

Defesa de Mestrado

Othon Luiz Teixeira de Oliveira

Escola Politécnica de Pernambuco – Poli — UPE

29 - Maio - 2017

PPGES



Programa de Pós-Graduação em Engenharia de Sistemas

DISSERTAÇÃO DE MESTRADO

Programa de Pós-Graduação em Engenharia de Sistemas

DISSERTAÇÃO DE MESTRADO

Modelo Preditivo para Sugestão de Roteamento de cargas considerando dados históricos, sócio-ambientais, e de redes sociais

Programa de Pós-Graduação em Engenharia de Sistemas

DISSERTAÇÃO DE MESTRADO

Modelo Preditivo para Sugestão de Roteamento de cargas considerando dados históricos, sócio-ambientais, e de redes sociais

Mestrando: Othon Luiz Teixeira de Oliveira

Orientador: Prof. Dr. Fernando Buarque de Lima Neto

Sumário

- 1 **Introdução**
 - Introdução
 - Objetivos

2 Background

3 Contribuição

4 Resultados

5 Conclusões

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

Resumo

O cenário

- O transporte de cargas no Brasil é feito principalmente pelas rodovias federais (BRs).

Resumo

O cenário

- O transporte de cargas no Brasil é feito principalmente pelas rodovias federais (BRs).
- Essas rodovias estão constantemente congestionadas nos perímetros urbanos.

Resumo

O cenário

- O transporte de cargas no Brasil é feito principalmente pelas rodovias federais (BRs).
- Essas rodovias estão constantemente congestionadas nos perímetros urbanos.
- Comunidades bloqueiam as rodovias para reivindicar, dos entes públicos, todo tipo de necessidades

Resumo

O cenário

- O transporte de cargas no Brasil é feito principalmente pelas rodovias federais (BRs).
- Essas rodovias estão constantemente congestionadas nos perímetros urbanos.
- Comunidades bloqueiam as rodovias para reivindicar, dos entes públicos, todo tipo de necessidades
- Em alguns trechos o traçado das rodovias está próximo a morros e florestas.

Organização

1 Introdução

- Introdução
- **Objetivos**

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

Proposição de um modelo preditivo

Proposição de um modelo preditivo

Objetivo Geral

Esta pesquisa teve como principal objetivo desenvolver um modelo de plataforma autoadaptável

Proposição de um modelo preditivo

Objetivo Geral

Esta pesquisa teve como principal objetivo desenvolver um modelo de plataforma autoadaptável que contemple predição do comportamento das rodovias federais pernambucanas,

Proposição de um modelo preditivo

Objetivo Geral

Esta pesquisa teve como principal objetivo desenvolver um modelo de plataforma autoadaptável que contemple predição do comportamento das rodovias federais pernambucanas, antecipando alguns eventos que nela possam ocorrer, apontando onde ocorrerão.

Objetivos Específicos

Objetivos Específicos

Quatro objetivos

- Caracterizar a problemática de cada rodovia;

Objetivos Específicos

Quatro objetivos

- Caracterizar a problemática de cada rodovia;
- Desenvolver um modelo preditivo dos fenômenos que envolvem as rodovias;

Objetivos Específicos

Quatro objetivos

- Caracterizar a problemática de cada rodovia;
- Desenvolver um modelo preditivo dos fenômenos que envolvem as rodovias;
- Desenvolver um ambiente de simulações interativas da estrutura viária em sua dinâmica

Objetivos Específicos

Quatro objetivos

- Caracterizar a problemática de cada rodovia;
- Desenvolver um modelo preditivo dos fenômenos que envolvem as rodovias;
- Desenvolver um ambiente de simulações interativas da estrutura viária em sua dinâmica
- Propor soluções para melhorar a experiência dos usuários que utilizam as rodovias pernambucanas

Sumário

1 Introdução

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição



Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

Dados originais da PRF

- Base de dados da Polícia Rodoviária Federal: acidente e interdições, entre 2007 a 2015.

