



UNIVERSIDADE DO MINHO
Mestrado em Engenharia Informática
SBS

Trabalho Prático Nº1

Pedro Ribeiro - PG42848

Professor Doutor César Analide Freitas Silva Costa Rodrigues
Professor Doutor Bruno Filipe Martins Fernandes

Mestrado em Engenharia Informática
2021

Conteúdo

1	Introdução	1
1.1	Conjuntos de dados	1
2	Metodologia	2
2.1	<i>CRISP-DM</i>	2
3	Exploração e modelação	3
3.1	Incidentes em Braga - 2019	3
3.2	Banco de Portugal - Campanhas de Marketing	4
4	Workflows Desenvolvidos	5
4.1	Incidentes Braga - 2019	5
4.1.1	Descrição do Tuning	7
4.2	Banco de Portugal - Campanhas de Marketing	7
4.2.1	Descrição do Tuning	8
5	Conclusões	10
	Bibliografia	11

Capítulo 1

Introdução

No âmbito da unidade curricular "Sistemas Baseados em Similaridade", do perfil "Machine Learning: Fundamentos e Aplicações" foi desenvolvido este trabalho prático que se foca na exploração de dois datasets, um fornecido pelos docentes da UC e outro escolhido pelo autor do trabalho prático.

1.1 Conjuntos de dados

O dataset fornecido consta com os registos de incidentes rodoviários na cidade de Braga no período de 2019, contendo dados como a magnitude de atraso causado perante as estradas afetadas, o atraso em segundos provocado pelo incidente, entre outros.

O objetivo é prever a magnitude de incidentes de cada registo com o resultado possível associado, ou seja, que englobe o mínimo de erro possível na previsão,

O segundo dataset foi criado por Sérgio Moro (ISCTE-IUL), Paulo Cortez (Univ. Minho) e Paulo Rita (ISCTE-IUL) em 2014, e consta com informações sobre campanhas de marketing de um Banco de Portugal, conduzidas particularmente através de chamadas de telemóvel.[2] Este conjunto de dados consta com um total de 21 colunas, com informações pessoais dos inquiridos, dados relativos à campanha no ato do registo e dados económicos como por exemplo a taxa de euribor naquele momento.

Capítulo 2

Metodologia

2.1 *CRISP-DM*

Para alcançar uma qualidade nos modelos produzidos foi utilizada a metodologia *Cross Industry Standard Process for Data Mining*, ou *CRISP-DM*. Publicado em 1999, tornou-se a metodologia mais comum para *data mining* e projetos de *data science*, até ao dia 30 de junho de 2020 este foi o modelo mais pesquisado no google por uma margem bastante grande.[1]

O CRISP-DM é dividido em 6 fases:

- **Compreensão do modelo de negócio**, qualquer projeto começa com uma compreensão profunda das necessidades do cliente. Projetos de mineração de dados não são exceção e o CRISP-DM reconhece isso.
- **Compreensão dos dados**, depois da compreensão do modelo de negócio o CRISP-DM direciona o foco para identificar, coletar e analisar o conjunto de dados que podem ajudar a cumprir os objetivos do projeto.
- **Preparação dos dados**, uma regra comum é que 80% do projeto é preparação de dados, esta tarefa consiste em selecionar, limpar, construir, integrar e formatar os dados.
- **Modelação**, o que é considerado o trabalho mais empolgante de *data science* também costuma ser a fase mais curta do projeto, nesta fase consiste em construir e avaliar vários modelos com base em várias técnicas de modelagem diferentes.
- **Avaliação dos resultados**, esta fase analisa de forma mais ampla qual modelo que apresenta os melhores resultados e são definidos os próximos passos.
- **Implementação**, dependendo dos requisitos, a fase de implantação pode ser tão simples quanto gerar um relatório ou tão complexa quanto implementar um processo de mineração de dados.

