



Relatório Fase 2 do Trabalho Prático da Unidade Curricular Aprendizagem e Decisão Inteligentes

Grupo 48

A93253 David Alexandre Ferreira Duarte

A90614 Pedro Aquino Martins de Araújo

A94166 Samuel de Almeida Simões Lira

Table of Contents

Introdução	3
Metodologia Escolhida	3
MADRID CLIMATE.....	3
Estudo dos Dados	3
Preparação dos dados	5
Modelação	8
Problema de Classificação.....	8
Avaliação dos resultados	10
Problema de Regressão	11
Avaliação dos Resultados.....	14
Problema de Classificação.....	14
Problema de Regressão	17
SALARY CLASSIFICATION	18
Estudo dos dados	18
Objetivos	18
Preparação dos dados	19
Modelação	19
Modelos com a utilização de árvores de decisão	20
Modelo com redes neuronais	22
Avaliação	23
Modelo por <i>Decision Tree, gini index</i> sem <i>pruning</i>	23
Modelo por <i>Decision Tree, gain ratio</i> com <i>pruning MDL</i>	24
.....	24
Modelo por <i>Artificial Neural Networks</i>	24
Sugestões e Recomendações	24

Formatado: Inglês (Reino Unido)

Introdução

No desenvolvimento deste projeto o grupo decidiu abordar os seguintes paradigmas de aprendizagem, com supervisão e sem supervisão, para cada um dos *dataset* escolhidos. No entanto são utilizadas diferentes técnicas de aprendizagem, escolhidas a fim de melhor explorar cada *dataset*.

Metodologia Escolhida

A metodologia escolhida pelo grupo é a CRISP-DM. O grupo decidiu seguir esta metodologia pois embora esta metodologia seja voltada para desenvolvimento de projetos em concreto, as etapas de Estudo dos Dados, Preparação dos Dados, Modelação e Avaliação do Modelo, estão dentro do âmbito desta cadeira e serão as etapas abordadas neste relatório.

MADRID CLIMATE

Este *dataset* corresponde a recolha de dados relacionados ao clima, observados na cidade de Madrid dentro o período de 01/07/2008 e 20/04/2019. Esta amostra contém o total de 3946 linhas. Nos próximos capítulos, iremos indicar os passos de desenvolvimento de modelos de *machinelearning* aplicados ao *dataset* em questão, seguindo a metodologia mencionada anteriormente.

Estudo dos Dados

No caso deste *dataset* o repositório que o continha, não tinha nenhum problema associado ao mesmo. Por essa razão o grupo decidiu analisar os dados e estudar o comportamento dos mesmos com o intuito de criar um problema para de seguida fazer o modelo que soluciona-se o mesmo.

Para atingir esse objetivo primeiramente foi observado quais são as variáveis existentes na amostra, demonstradas nas seguintes imagens e explicadas a seguir:

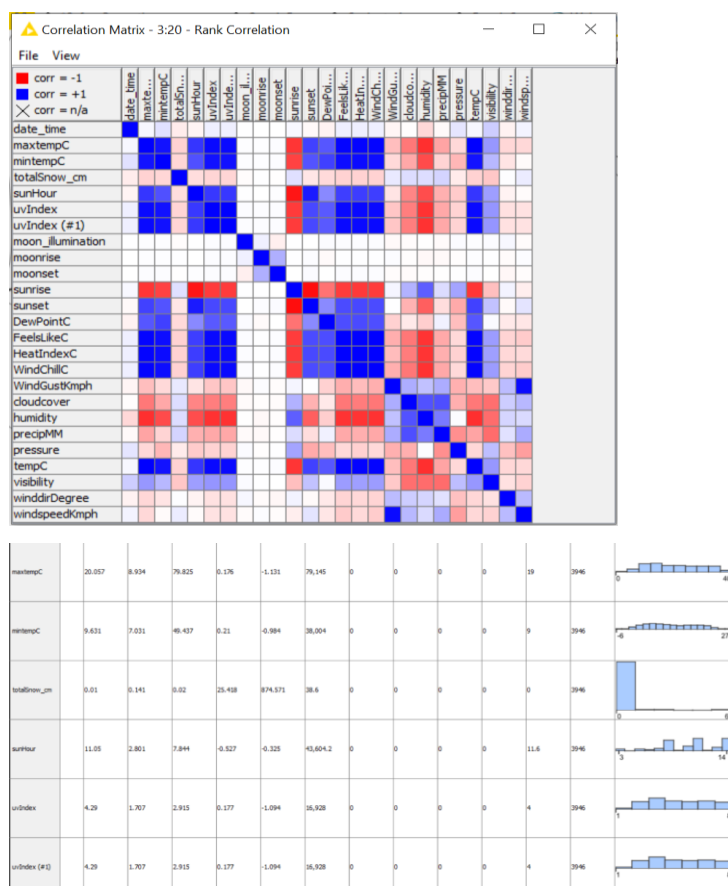
[S] date_time	[I] maxtempC	[I] mintempC	[D] totalSnow_cm	[D] sunHour	[I] uvIndex	[I] uvIndex(≠1)	[I] moon_illumination	[S] moonrise	[S] moonset	[S] sunrise	[S] sunset	[I] DewPointC
2008-07-01	33	17	0	14.5	6	6	4	04:32 AM	08:41 PM	06:48 AM	09:49 PM	6

[I] FeelslikeC	[I] HeatIndexC	[I] WindChillC	[I] WindGustKmph	[I] cloudcover	[I] humidity	[D] precipMM	[I] pressure	[I] tempC	[I] visibility	[I] winddirDegree	[I] windspeedKmph
28	28	28	11	5	28	0	1014	33	10	203	8

- date_time : representa a data exata que foi coletado os dados
- TotalSnow_cm: total de neve no dia
- MintempC: mínimo de temperatura coletada o dia
- MaxtempC: máximo de temperatura coletada no dia
- FeelslikeC: temperatura ambiente
- Humidity: humidade registada no dia
- PrecipMM: precipitação de chuva no dia
- Moon_illumination: Índice de iluminação da lua

- DewPointC: Corresponde aos Pontos Orvalho, que é a temperatura no qual o vapor de água no ar passa ao estado líquido
- HeatIndexC: Índice de calor, uma medida importante para identificar o calor sentido no dia
- WindChillC: Temperatura sentida de acordo com o vento
- CloudCover: Índice de cobertura no céu que as nuvens tiveram
- Pressure: Pressão do ar
- TempC: temperatura registada no dia
- Visibility: Índice de visibilidade do ambiente
- WindspeedKmph : velocidade do vento em km/h
- UvIndex: índice dos raios ultravioleta

Após verificar as variáveis que o dataset possui, foi verificado também qual é o seu comportamento. Para este fim, foi utilizado o nodo “*Rank Correlation*” para verificação do relacionamento entre as variáveis, e também o nodo “*Statistics*” com o objetivo de verificar as estatísticas das mesmas.