Dados originais da PRF

- Base de dados da Polícia Rodoviária Federal: acidente e interdições, entre 2007 a 2015.
- Rede Social Twitter com 3200 tweets (limite permitido), de 2017 a 2014

Planilha com os dados originais

Planilha PRF

Figura: Planilha para Preprocessamento

Table "default" - Rows: 85209			Spec - Columns: 22		Properties	Flow Variables			
Row ID	BR	KM	\$ Latitude	\$ Longitude	\$ Condi...	\$ Restri...	\$ TipoAcid...	\$ CausaAcide...	\$
Row0	101	79	-7.56256206749...	-34.990451894499...	Seca	Inexistente	Colisão late...	Falta de atenç...	De
Row1	101	210	-7.66697808649...	-34.933297438999...	Seca	Inexistente	Colisão Tra...	Outras	De
Row2	101	416	-7.83434453699...	-34.912533977999...	Seca	Inexistente	Colisão Tra...	Falta de atenç...	Cre
Row3	101	416	-7.83434453699...	-34.912533977999...	Seca	Inexistente	Colisão Tra...	Falta de atenç...	Cre
Row4	101	474	-7.86058770949...	-34.908337135499...	Seca	Inexistente	Colisão com...	Outras	Cre
Row5	101	511	-7.90212986999...	-34.900132434499...	Seca	Inexistente	Atropelame...	Outras	Cre
Row6	101	512	-7.91819106999...	-34.896062575499...	Seca	Inexistente	Colisão Tra...	Falta de atenç...	Cre
Row7	101	512	-7.91819106999...	-34.896062575499...	Seca	Inexistente	Colisão Tra...	Falta de atenç...	Cre
Row8	101	545	-7.94363545599...	-34.904202790499...	Seca	Inexistente	Colisão late...	Falta de atenç...	Cre
Row9	101	545	-7.94363545599...	-34.904202790499...	Seca	Inexistente	Colisão late...	Falta de atenç...	Cre
Row10	101	589	-7.97661785749...	-34.925137692499...	Seca	Inexistente	Colisão late...	Velocidade inc...	De
Row11	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Colisão late...	Falta de atenç...	Cre
Row12	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Colisão late...	Falta de atenç...	Cre
Row13	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Colisão late...	Ultrapassagem...	Cre
Row14	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Colisão late...	Ultrapassagem...	Cre
Row15	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Colisão tras...	Falta de atenç...	Cre
Row16	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Colisão tras...	Falta de atenç...	Cre
Row17	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Queda de m...	Outras	Cre
Row18	101	680	-8.87953132499...	-35.628922247499...	Seca	Inexistente	Queda de m...	Velocidade inc...	De
Row19	101	830	-8.17277505899...	-34.939988807999...	Seca	Inexistente	Colisão late...	Falta de atenç...	De
Row20	101	830	-8.17277505899...	-34.939988807999...	Seca	Inexistente	Colisão tras...	Falta de atenç...	De
Row21	101	830	-8.17277505899...	-34.939988807999...	Seca	Inexistente	Colisão tras...	Falta de atenç...	De
Row22	101	830	-8.17277505899...	-34.939988807999...	Seca	Inexistente	Queda de m...	Outras	De
Row23	101	865	-8.18644581449...	-34.971397555499...	Seca	Inexistente	Colisão tras...	Falta de atenç...	De
Row24	101	1018	-8.26668917249...	-35.041692694499...	Seca	Inexistente	Colisão late...	Falta de atenç...	Cre
Row25	101	1018	-8.26668917249...	-35.041692694499...	Seca	Inexistente	Colisão late...	Falta de atenç...	Cre
Row26	101	1020	-8.28279581549...	-35.057672119999...	Com buro...	Inexistente	Colisão late...	Outras	Cre
Row27	101	1020	-8.28279581549...	-35.057672119999...	Com buro...	Inexistente	Colisão late...	Outras	Cre

Entendendo a base de dados da PRF

Atributos e Instâncias – iniciais

- Dados de 2007 à 2015 – BRs de Pernambuco

Entendendo a base de dados da PRF

Atributos e Instâncias – iniciais

- Dados de 2007 à 2015 – BRs de Pernambuco
- 85.209 Instâncias – 27 Atributos

Entendendo a base de dados da PRF

Atributos e Instâncias – iniciais

- Dados de 2007 à 2015 – BRs de Pernambuco
- 85.209 Instâncias – 27 Atributos
- Dentre eles: Km, Latitude, Longitude, Condições da Pista, Causa do Acidente, Município, Data, Hora, Tipo de Veículo, Quantidade de Mortos, etc.