Capítulo 3

Exploração e modelação

3.1 Incidentes em Braga - 2019

O conjunto de dados fornecidos pelos docentes da UC vesam para os incidentes rodoviários na cidade de Braga no ano de 2019, onde consta com as seguintes colunas e o respectivo tratamento de dados na respetiva coluna:

- city_name
 - Retirado devido a ser uma variável com pouca importância,o seu valor é sempre a respetiva cidade de Braga.
- magnitude.of.delay
 - devido a baixa incidência de casos "*Minor*" e "*UNKNOWN_DELAY*" estas foram mapeadas para um caso moderado
- delay_in.seconds
- affected.roads
 - Foi transformada em uma Lista e de seguida foi executado uma contagem do numero total de ruas afectadas.
- record.date
 - foi mapeada para um campo de data e hora e de seguida apenas foi extraído o numero do mês, o dia do ano , o numero do dia da semana e a Hora; de seguida o campo é retirado por já ter sido extraído toda a informação necessária
- luminosity
- avg.temperature
- avg.humidity
- avg.wind.speed
- avg.atm.pressure, avg.precipitation, avg.rain
 - Retidados por falta de importância para o modelo
- accidents

Modelação

3.2 Banco de Portugal - Campanhas de Marketing

O conjunto de dados escolhido pelo autor do trabalho incide em resultados de campanhas de marketing do **Banco de Portugal**, onde estas são resultantes principalmente através de chamadas telefónicas, oferecendo ao cliente um termo de depósito. A variável alvo categórica de decisão consta com 2 hipóteses: "yes", "no

. Este conjunto de dados foi recolhido através da plataforma online *Kaggle*, uma comunidade focada na resolução de problema de *data science* e de *Machine Learning*. [2]

Este conjunto consta com as seguinte colunas e os seus respetivos pré-processamentos:

- age
- day_of_week
- duration
- campaign
- emp.var.rate: employment variation rate
- cons.price.idx: consumer price index
- cons.conf.idx: consumer confidence index
- euribor3m: euribor 3 month rate
- nr.employed
- education, marital, job
 - tratamento dos valores "Unknown" e a aplicação do *One Hot Encoding*
- y, loan, housing, contact, month, default, previous, pdays, poutcome
 - transformação dos valores nominais para valores numéricos

Modelação

Capítulo 4

Workflows Desenvolvidos

4.1 Incidentes Braga - 2019

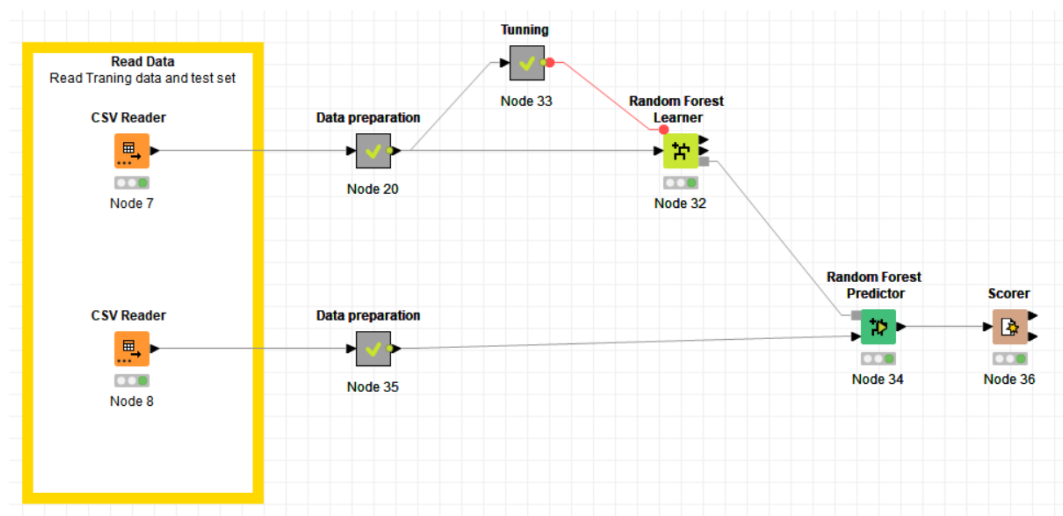


Figura 4.1: Visão geral do workflow gerado

- Read Data:
leitura do dataset de treino e do dataset de teste nos *NODE 7 e 8*
- Data preparation
Os *NODE 20 e 35* consistem no mesmo processamento de dados porque a estrutura dos dois conjuntos é a mesma