Após a análise do comportamento de cada variável, é possível identificar quais possuem aspetos mais interessantes, variação de valores consideráveis e correlação com outras variáveis. Seguindo estes critérios o grupo atribuiu uma maior relevância para o decorrer do projeto às seguintes variáveis:

- date_time
- MinTempC
- MaxTempC
- FeelslikeC
- Humidity
- PrecipMM
- DewPointC
- WindChillC
- UvIndex
- TempC
- HeatIndex
- WindSpeedkmph

Preparação dos dados

Após ter sido realizada a análise dos dados, foi realizado o tratamento destes a fim de ser possível criar modelos que obtenham os melhores resultados possíveis.

Primeiramente é necessário verificar se há valores em falta na tabela, os chamados *missing values*, se houver retirá-los pois são danosos para os resultados. Na amostra escolhida não se observou nenhum *missing value* em nenhum registo, observado no nodo “*Statistics*” anterior.

Observou-se pela análise dos dados que os parâmetros **minTempC** e **maxTempC** possuem uma diferença grande registada em cada dia, assumindo a hipótese de em Madrid existir grandes amplitudes térmicas. A fim de verificar com maior detalhe esta característica, foi criado uma nova coluna ‘Amplitude térmica’ com auxílio do nodo “*Math Formula*”.

File Edit History Navigation View																						
Table "Default" - Rows: 12 Spec - Columns: 11 Properties Flow Variables																						
	I	Month	D	Mean(festIndexC)	D	Mean(humidity)	D	Mean(windChillC)	D	Mean(uvIndex)	D	Mean(Amplitude térmica)	D	Mean(temperC)	D	Mean(headIndexC)	D	Mean(DewPointC)	D	Mean(humidtempC)	D	Mean(humidtempC)
1		4.39		73.759		4.39		2.381		7.692		9.836		8.57		1.463		9.836		2.144		2.144
2		4.819		66.803		4.819		2.403		8.658		11.181		7.055		1.081		11.181		2.523		2.523
3		8.475		64.783		8.475		3.211		10.519		14.9		10.106		2.862		14.9		4.381		4.381
4		11.888		63.294		11.888		3.725		11.381		18.094		13.025		5.244		18.094		6.712		6.712
5		16.997		54.152		17.01		4.71		12.016		22.6		17.519		6.987		22.6		10.584		10.584
6		23.213		42.813		23.407		5.87		12.18		28.327		23.337		8.643		28.327		16.147		16.147
7		27.05		33.199		27.462		6.625		12.935		32.288		27.073		9.226		32.288		19.273		19.273
8		27.035		33.622		27.443		6.625		13.018		32.276		27.012		8.437		32.276		19.258		19.258
9		22.227		43.579		22.303		5.603		11.842		27.024		22.336		8.048		27.024		15.182		15.182
10		16.513		57.724		16.504		4.49		9.768		20.768		16.938		7.284		20.768		11		11
11		8.933		71.27		8.933		3.17		7.661		13.445		10.258		4.764		13.445		5.785		5.785
12		5.61		73.985		5.61		2.513		7.651		10.334		7.141		2.235		10.334		2.683		2.683

Após realizar a análise das tabelas acima, concluiu-se que as variáveis em cada ano no geral não variaram muito e possuem valores muito próximos. Em relação ao agrupamento realizado por mês, foi verificado um comportamento mais interessante, existe uma variação entre as variáveis dependendo mês registrado. Outra análise verificada acerca da tabela dos agrupamentos por mês, foi de que os valores têm tendência a variarem de acordo com a época do ano que estão. A partir disso, para ser possível analisar o registo de acordo com a estação do ano, o que é mais interessante no caso estudado do que uma data específica, foi criado uma nova coluna para indicar a correspondente estação do ano de cada registo. Para atingir tal objetivo, foi utilizado o Nodo "Rule Engine".

totalSnow_cm

sunHour

uvIndex

uvIndex (#1)

moonIllumination

moonrise

moonset

sunrise

sunset

Flow Variable List

krime.workspace

?

>

>=

<

<=

AND

OR

NOT

LIKE

MATCHES

MISSING

Expression

4

// TRUE => "default outcome"

5

\$Month (number)\$ = 9 AND \$Day of month\$ >= 23 => "outono"

6

\$Month (number)\$ > 9 AND \$Month (number)\$ < 12 => "outono"

7

\$Month (number)\$ = 12 AND \$Day of month\$ < 22 => "outono"

8

\$Month (number)\$ = 12 AND \$Day of month\$ >= 22 => "inverno"

9

\$Month (number)\$ >= 1 AND \$Month (number)\$ < 3 => "inverno"

10

\$Month (number)\$ = 3 AND \$Day of month\$ < 20 => "inverno"

11

\$Month (number)\$ = 3 AND \$Day of month\$ >= 20 => "primavera"

Append Column:

estacao

Replace Column:

Month (number)

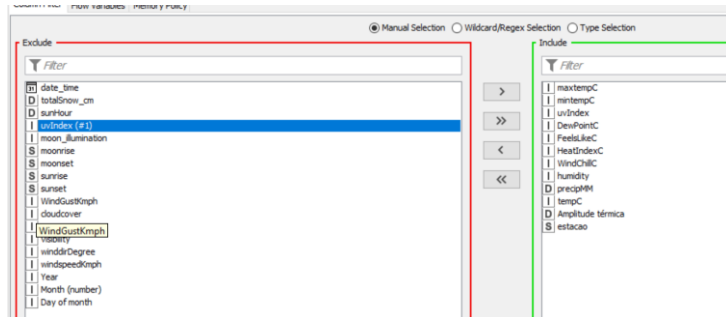
OK

Apply

Cancel

?