Entendendo a base de dados da PRF

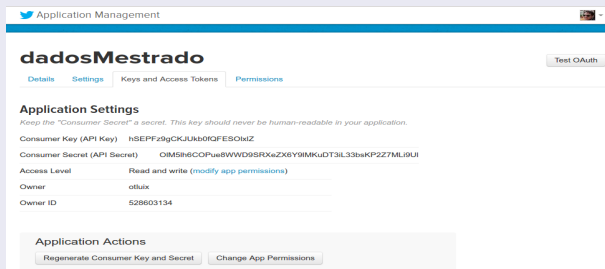
Atributos e Instâncias – iniciais

- Dados de 2007 à 2015 – BRs de Pernambuco
- 85.209 Instâncias – 27 Atributos
- Dentre eles: Km, Latitude, Longitude, Condições da Pista, Causa do Acidente, Município, Data, Hora, Tipo de Veículo, Quantidade de Mortos, etc.
- 40% missing data – 70% do tempo total para tratar

O Tweeker

Acesso aos dados do Tweeker

Figura: Registro de App para para acessar e baixar dados



Entendendo a base de dados do Tweeter

Atributos e Instâncias – iniciais

- Dados de 2017 – 2014 — canal @PRF191PE

Entendendo a base de dados do Tweeter

Atributos e Instâncias – iniciais

- Dados de 2017 – 2014 — canal @PRF191PE
- 2864 Instâncias – 16 Atributos

Entendendo a base de dados do Tweeter

Atributos e Instâncias – iniciais

- Dados de 2017 – 2014 — canal @PRF191PE
- 2864 Instâncias – 16 Atributos
- Dentre eles: 'text', 'favorited', 'favoriteConunt', 'created', 'ID', 'statusSource', 'screenName', 'retweetCount',

Entendendo a base de dados do Tweeter

Atributos e Instâncias – iniciais

- Dados de 2017 – 2014 — canal @PRF191PE
- 2864 Instâncias – 16 Atributos
- Dentre eles: 'text', 'favorited', 'favoriteConunt', 'created', 'ID', 'statusSource', 'screenName', 'retweetCount',
- 'isRetweet', 'retweeted', dentre outros.

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

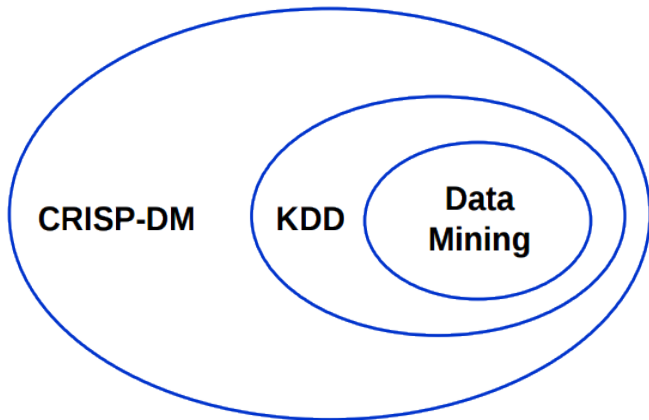
3 Contribuição

4 Resultados

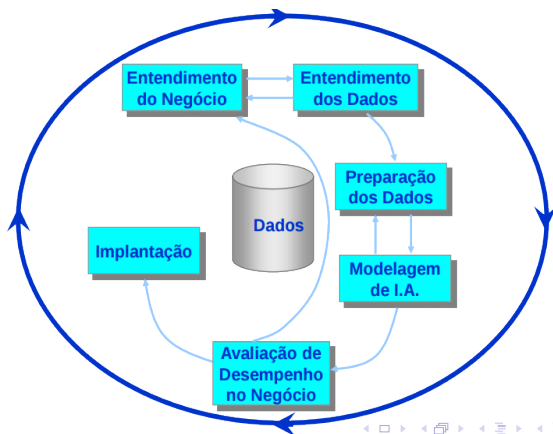
5 Conclusões

6 Extras

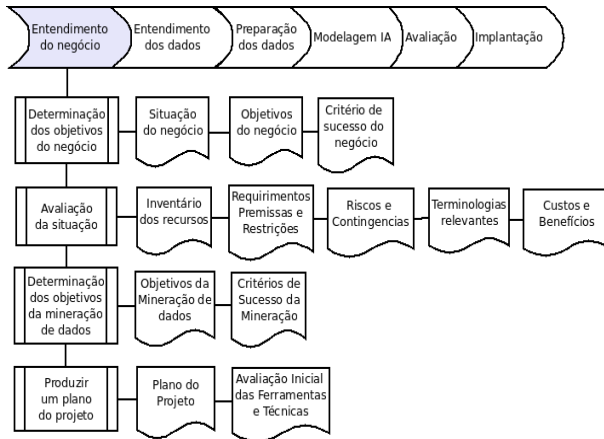
Domínio das técnicas de Mineração de Dados