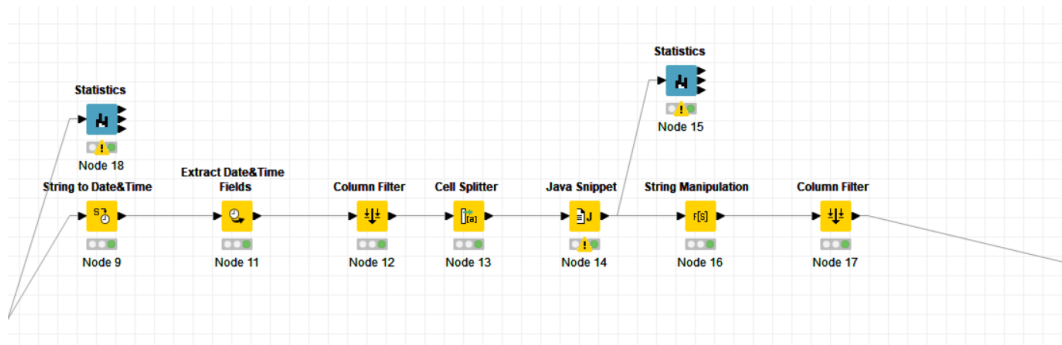


Figura 4.2: Pre-processamento

- Tuning

De forma a obter o melhor modelo possível é executado um tuning para recorrer a loops em que estes testam varias combinações possíveis entre os critérios de divisão de Random Forests e múltiplos valores para o número de modelos e *TreeDepth*. Desta forma é obtido melhor modelo possível perante as opções de modificações

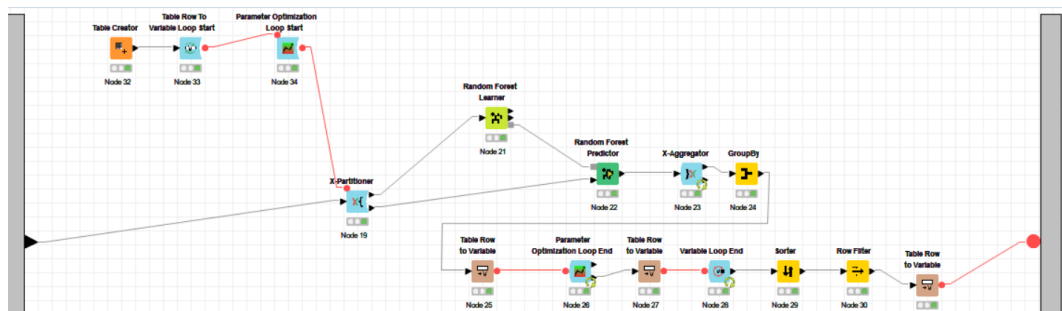


Figura 4.3: Tuning

- Previsões

Após as otimizações é feito o treino utilizando uma Random Forest Learner e o respectivo teste do conjunto de teste através do Random Forest Predictor.

Através do *Scorer* é obtido a matriz de confusão que devolve como as combinações entre Falso e Positivos relativamente ao conjunto de teste, podemos observar na figura 4.4.

Confusion Matrix - 3:36 - Scorer			
luminosity ...	DARK	LIGHT	LOW_LIGHT
DARK	2347	0	0
LIGHT	0	2437	0
LOW_LIGHT	0	0	216
Correct classified: 5.000			
Wrong classified: 0			
Accuracy: 100 %			
Error: 0 %			
Cohen's kappa (κ) 1			

Figura 4.4: Matriz de confusão conjunto de teste

4.1.1 Descrição do Tuning

Após o loop de optimização foram obtidos os seguintes valores (figura 4.5):

i	maxIterations	10
i	currentIteration	9
i	Loop-Execute	
i	Loop (0)	
i	iteration	24
i	depth	30
i	NumMaxModels	800
i	Loop-Execute	
i	Loop (1)	
s	Criteria	Gini
s	RowID	Row2
i	currentIteration	2
i	maxIterations	3
i	Loop-Execute	
i	Loop (2)	
s	knime.workspace	D:\Github\SBS