Após foi realizado o filtro das colunas que se mostraram mais relevantes no caso estudado, após análises estatísticas, correlações e das que foram consideradas importantes para criação dos modelos. Vale informar que para chegar a estas colunas houve diversas mudanças a partir de testes realizados aos modelos.



Modelação

Após se ter realizado a análise e tratamento dos dados, foi possível criar modelos de aprendizagem. Entretanto, como o *dataset* não veio associado nenhum problema, como mencionado anteriormente, houve a necessidade e a possibilidade de criar problemas. Com o objetivo de diversificar e colocar em prática uma maior quantidade de conhecimentos aprendidos nas aulas, o grupo criou um problema de classificação, e outro de regressão.

Problema de Classificação

O grupo definiu a previsão do campo de estação de ano como problema de classificação. Chegou-se a este problema através da análise realizada anteriormente das variáveis que tinham valores distintos em determinados meses e épocas do ano. Realizou-se uma outra análise para verificar com maior detalhe a variação dos dados por estação do ano, e por fim foi utilizado agrupamentos de dados a partir do nodo “*GroupBy*”.

estacao	Mean(FeelsLikeC)	Mean(humidity)	Mean(WindChillC)	Mean(HeatIndexC)	Mean(Amplitude térmica)	Mean(maxtempC)	Mean(mintempC)	Mean(uvIndex)	Mean(DewPointC)	Mean(Amplitude térmica)
...inverno	5.476	70.053	5.476	7.388	8.599	11.333	2.734	2.629	1.653	11.333
...outono	11.752	65.317	11.748	12.723	8.644	16.226	7.582	3.653	5.396	16.226
...primavera	15.213	56.963	15.24	16.005	11.596	21.022	9.426	4.377	6.235	21.022
...verao	25.983	35.476	26.313	25.998	12.78	31.118	18.338	8.397	8.262	31.118

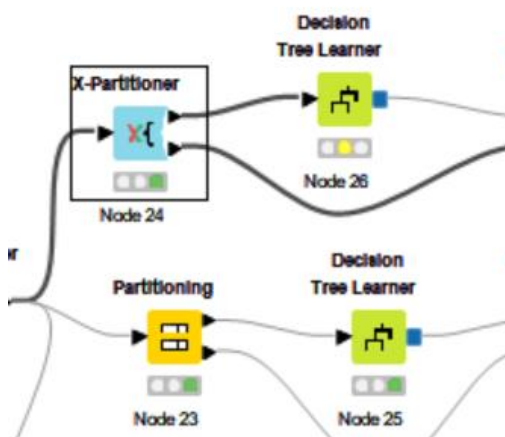
A partir da tabela acima, é possível verificar que no verão e no inverno as variáveis possuem comportamentos muito distintos, e a primavera e o outono têm valores intermédios e próximos um do outro. Como todo problema de *machine learning* surge de uma hipótese criada a partir de observações e estudos, o grupo chegou à seguinte hipótese: “A cidade de Madrid possui 4 estações bem definidas?”.

Com o objetivo de testar esta hipóteses definida foram criados modelos a partir de duas técnicas de classificação: *Decision Tree Learner* e *Clustering*.

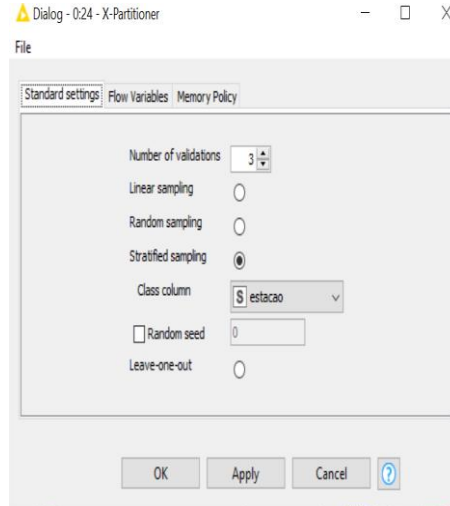
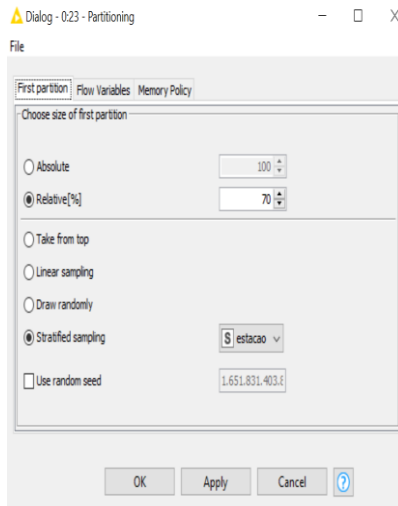
- **Decision Tree Learner**

Nesta técnica foi utilizado o modelo com supervisão, e com o objetivo de validar e analisar se o modelo obteve resultados satisfatórios, é necessário ter dados de treino, que servirão para a criação do modelo, e dados de teste, que terão a utilidade de testar o modelo criado.

A fim de cumprir este requisito, foi utilizado duas técnicas de separação de dados, *X-Partitioner* e *Partitioning*.



Foi também utilizado como critérios de separação no nodo *X-partitioner* 3 validações, pelo motivo de ser pouca quantidade de dados e ter sido o que apresentou melhores resultados, e para o modo com que é escolhido os dados foi utilizado o *Stratified Sampling* para a coluna 'estação', de forma que os dados sejam escolhidos de igual forma para cada estação, para assim evitar enviesamento, o chamado *bias* problema muito encontrado no estudo de *machine learning*. E para o *Partitioning* foi utilizado a taxa de 70% para dados de treino e o restante para teste, com a justificação de apresentar melhores resultados e também por serem poucos dados.



Avaliação dos resultados

Nos parâmetros utilizados no nodo “*Decision Tree Learner*” foi utilizado como *target* o campo estação, por ser a coluna objetivo, ou seja, a coluna que se pretende prever, e como *quality measure*, o *Gain ratio* e *pruning MDL*, uma vez que conseguem produzir melhores resultados.

