# TEXT MINING: OUTPUTS

Texto Cru  
(raw text)

Text Mining

**Números**  
**Dados**  
**Insights**  
**Organização**  
**Seleção**  
**Classificação**  
**Combinação**  
**Informação**  
**Predição**  
**etc.**

# Etapas da Mineração

# ETAPAS DA MINERAÇÃO DE TEXTO

- Coleta do Corpus Textual
- Pré-processamento Textual (Processamento de Linguagem Natural)
- Processamento Textual (inclui a vetorização do texto)
- Pós-Processamento Textual

# Vetorização do Texto



# BAG OF WORDS

- O conjunto de pesos para termos de um documento é conhecido na literatura como **Modelo Bag of Words**, ou Saco de Palavras
- A essência do modelo Bag Of Words é converter documentos-texto em vetores que representem a **frequência das palavras distintas existentes nos documentos**

$$D = \{ wD_1, wD_2, \dots, wD_n \}$$

- onde  $w_{Dn}$  denota o peso da palavra  $n$  no documento  $D$

# FREQUÊNCIA DE TERMOS

- Term Frequency (*tf*)
- Um peso é associado para cada termo de um documento
- O peso depende do número de ocorrências do termo no documento
- A abordagem associa o peso à quantidade de ocorrências do termo
- Para o número de ocorrências de um termo *t* em determinado documento *d*, temos a seguinte representação

$$tf_{t, d}$$

# INVERSO DA FREQUÊNCIA NOS DOCUMENTOS

- A frequência de termos possui um grande problema: todos os termos são considerados equanimemente
- Isso pode se tornar um problema crítico, dependendo da aplicação trabalhada

# INVERSO DA FREQUÊNCIA NOS DOCUMENTOS

- Por exemplo, é possível que uma coleção de documentos da indústria automobilística ocorra muito frequentemente a ocorrência dos termos carro, automóvel e veículo
- Neste caso, é necessário se criar um mecanismo para atenuar esta situação, de maneira a **não superdimensionar a ocorrência destes termos** no modelo Bag of Words

# INVERSO DA FREQUÊNCIA NOS DOCUMENTOS

- Para este fim, costuma-se considerar a frequência de documentos  $df_t$ , ou seja, o número de documentos do corpus que contêm o termo  $t$
- Denotando-se por  $N$  o número de documentos de determinado corpus, o inverso da frequência de documentos ( $idf$ ) para um termo  $t$  é definido por

$$idf_t = \log \frac{N}{df_t}$$

# INVERSO DA FREQUÊNCIA NOS DOCUMENTOS

- Exemplo:

term	$df_t$	$idf_t$
car	18,165	1.65
auto	6723	2.08
insurance	19,241	1.62
best	25,235	1.5

## TF-IDF

- O *tfidf* é o produto das duas métricas anteriores, e pode ser definido por

$$tfidf = tf \times idf$$

- Este cálculo é capaz de estabelecer uma **ponderação aos termos presentes em um documento**, considerando **tanto o próprio documento quanto seu relacionamento com o corpus textual analisado**

# TF-IDF

- É importante ressaltar, no entanto, que embora o cálculo do *tfidf* crie as representações para os textos, ainda é preciso considerar que os documentos de um mesmo corpus podem possuir tamanhos diferentes
- Neste caso, para corrigir esta distorção, é comum a realização de um procedimento de normalização, dando origem à medida *tfidf normalizada*
- O cálculo do TF-IDF é utilizado no processo de vetorização para preencher a Bag of Words



# Exemplo de Vetorização

# EXEMPLO

- Seja  $\Psi = \{ doc_1, doc_2, \dots, doc_n \}$  um determinado corpus de documentos

- Vamos considerar a situação na qual  $n = 3$ , portanto,

$$\Psi = \{ doc_1, doc_2, doc_3 \}$$

<b>doc1</b>	Ambiente agradável e tranquilo. Comida e música com qualidade. Adoramos o Filé à Parmegiana.
<b>doc2</b>	O Filé à Parmegiana da cidade. Ambiente agradável e qualidade no atendimento.
<b>doc3</b>	O Filé à Parmegiana com fritas é uma delícia.