Figura 4.5: Matriz de confusão conjunto de teste

4.2 Banco de Portugal - Campanhas de Marketing

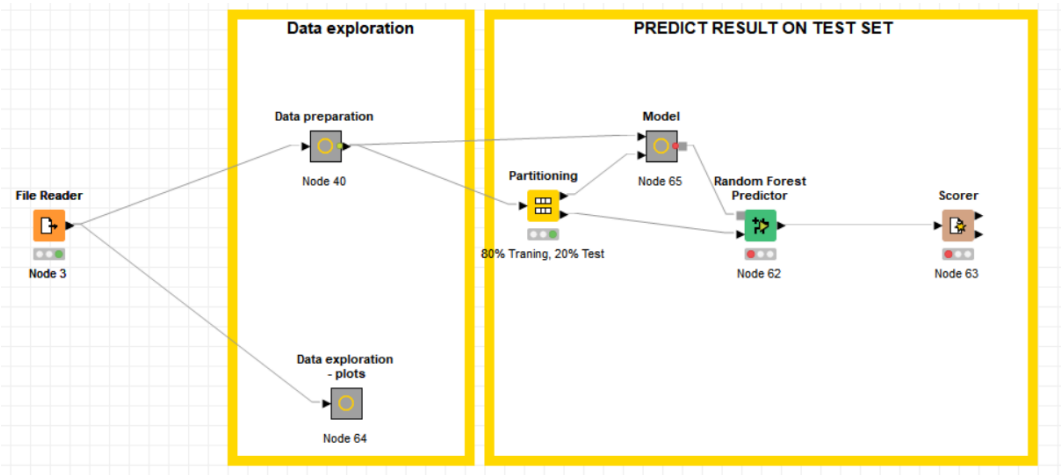
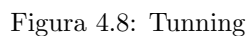


Figura 4.6: Visão geral do workflow gerado

- Read Data:
leitura do dataset de treino e do dataset de teste nos *NODE 7* e *8*
- Data preparation
Os *NODE 20* e *35* consistem no mesmo processamento de dados porque a estrutura dos dois conjuntos é a mesma



De forma a obter o melhor modelo possível é executado um tuning para recorrer a loops em que estes testam varias combinações possíveis entre os critérios de divisão de Random Forests e múltiplos valores para o número de modelos e *TreeDepth*. Desta forma é obtido melhor modelo possível perante as opções de modificações



Após as otimizações é feito o treino utilizando uma Random Forest Learner e o respectivo teste do conjunto de teste através do Random Forest Predictor.

Após o loop de otimização foram obtidos os seguintes valores (figura 4.9):

Name	Value
i maxIterations	10
i currentIteration	3
i Loop-Execute	
i Loop (0)	
i iteration	45
i numModels	500
i treeDepth	10
i stopingCriteria	2
i Loop-Execute	
i Loop (1)	
s splitCriterion	InformationGain
s RowID	Row0
i currentIteration	0
i maxIterations	3
i Loop-Execute	
i Loop (2)	
s knime.workspace	D:\Github\SBS

Figura 4.9: Configurações utilizados para o treino do modelo

Capítulo 5

Conclusões

Uma forma mais intuitiva e mais útil é a utilização da metodologia CRISP-DM, ao longo deste projeto foi um ponto fulcral devido a correlação do aumento da complexidade do trabalho e o aumento do conhecimento sobre a UC. Este foi essencial não só para garantir a qualidade dos modelos como também para estruturar e organizar o desenvolvimento do projeto.

O CRISP-DM foi essencial quer a a exploração do dataset relativo aos incidentes em Braga em 2019, quer na exploração do dataset relativo Banco de Portugal.

Relativamente modelos criados, de uma forma geral o resultado é satisfatório com os resultados obtidos. Foi possível uma precisão alta com o modelo gerado para prever a magnitude de incidentes rodoviários na cidade de Braga.

Relativamente ao modelo desenhado para prever o Banco de Portugal, há uma necessidade de um reforço no modelo para poder obter resultados em concreto, alguma falta de conhecimento na utilização da ferramenta KNIME tornou a solução difícil de conseguir.

Sem dúvida há muito ponto que poderiam ser trabalhados/modificados de forma a conseguir resultados ainda melhores e ter um trabalho ainda mais completo.

Bibliografia

- [1] Data Science Project Management
<https://www.datascience-pm.com/crisp-dm-2/>
- [2] Dataset Banco de Portugal
<https://www.kaggle.com/volodymyrgavrysh/bank-marketing-campaigns-dataset>