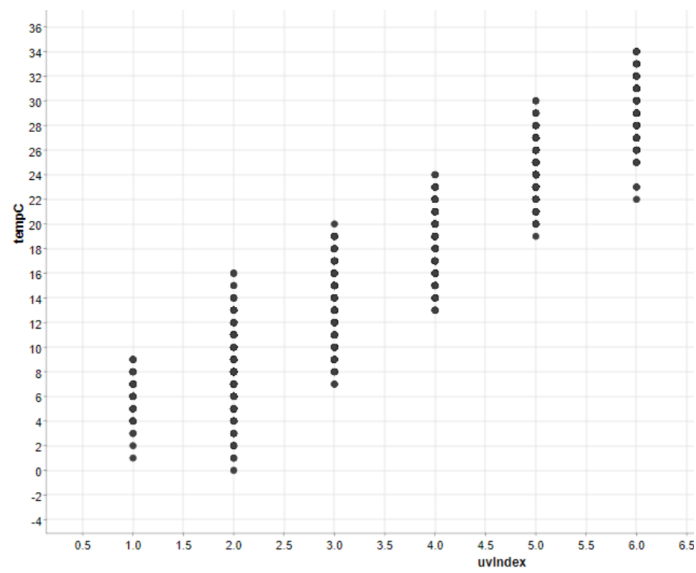
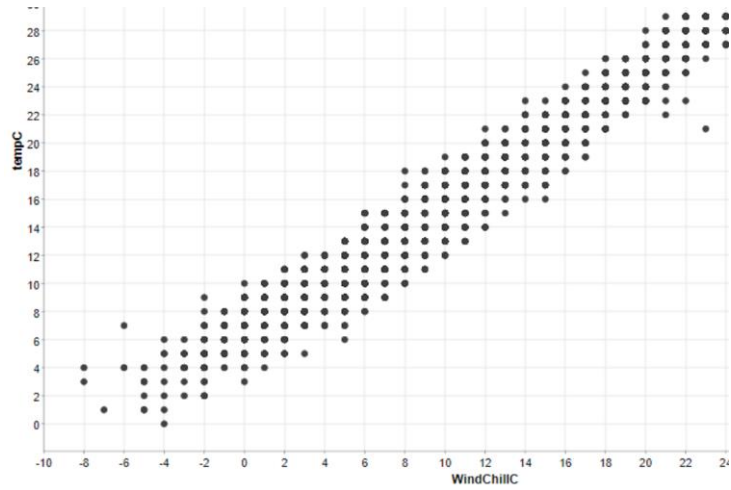
- **Clustering**

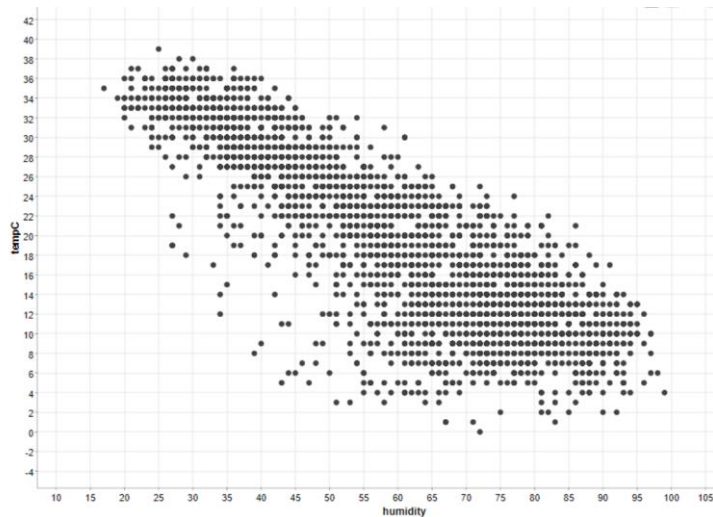
Foi utilizado a técnica *clustering* com aprendizagem sem supervisão, com objetivo de verificar se é possível categorizar os dados e separar eles de acordo com seu comportamento.

Dentre as técnicas de *clustering* foram utilizadas: *k-means* e *Hierarchical Clustering*. Especificou-se 4 clusters em cada modelo, sabendo que havia 4 estações, e assim é possível verificar se os dados estão bem distribuídos, e assim comprovar se é possível criar 4 categorias com comportamento diferente.

Foi utilizado o “*Partitioning*” com as mesmas configurações utilizadas anteriormente, a fim de testar o modelo *k-means* a partir do nodo “*cluster assigner*”. Para verificar se o cluster atribuído teve sucesso ou não de acordo com a estação do ano, utilizou-se o “*rule engine*” para atribuir ao nome do *cluster* a estação do ano correspondente, para assim conseguir comparar de forma efetiva os parâmetros de treino e teste.

Na imagem a seguir mostra como ficou o *pipeline*.

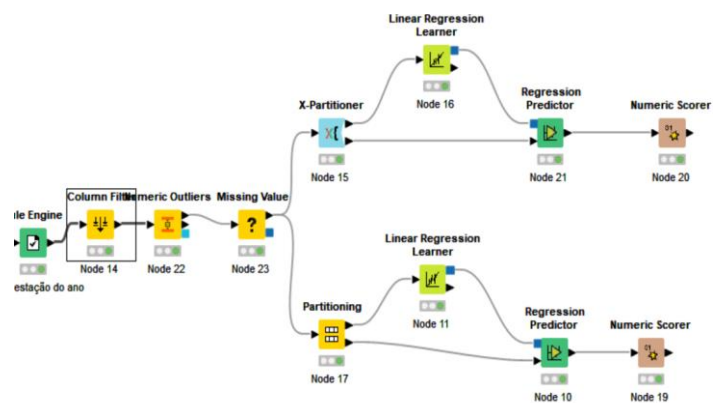




A partir destas observações e análises, observou-se que o campo 'tempC' possui um comportamento linear em detrimento da variação de outras variáveis. Dito isto, criou-se um modelo *Linear Regression*, a fim de verificar se um modelo baseado na equação linear consegue atingir índices de previsão satisfatórios para a variável fixa, no caso o campo 'tempC'.

Para atingir este objetivo o grupo realizou mais alguns tratamentos nos dados na tentativa de obter os melhores resultados no modelo, o que inclui a remoção do *outlier* e em seguida remoção dos valores em falta. Foi também utilizado as técnicas de separação de dados ie *partitioning* para o teste do modelo, utilizando basicamente as mesmas configurações dos modelos já vistos anteriormente.