# EXEMPLO

- **Lista de tokens** obtida após a execução da análise lexical sobre o corpus

<b>doc1</b>	ambiente   agradável   e   tranquilo   comida   e   música   com   qualidade   adoramos   o   filé   à   parmegiana
<b>doc2</b>	o   filé   à   parmegiana   da   cidade   ambiente   agradável   e   qualidade   no   atendimento
<b>doc3</b>	o   filé   à   parmegiana   com   fritas   é   uma   delícia

# EXEMPLO

- Lista de tokens obtida **após a remoção das stopwords**

<b>doc1</b>	ambiente   agradável   tranquilo   comida   música   qualidade   adoramos   filé   parmegiana
<b>doc2</b>	filé   parmegiana   cidade   ambiente   agradável   qualidade   atendimento
<b>doc3</b>	filé   parmegiana   fritas   delícia

# EXEMPLO

- Lista resultante após aplicação da **radicalização**

<b>doc1</b>	ambient   agrad   tranquil   com   músic   qualidad   ador   fil   parmegian
<b>doc2</b>	fil   parmegiana   cidad   ambient   agrad   qualidad   atend
<b>doc3</b>	fil   parmegian   frit   delíci

# EXEMPLO

- **Dicionário de termos** obtidos

termo 1	ador	termo 8	fil
termo 2	agrad	termo 9	frit
termo 3	ambient	termo 10	músic
termo 4	atend	termo 11	parmegian
termo 5	ciudad	termo 12	qualidad
termo 6	com	termo 13	tranquil
termo 7	delíc		

# EXEMPLO

- Conjunto de dados  $\Psi$  considerando uma **representação binária**

	ador	agrad	ambient	atend	cidad	com	delíc	fil	frit	músic	parmegian	qualidad	tranquil
$\Psi$	$wte_1$	$wte_2$	$wte_3$	$wte_4$	$wte_5$	$wte_6$	$wte_7$	$wte_8$	$wte_9$	$wte_{10}$	$wte_{11}$	$wte_{12}$	$wte_{13}$
$doc_1$	1	1	1	0	0	1	0	1	0	1	1	1	1
$doc_2$	0	1	1	1	1	0	0	1	0	0	1	1	0
$doc_3$	0	0	0	0	0	0	1	1	1	0	1	0	0

# EXEMPLO

- Conjunto de dados  $\Psi$  considerando uma **representação *tfidf***

	ador	agrad	ambient	atend	ciudad	com	delíc	fil	frit	músic	parmegian	qualidad	tranquil
$\Psi$	$wte_1$	$wte_2$	$wte_3$	$wte_4$	$wte_5$	$wte_6$	$wte_7$	$wte_8$	$wte_9$	$wte_{10}$	$wte_{11}$	$wte_{12}$	$wte_{13}$
$doc_1$	1,58	0,58	0,58	0,00	0,00	1,58	0,00	0,00	0,00	1,58	0,00	0,58	1,58
$doc_2$	0,00	0,58	0,58	1,58	1,58	0,00	0,00	0,00	0,00	0,00	0,00	0,58	0,00
$doc_3$	0,00	0,00	0,00	0,00	0,00	0,00	1,58	0,00	1,58	0,00	0,00	0,00	0,00



# EXEMPLO

- Conjunto de dados  $\Psi$  considerando uma **representação *tfidf* normalizada**

	ador	agrad	ambient	atend	cidad	com	delíc	fil	frit	músic	parmegian	qualidad	tranquil
$\Psi$	$wte_1$	$wte_2$	$wte_3$	$wte_4$	$wte_5$	$wte_6$	$wte_7$	$wte_8$	$wte_9$	$wte_{10}$	$wte_{11}$	$wte_{12}$	$wte_{13}$
$doc_1$	0,18	0,06	0,06	0,00	0,00	0,18	0,00	0,00	0,00	0,18	0,00	0,06	0,18
$doc_2$	0,00	0,08	0,08	0,23	0,23	0,00	0,00	0,00	0,00	0,00	0,00	0,08	0,00
$doc_3$	0,00	0,00	0,00	0,00	0,00	0,00	0,40	0,00	0,40	0,00	0,00	0,00	0,00

# Text Mining: Aplicações

# TEXT MINING: APLICAÇÕES

- Classificação de Documentos
- Clusterização (ou Categorização) de Documentos
- Sumarização de Documentos
- Similaridade de Documentos
- Reconhecimento de Entidades Nomeadas (ou Mencionadas)
- Análise de Sentimentos (Mineração de Opiniões)
- Sistemas de Pergunta e Resposta
- Chatbots
- entre outras

# Desafios e Oportunidades

# DESAFIOS

- Mineração de texto em outras línguas que não o inglês
- Padrões para os mineradores em língua portuguesa
- Relações léxico-semânticas

# OPORTUNIDADES

- Muitos softwares disponíveis
- Pacotes, "bibliotecas", APIs (Application Programming Interfaces)
- Poucas linhas de código em Python, Java ou R, por exemplo
- Percorrimento da curva de aprendizagem é bastante rápido
- Cinema, publicação científica e inúmeras outras áreas de aplicação