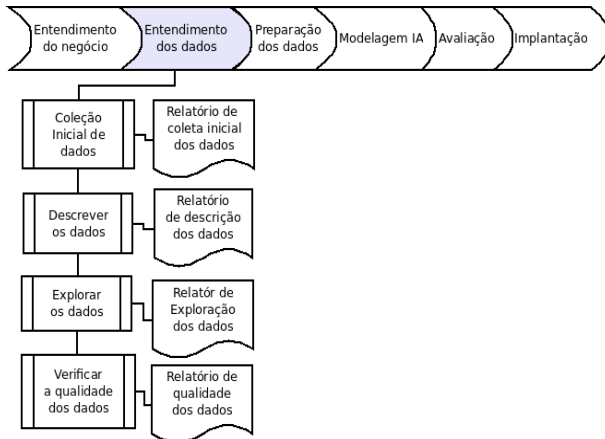
Cross Industry Standard Process for Data Mining – CRISP - DM



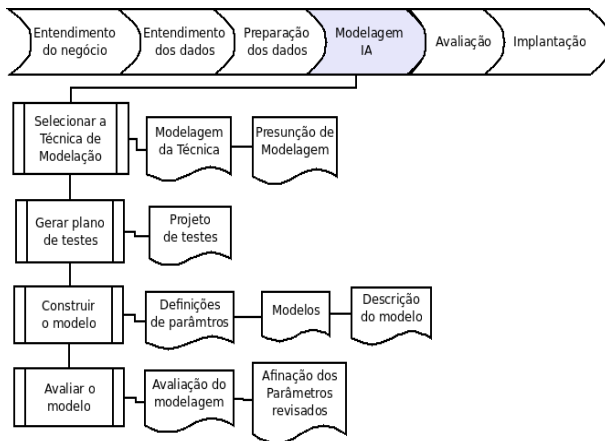
Etapas CRISP-DM resumo



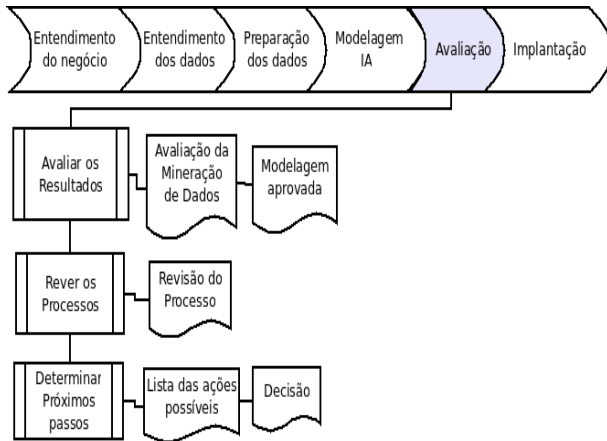
Etapas CRISP-DM resumo



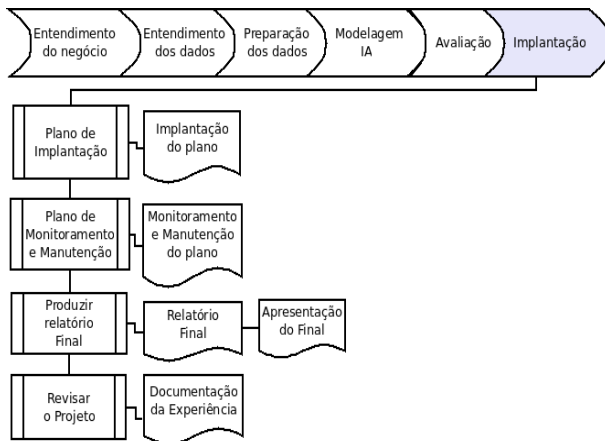
Etapas CRISP-DM resumo



Etapas CRISP-DM resumo



Etapas CRISP-DM resumo



Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- **Minerando os dados do problema**
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