Obtendo-se assim o seguinte pipeline:



Avaliação dos Resultados

A análise dos resultados produzidos pelos modelos serão divididos separadamente por problema, o problema de Classificação e de Regressão.

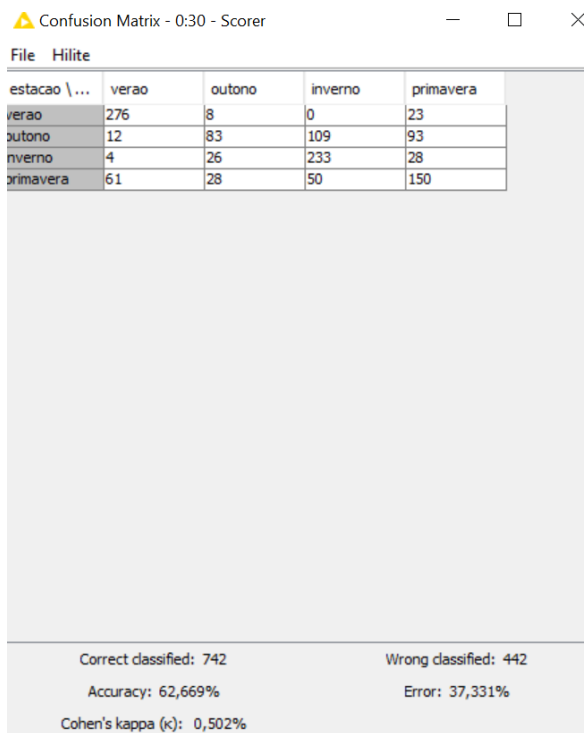
Problema de Classificação

- Resultados**

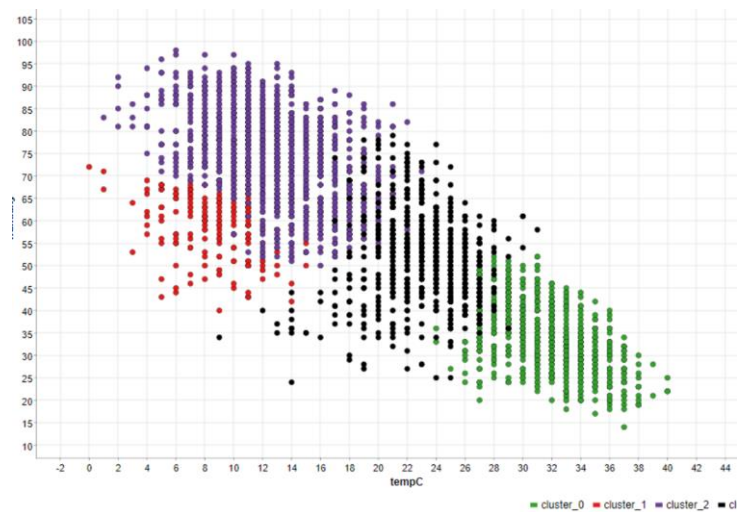
Os resultados produzidos pelos modelos de Decision Tree, a partir dos nodos X-partitioner e Partitioning consequentemente:

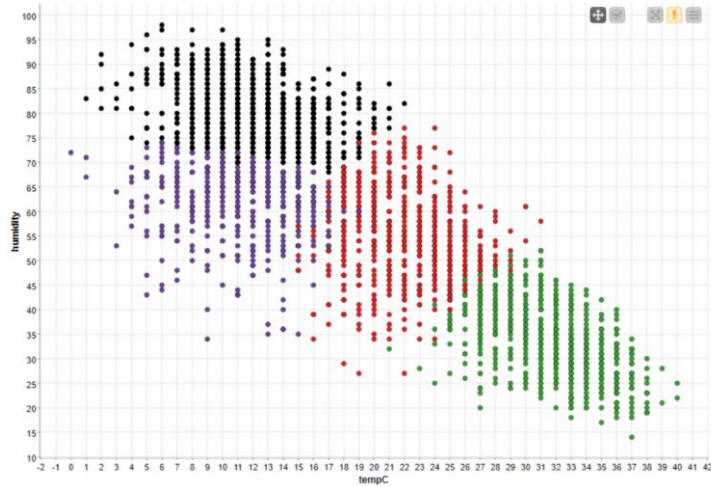
Confusion Matrix - 0:31 - Scorer				
File	Hilite			
estacao \...	verao	outono	inverno	primavera
verao	321	0	1	20
outono	29	54	164	83
inverno	0	10	294	19
primavera	69	20	63	169

Correct classified: 838	Wrong classified: 478
Accuracy: 63,678%	Error: 36,322%
Cohen's kappa (κ): 0,516%	



Gráficos produzidos pelos modelos de clustering utilizando as variáveis 'tempC' e 'humidity', por hierarchical clustering e k-means, consequentemente:





E o resultado produzido pelo scorer a partir do modelo k-means:

Scorer View

Confusion Matrix

	inverno (Predicted)	outono (Predicted)	primavera (Predi...	verao (Predicted)	
inverno (Actual)	148	122	21	0	50.86%
outono (Actual)	132	48	94	23	16.16%
primavera (Actual)	51	30	156	52	53.98%
verao (Actual)	1	0	29	277	90.23%
	44.58%	24.00%	52.00%	78.69%	

Overall Statistics

Overall Accuracy	Overall Error	Cohen's kappa (κ)	Correctly Classified	Incorrectly Classified
53.13%	46.88%	0.375	629	555