# MATERIAL CONSULTADO

- BUDER, Emily. **Can AI Predict a Movie's Success? Algorithmic Screenplay Service 'Scriptbook' Causes Major Backlash.** Disponível em: <<https://nofilmschool.com/2017/04/scriptbook-black-list-screenwriting-ai-algorithm>>. Acesso em: 01 maio 2018.
- GLADWELL, Malcolm. **The Formula:** What if you built a machine to predict hit movies?. Disponível em: <<https://www.newyorker.com/magazine/2006/10/16/the-formula>>. Acesso em: 01 maio 2018.
- HALL, Jacob. **WTF:** ScriptBook Will Predict If a Screenplay Will Make Money; Will Also Possibly, Maybe Trigger Nuclear Armageddon. Disponível em: <<http://www.slashfilm.com/scriptbook-software/>>. Acesso em: 01 maio 2018.
- WIKIPÉDIA. **Epagogix.** Disponível em: <<https://en.wikipedia.org/wiki/Epagogix>>. Acesso em: 01 maio 2018.
- WESTBROOK, Adam. **Robots are reading Hollywood scripts – and writers aren't happy.** Disponível em: <<https://www.thememo.com/2017/04/20/scriptbook-blacklist-robots-reading-hollywood-scripts-writers-not-happy/>>. Acesso em: 01 maio 2018.
- WIKIPÉDIA. **Gigli.** Disponível em: <<https://pt.wikipedia.org/wiki/Gigli>>. Acesso em: 01 maio 2018.
- WIKIPÉDIA. **The Hero with a Thousand Faces.** Disponível em: <[https://pt.wikipedia.org/wiki/The\\_Hero\\_with\\_a\\_Thousand\\_Faces](https://pt.wikipedia.org/wiki/The_Hero_with_a_Thousand_Faces)>. Acesso em: 01 maio 2018.
- WILCOCK, David. **The Synchronicity Key:** The Hidden Intelligence Guiding the Universe and You. New York: Dutton - Est. 1852, 2013. 528 p.
- PRESSBERG, Matt. **Swimming With Algorithms:** Can Software Predict The Next Hollywood Hit Based On The Script Alone?. Disponível em: <<http://www.ibtimes.com/swimming-algorithms-can-software-predict-next-hollywood-hit-based-script-alone-2266229>>. Acesso em: 01 maio 2018.

# MATERIAL CONSULTADO

- APPLE INC. **Dicionário**. 2. ed. Cupertino (CA), EUA: macOS High Sierra, 2017.
- DIAS, Maria; MALHEIROS Marcelo. (2005). **Estudo de técnicas de radicalização para a Língua Portuguesa**. Disponível em: <[https://www.researchgate.net/publication/242193490\\_Estudo\\_de\\_tecnicas\\_de\\_radicalizacao\\_para\\_a\\_Lingua\\_Portuguesa](https://www.researchgate.net/publication/242193490_Estudo_de_tecnicas_de_radicalizacao_para_a_Lingua_Portuguesa)>. Acesso em 20 mar. 2018.
- GANTZ, John; REINSEL, David. Extracting value from chaos. IDC iView, v. 1142, n. 2011, p. 1-12, 2011. Disponível em: <<http://www.emc.com/collateral/analyst-reports/idc-extracting-value-from-chaos-ar.pdf>>. Acesso em: 27 nov. 2018.
- INGERSOLL, Grant S.; MORTON, Thomas S.; FARRIS, Andrew L.. **Taming Text: How to find, organize and manipulate it**. Shelter Island, NY (USA): Manning Publications Co., 2013. 298 p.
- MANNING, Christopher D.; RAGHAVAN, Prabhakar; SCHÜTZE, Hinrich. **Introduction to Information Retrieval**, Cambridge (UK): Cambridge University Press. 2008.
- NEVES, P. I. ; CORREA, D. A. ; GOLDSCHMIDT, R. R. ; MOURA, Ana Maria Carvalho ; CAVALCANTI, M. C. . UMA ANÁLISE SOBRE ABORDAGENS E FERRAMENTAS PARA EXTRAÇÃO DE INFORMAÇÃO. **Revista Militar de Ciência e Tecnologia** , v. 30, p. 32-58, 2013.
- SARKAR, Dipanjar. **Text Analytics with Python: A Practical Real-World Approach to Gaining Actionable Insights from Your Data**. Middletown, DE (USA): Apress, 2016. 385 p.
- SOUSA, Afonso F. (Org.). **Máximas e Mínimas do Barão de Itararé**. Rio de Janeiro: Record, 1985.
- WEISS, Sholom M.; INDURKHYA, Nitin; ZHANG, Tong. **Fundamentals of Predictive Text Mining**. New York: Springer, 2010. 226 p. (Texts in Computer Science).



OBRIGADO!