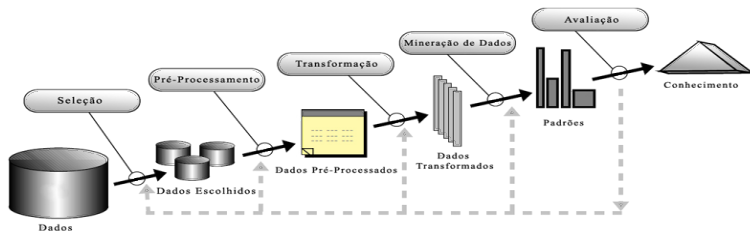
3 Contribuição

4 Resultados

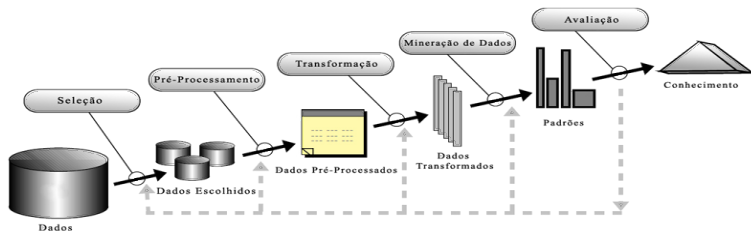
5 Conclusões

6 Extras

A descoberta de conhecimento (KDD) nas bases da PRF



A descoberta de conhecimento (KDD) nas bases da PRF



Diferença CRISP-DM — KDD

O CRISP-DM difere do KDD principalmente pelas fases do entendimento do negócio (anterior ao KDD) e da implantação (posterior ao KDD)

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- **Árvore de Decisão**
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

A escolha

Pontos fortes

- Custo benefício – Uma árvore é gerada de maneira simples com resultados abrangente e facilidade de interpretação.

A escolha

Pontos fortes

- Custo benefício – Uma árvore é gerada de maneira simples com resultados abrangente e facilidade de interpretação.
- Permite fazer predição e classificação instantâneas

A escolha

Pontos fortes

- Custo benefício – Uma árvore é gerada de maneira simples com resultados abrangente e facilidade de interpretação.
- Permite fazer predição e classificação instantâneas

Pontos fracos

- Nó raiz não é facilmente identificável

A escolha

Pontos fortes

- Custo benefício – Uma árvore é gerada de maneira simples com resultados abrangente e facilidade de interpretação.
- Permite fazer predição e classificação instantâneas

Pontos fracos

- Nó raiz não é facilmente identificável
- Vários testes até se conseguir resultados satisfatórios

Cálculo da Entropia e Entropia Condicional

Equação da entropia

$$H_x = - \sum_{\forall x \in X} P(x) \log_2 P(x) \quad (1)$$

Equação da entropia condicional

$$H_{Y|X} = \sum_x P(x) H(Y|X = x) = - \sum_{\forall x \in X} P(x) \sum_{\forall y \in Y} P(y|x) \log_2 P(y|x) \quad (2)$$

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- **Naïve Bayes**
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

Teorema de Bayes e Aprendizagem bayesiana

Teorema de Bayes: forma tradicional

$$p(C_k|x) = \frac{p(k)p(x|C_k)}{p(x)} \quad (3)$$

Teorema de Bayes: forma simplificada

$$p(\textit{posteriori}) = \frac{p(\textit{priori}) * \textit{verossimilhança}}{\textit{evidência}} \quad (4)$$

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- **Redes Neurais**
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

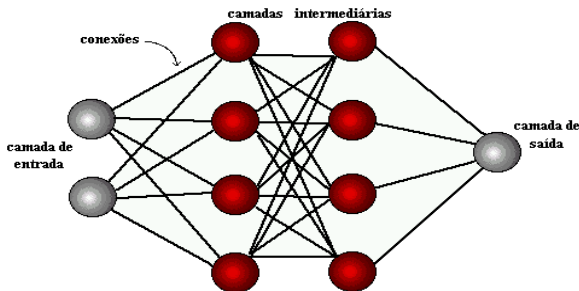
5 Conclusões

6 Extras

Classificação com Redes Neurais

Rede Neural com três camadas identificadas

- Camada de entrada: apresenta-se os padrões à rede
- Camada intermediária (ocultas): realiza a maior parte do processamento por conexões ponderadas
- Camada de saída: conclusão e apresentação do resultado.



Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

TF – IDF

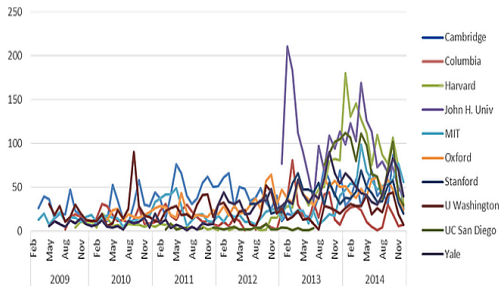
Text Frequency X Inverse Document Frequency: $TF * IDF$

$$idf(term) = \ln\left(\frac{n_{documents}}{n_{documents \text{ containing term}}}\right) \quad (5)$$

TF – IDF

Text Frequency X Inverse Document Frequency: $TF * IDF$

$$idf(term) = \ln\left(\frac{n_{documents}}{n_{documents \text{ containing term}}}\right) \quad (5)$$



Análise de Desempenho aplicados à mineração

A área sob a curva AUC (Area under ROC – Receiver Operating Characteristic) é uma métrica que determina a qualidade do classificador. Ela é calculada cruzando-se verdadeiros positivos com os falso positivos. Esta área varia entre zero e um, quanto mais próximo de 1(um) o classificador conseguiu acertar mais vezes do que errar. A seguir a Matriz de Contingência (ou confusão) que agrupa essas atributos.

Tabela: Matriz de Confusão

	Predito	
Real	TP FN	Positive – POS
Real	FP TN	Negative – NEG
—	PP PN	—

Análise de Desempenho aplicados à mineração

A Matriz Modelo de Confusão sintetiza a Matriz de Confusão

Tabela: Matriz modelo de Confusão

	Y	\bar{Y}	
X	$P(X, Y)$	$P(X, \bar{Y})$	Positive – POS
\bar{X}	$P(\bar{X}, Y)$	$P(\bar{X}, \bar{Y})$	Negative – NEG
—	$P(Y)$	$P(\bar{Y})$	—

Fonte: Bradley – 1997

Para construir a curva ROC utiliza-se as probabilidades condicionais cruzando-se a taxa de verdadeiros positivos ($tpr = P(Y|X)$) probabilidade de falsos alarmes ou taxa de falsos positivos será ($fpr = P(Y, \bar{X})$),

Sumário

- 1 Introdução
- 2 Background
- 3 Contribuição**
- 4 Resultados
- 5 Conclusões
- 6 Extras

Sumário

- 1 Introdução
- 2 Background
- 3 Contribuição
- 4 Resultados**
- 5 Conclusões
- 6 Extras

Sumário

- 1 Introdução
- 2 Background
- 3 Contribuição
- 4 Resultados
- 5 Conclusões**
- 6 Extras

Sumário

- 1 Introdução
- 2 Background
- 3 Contribuição
- 4 Resultados
- 5 Conclusões
- 6 Extras

- Trabalhos futuros
- Artigos aprovados em eventos

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

Aplicativo

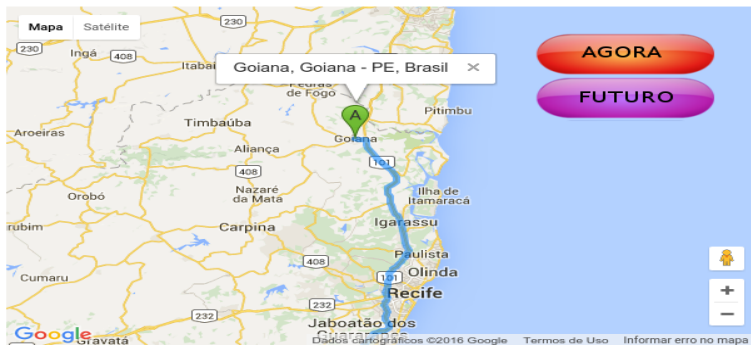


Figura: Interface Gráfica Proposta

Organização

1 Introdução

- Introdução
- Objetivos

2 Background

- Dados a minerar
- Metodologias de Mineração Dados — CRISP – DM
- Minerando os dados do problema
- Árvore de Decisão
- Naïve Bayes
- Redes Neurais
- Tecnologias empregues na Mineração em Textos

3 Contribuição

4 Resultados

5 Conclusões

6 Extras

Artigo aprovado em evento internacional

ANÁLISE DE COMPORTAMENTO RODOVIÁRIO NA REGIÃO METROPOLITANA DO RECIFE-BRASIL – UMA ABORDAGEM DE MINERAÇÃO DE DADOS



12^a Conferência Ibérica de
Sistemas e Tecnologias
de Informação

21 a 24
JUNHO
2017
Lisboa
Portugal

Figura: CISTI – Conf. Ibérica de Sistemas e Tecnologia da Informação