- **Análise**

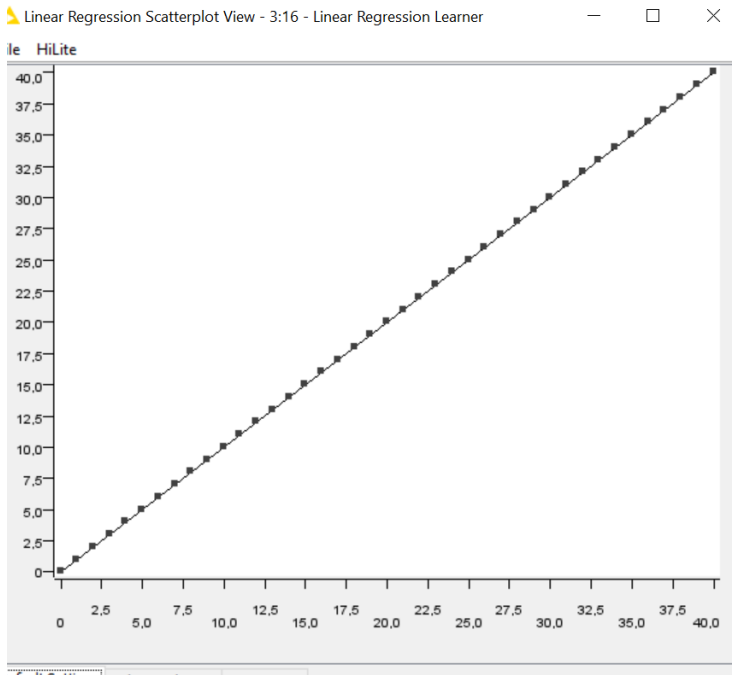
A partir dos resultados apresentados é possível inferir alguns aspetos acerca do clima de Madrid e refutar a hipótese do problema realizado anteriormente. Observa-se que nas estações de inverno e verão há uma grande distinção de variáveis climáticas, deduz-se a partir do alto índice de acerto nas previsões e partir da análise de dados feitas. Entretanto, já as estações de primavera e outono possuem valores semelhantes as outras estações sendo de difícil acerto, principalmente o outono onde obteve-se índices de acerto muito baixos.

Conclui-se, que a partir dos dados utilizados deste *dataset* do clima de Madrid, que a hipótese que Madrid possui 4 estações bem definidas não se confirma, entretanto, apresenta 2 estações muito bem definidas, sendo que as outras estações possuem valores mais semelhantes às demais.

Problema de Regressão

- Resultados

Os resultados produzidos pelo modelo linear regression:



S...

File	
R ² :	1
Mean absolute error:	0
Mean squared error:	0
Root mean squared error:	0
Mean signed difference:	0
Mean absolute percentage error:	0
Adjusted R ² :	1

Vale destacar que os resultados tanto do X-Partitioning quanto do Partitioning foram iguais.

- **Análise**

Verifica-se pelos resultados que o modelo teve 100% de eficácia com o R-Square Value = 1, correspondendo ao entendimento da variância de todos os campos, e com o Root Mean Square Error (RMSE) = 0, ou seja, não houve nenhum erro apresentado.

Conclui-se que a análise realizada da linearidade e correlação das variáveis selecionadas em detrimento ao campo tempC, estavam corretas, e que desta forma implica-se que a temperatura de Madrid possui valores que são relacionados a outros fatores climáticos também.

SALARY CLASSIFICATION

Para a análise dos dados deste dataset, o grupo optou por seguir a metodologia CRISP-DM. Por se tratar de um projeto com fins acadêmicos, decidiu-se dar prioridade às seguintes etapas da metodologia: estudo dos dados, preparação dos dados, modelação e avaliação.

Estudo dos dados

O *dataset salary classification* possui dados referentes a 48842 registos distintos, e contém 15 campos de informação (colunas) são estas:

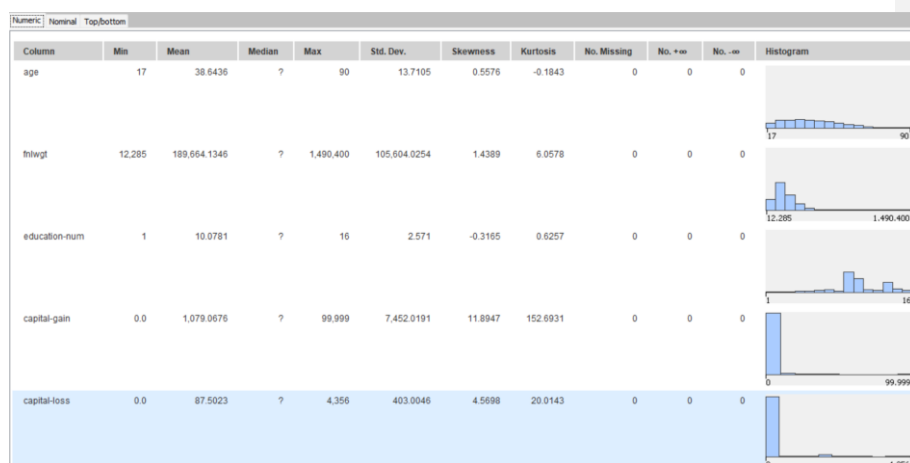
- *age*: idade de um indivíduo. Inteiro maior que 0.
- *workclass*: termo geral que representa a categoria de trabalho de um indivíduo.
- *fnlwgt*: peso final. Em outras palavras, número de pessoas as quais esse registo representa. Inteiro maior que 0.
- *education*: o nível mais alto de educação alcançado por um indivíduo.
- *education-num*: o nível mais alto de educação alcançado por um indivíduo em forma numérica. Inteiro maior que 0.
- *maritalstatus*: marital status of an individual.
- *occupation*: tipo geral de ocupação profissional de um indivíduo.
- *relationship*: representa o que esse indivíduo é em relação a outros. Cada entrada tem apenas um atributo de relacionamento e se confunde em muito com a coluna *maritalstatus*.
- *race*: descrição da etnia de um indivíduo.
- *sex*: sexo biológico de um indivíduo. 2 valores únicos possíveis.
- *capitalgain*: ganho de capital para um indivíduo.
- *capitalloss*: Perda de capital para um indivíduo.
- *hoursperweek*: horas de trabalho um indivíduo relatou ter trabalhado.
- *nativecountry*: país de origem de um indivíduo.
- *the label*: Rótulo que indica se um indivíduo ganha anualmente mais de \$50,000. Os únicos valores possíveis são $\leq 50k$ e $> 50k$.

Objetivos

Com a informação disponibilizada o grupo decidiu aplicar os conhecimentos aprendidos sobre o paradigma de aprendizagem com supervisão, com o objetivo de classificar se um indivíduo ganha anualmente mais de \$50,000 ou não.

Preparação dos dados

O *dataset* não possui valores em falta. No entanto, apresenta algumas colunas redundantes, como as colunas *education* e *education-num* que indicam a mesma informação, sendo que a primeira de forma textual e a segunda de forma numérica. Com a limpeza dos dados em mente, filtrou-se a coluna *education*. Algo semelhante ocorre com as colunas *relationship* e *marital-status*, onde se utilizou o nodo “Rule Engine” para se alterar os valores da coluna *marital status* de forma que a coluna *relationship* se torne completamente redundante, portanto descartável considerando o objetivo descrito acima. Ao analisar os histogramas das colunas numéricas, através do nodo “Statistics”, notou-se uma particularidade com as colunas ‘*capital-gain*’ e ‘*capital-loss*’ ambas tinham como maioria absoluta o zero como valor, e níveis de *skewness* muito elevados, 11.8947 e 4.5698, respetivamente, como se pode observar na figura abaixo. As colunas filtradas até ao momento são: ‘*education*’, ‘*education-num*’, ‘*capital-gain*’ e ‘*capital-loss*’.



Por se tratar de apenas um ficheiro *csv*, portanto apenas uma fonte de dados, não houve a necessidade se aplicar técnicas de integração.

Os dados também passam pelos nodos *Category to Number* e *Normalizer*, que respectivamente transformam os valores textuais que representam categorias em números e normalizam os valores da categoria *fnlwgt* para melhor processamento dos mesmos pelos nodos de aprendizagem.

Modelação

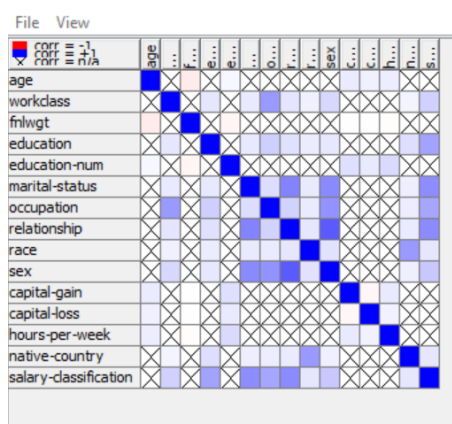
Para este *dataset* foi tomada a decisão de se utilizar duas técnicas de aprendizagem. Um *workflow* *Knime* foi criado utilizando árvores de decisão de tipo discreto e o outro utilizando redes neurais.

Modelos com a utilização de árvores de decisão

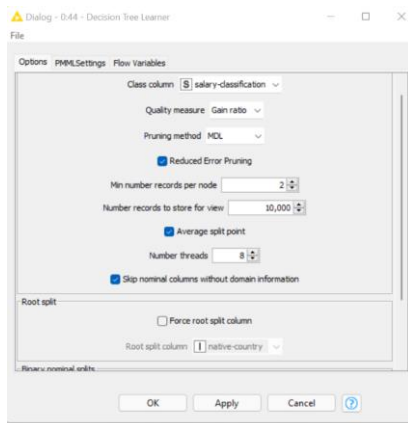
O grupo aplicou a técnica de árvores de decisão de tipo discreto por ter em vista que o atributo de decisão (*salary-classification*) se divide em duas possíveis categorias se um indivíduo ganha ou não anualmente mais de \$50,000.

Esta técnica de aprendizagem segue o paradigma *bottom up*, ou seja, modelo é construído pela identificação das relações entre os atributos do *dataset*. Esta técnica foi escolhida pela sua fácil compreensão e configuração, e atrelado ao fato de haver poucos atributos com níveis significativos de correlação entre si (ver figura abaixo), não apresenta vulnerabilidades muito expressivas.

Correlation Matrix - 0:42 - Linear Correlation

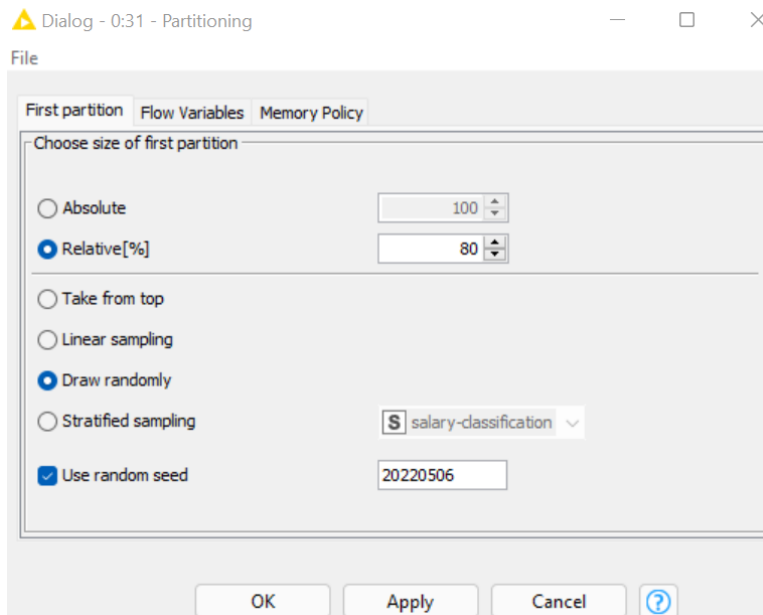


O grupo decidiu fazer dois modelos utilizando árvores de decisão, um utilizando como métrica de qualidade *gain ratio* e como método *MDL* de *pruning* e o outro utilizando o *gini index* sem métodos de *pruning*. Através do pruning o tamanho das árvores de decisão é reduzido ao remover partes da árvore que são redundantes e não importantes para aprendizagem, tornando-se assim possível prevenir *overfitting* do modelo aos dados aumentando a precisão da capacidade de previsão.

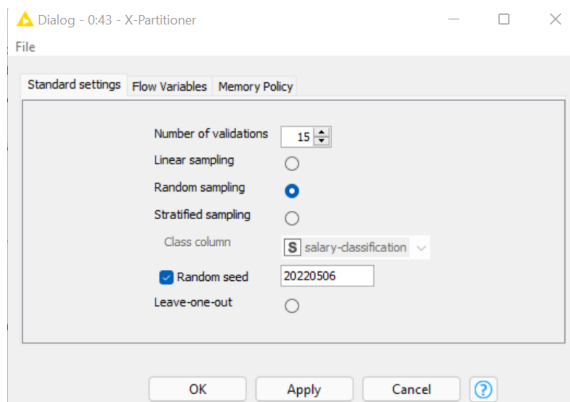


Após a preparação dos dados, o grupo recorreu a duas técnicas de *partitioning*, o nodo “*Partitioner*” e o nodo “*X-Partitioner*”.

O nodo “*Partitioner*” divide os dados em duas partições, uma com 80% dos dados de forma aleatória com recurso a uma *random seed* que garante a possibilidade de refazer as partições com os mesmos dados, visto que o tamanho se mantenha. A partição com 80% dos dados é então passada para o nodo “*Decision Tree Learner*” como mencionado anteriormente, enquanto o restante dos dados é passado junto com modelo gerado pelo nodo “*Decision Tree Learner*” para o nodo “*Decision Tree Predictor*” que testa o modelo utilizando os dados recebidos.



Já o nodo “X-Partitioner”, efetua a validação por cruzamento de dados, dividindo o conjunto de dados em k folds, onde durante k interações e em cada interação são utilizadas $k-1$ folds para treino e 1 para teste, neste caso em concreto o valor de k é 15, ver figura abaixo para a configuração do nodo.



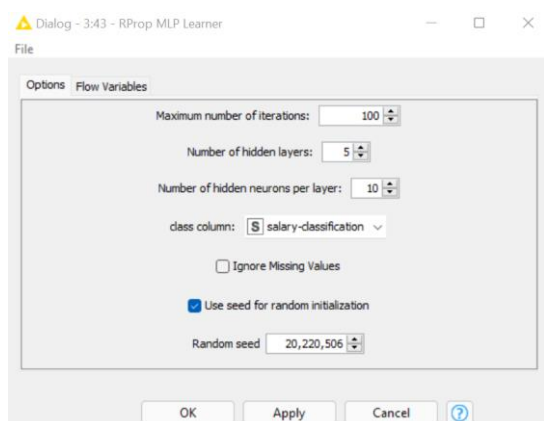
Modelo com redes neuronais

Uma Rede Neuronal Artificial é um sistema computacional de base conexionista inspirado no sistema nervoso central dos seres humanos, definido por uma estrutura interligada de unidades computacionais com capacidade de aprendizagem, tais unidades são designadas neurónios. Redes neuronais podem ser construídas de acordo com 3 tipos de arquitetura distintos, para este modelo foi escolhido a arquitetura *Feed-Forward* multicamada, arquitetura organizada por camadas, podendo

possuir 1 ou mais camadas intermédias, com conexões unidirecionais, a cada camada adicional maior é a capacidade da rede em modelar funções de maior complexidade.

Os dados são particionados com o nodo “*Partitioner*” já previamente explicado, devido a simplicidade deste tipo de validação tendo em vista o aumento na complexidade da técnica de aprendizagem, o nodo é configurado de forma semelhante ao modelo anterior.

O nodo “*RProp MLP Learner*” foi configurado de forma a permitir um extensivo processo de aprendizagem com um razoável número de camadas e número de neurónios por camadas, ao passo que o tempo de execução do nodo não se estenda demasiadamente.



Avaliação

Para avaliação dos resultados o nodo “*Scorer(Javascript)*” foi utilizado para uma mais agradável visualização da matriz de confusão e das estatísticas gerais do modelo.

A matriz de confusão permite um melhor entendimento da precisão do modelo.

Modelo por *Decision Tree*, *gini index* sem *pruning*

A matriz de confusão produzida por este modelo apresenta uma alta precisão ao classificar corretamente dados que ganham até \$50,000 por ano, no entanto essa precisão cai ao classificar dados que ganhem mais de \$50,000 por ano. No entanto produzindo uma precisão geral razoável de 79.19% da capacidade do modelo. *Cohen’s kappa* é uma outra métrica que pode ser utilizada para se medir a performance de um modelo, com uma gama de valores entre -1 e 1, para dados mais balanceados o *Cohen’s kappa* tende a aumentar.

Scorer View

Confusion Matrix

	<=50K (Predicted)	>50K (Predicted)	
<=50K (Actual)	6482	991	86.74%
>50K (Actual)	1042	1254	54.62%
	86.15%	55.86%	

Overall Statistics

Overall Accuracy	Overall Error	Cohen's kappa (κ)	Correctly Classified	Incorrectly Classified
79.19%	20.81%	0.417	7736	2033

Modelo por *Decision Tree*, *gain ratio* com *pruning MDL*

A matriz de confusão produzida por este modelo apresenta uma alta precisão ao classificar corretamente dados que ganham até \$50,000 por ano, no entanto essa precisão cai ao classificar dados que ganhem mais de \$50,000 por ano, contudo menos do que o apresentado pelo modelo anterior que não aplica um método de *pruning*. Não obstante produzindo uma precisão geral razoável de 83.39% para este modelo, o melhor resultado entre os modelos utilizados. Também possui o maior *Cohen's Kappa* dentre os 3 modelos utilizados.

Scorer View

Confusion Matrix

	<=50K (Predicted)	>50K (Predicted)	
<=50K (Actual)	2223	220	90.99%
>50K (Actual)	321	493	60.57%
	87.38%	69.14%	

Overall Statistics

Overall Accuracy	Overall Error	Cohen's kappa (κ)	Correctly Classified	Incorrectly Classified
83.39%	16.61%	0.538	2716	541

Modelo por *Artificial Neural Networks*

A matriz de confusão produzida por este modelo apresenta uma maior diferença ao classificar corretamente os valores, tendo precisão muito baixa, 42.86% ao classificar corretamente se um registo ganha mais de \$50,000 por ano. No entanto produziu uma precisão geral razoável de 81.77% de capacidade de previsão do modelo.

Scorer View

Confusion Matrix

	<=50K (Predicted)	>50K (Predicted)	
<=50K (Actual)	7004	469	93.72%
>50K (Actual)	1312	984	42.86%
	84.22%	67.72%	

Overall Statistics

Overall Accuracy	Overall Error	Cohen's kappa (κ)	Correctly Classified	Incorrectly Classified
81.77%	18.23%	0.419	7988	1781

Sugestões e Recomendações

Após análise dos resultados obtidos com os modelos envolvendo o *dataset salary_classification*, a técnica de aprendizagem mais robusta, redes neuromais, não necessariamente produziu o melhor resultado. A fim de melhorar os resultados obtidos nomeadamente com as redes neuronais podia-se ter recorrido ao uso de técnicas de *feature selection*, que iriam indicar as features mais bem qualificadas para aprendizagem autónoma presentes neste dataset.