



Universidade do Minho

Escola de Ciências da Universidade do Minho

Departamento de Informática

Mestrado em Matemática e Computação

Mestrado Integrado em Engenharia Informática

Redes Neurais Recorrentes para previsão do fluxo de tráfego rodoviário

Alunos:

Andreia Costa (PG37013)

Henrique Faria (A82200)

Paulo Barbosa (PG40160)

Rui Teixeira (PG37021)

Docentes:

Bruno Fernandes

Victor Alves

Unidade Curricular: Classificadores e Sistemas Conexionistas

Maio
2020

Conteúdo

1	Introdução	1
2	<i>Dataset</i>	2
2.1	<i>Traffic Flow Braga</i>	2
2.2	<i>Traffic Incidents Braga</i>	3
2.3	<i>Weather Braga Descriptions</i>	3
2.4	<i>Weather Braga</i>	4
2.5	Preparação dos dados	4

1 Introdução

2 *Dataset*

Aquando da apresentação do presente trabalho foram disponibilizados dados referentes a duas cidades: Braga e Porto, sendo que o grupo escolheu os dados relativos à cidade de Braga para trabalhar.

Os dados encontram-se distribuídos em 4 *datasets*:

- *Traffic Flow Braga Until 20191231*;
- *Traffic Incidents Braga Until 20191231*;
- *Weather Braga Descriptions Until 20191231*;
- *Weather Braga Until 20191231*.

Todos os *datasets* contêm dados relativos ao período entre 15 Janeiro 2019 e 31 Dezembro 2019.

2.1 *Traffic Flow Braga*

O *dataset* "Traffic Flow Braga" é constituído pelos seguintes atributos:

- *city_name*;
- *road_num*;
- *road_name*;
- *functional_road_class_desc*;
- *current_speed*;
- *free_flow_speed*;
- *speed_diff*;
- *current_travel_time*;
- *free_flow_travel_time*;
- *time_diff*;
- *creation_date*.

2.2 *Traffic Incidents Braga*

- *city_name*;
- *description*;
- *cause_of_incident*;
- *from_road*;
- *to_road*;
- *affected_roads*;
- *incident_category_desc*;
- *magnitude_of_delay_desc*;
- *length_in_meters*;
- *delay_in_seconds*;
- *incident_date*;
- *latitude*;
- *longitude*.

2.3 *Weather Braga Descriptions*

- *city_name*;
- *cloudiness*;
- *atmosphere*;
- *snow*;
- *thunderstorm*;
- *rain*;
- *sunrise*;
- *sunset*;
- *creation_date*.

2.4 *Weather Braga*

- *city_name*;
- *temperature*;
- *atmospheric_pressure*;
- *humidity*;
- *wind_speed*;
- *clouds*;
- *precipitation*;
- *current_luminosity*;
- *sunrise*;
- *sunset*;
- *creation_date*.

2.5 Preparação dos dados

Após análise dos quatro *datasets* concluiu-se que, antes de se desenvolver o modelo para a previsão da *feature speed_diff*, era necessário fazer uma prévia preparação dos dados.

Começou-se por fazer um prévio tratamento do *dataset Traffic_Incidents*. Para isso, quadriplicou-se esse *dataset*, com o intuito de atribuir todas as ruas em estudo a todos os incidentes, para que posteriormente fosse possível avaliar a distância entre os incidentes e as ruas em estudo.

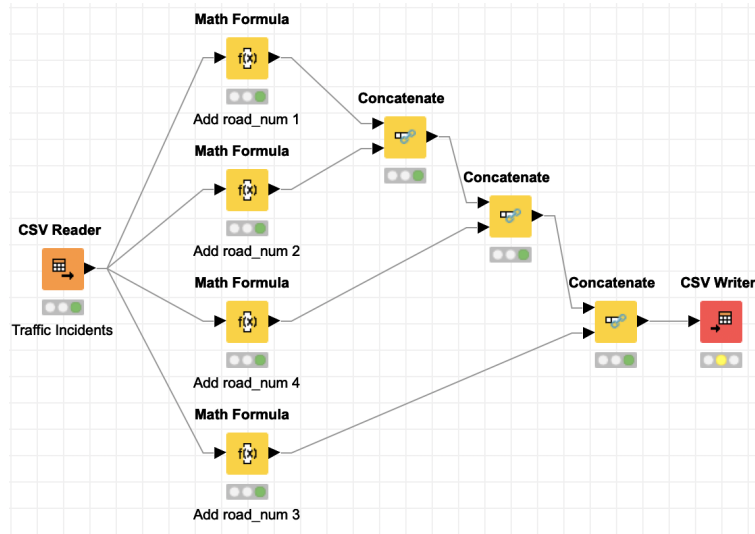


Figura 1: Preparação do *dataset Traffic Incidents*.

De seguida, recorrendo à latitude e longitude dos diferentes acontecimentos, calculou-se a distância dos incidentes a cada uma das ruas, para perceber quais os incidentes que podiam influenciar o *speed_diff* de uma determinada rua.

```

1 import pandas as pd
2 from math import radians, sin, cos, atan2, sqrt
3
4 df = pd.read_csv('Traffic_Incidents.csv', delimiter = ',',
5                 error_bad_lines = False, encoding = 'ISO-8859-1')
6
7 def distance(p1, n):
8     R = 6371.0
9     if n == 1:
10         lat2 = radians(41.548331)
11         lon2 = radians(-8.421298)
12     elif n == 2:
13         lat2 = radians(41.551356)
14         lon2 = radians(-8.420001)
15     elif n == 3:
16         lat2 = radians(41.546639)
17         lon2 = radians(-8.433517)
18     else:
19         lat2 = radians(41.508849)
20         lon2 = radians(-8.462299)
21     lat1, lon1 = radians(p1[0]), radians(p1[1])
22     dlon = lon2 - lon1
23     dlat = lat2 - lat1
24     a = sin(dlat / 2)**2 + cos(lat1) * cos(lat2) * sin(dlon /

```

```

2) **2
24 c = 2 * atan2(sqrt(a), sqrt(1 - a))
25 distance = R * c
26 return distance
27
28 df['Distance'] = df.apply(lambda row: distance((row['latitude
    '], row['longitude']), row['road_num']), axis=1)

```

Após calculadas todas as distâncias fez-se um tratamento estatístico, tendo-se obtido os seguintes resultados:

- $max = 6313,251$;
- $min = 0,0228$;
- $mean = 4,507$;
- $standard\ deviation = 81,789$.

Através dos resultados obtidos é possível verificar que existem *outliers*, uma vez que, sendo os dados recolhidos referentes apenas à cidade de Braga era impossível que a distância máxima dos incidentes às ruas fosse de cerca de 6313 km. Devido a este facto, optou-se por remover alguns dados do *dataset*. Uma vez que a distância é medida em linha reta utilizou-se como *threshold*, para remover dados, vários valores, nomeadamente, 0,5, 1 e 1,5.

Após feito este tratamento procedeu-se à preparação dos dados referentes aos restantes *datasets*, com o intuito de se obter, no final, um único *dataset*.

Começou-se por fazer o tratamento do *dataset Weather_Descriptions_Braga*, tendo-se removido as colunas: *city_name*, *snow* e *cloudiness*. A coluna *snow* apresentava apenas *missing values*, daí se ter optado pela sua remoção. Relativamente à coluna *cloudiness*, optou-se por fazer a remoção da mesma, uma vez que existe uma coluna que está diretamente relacionada com esta, a coluna *cloud*, e que não apresenta *missing values*.

De seguida, procedeu-se à remoção das colunas *city_name* e *precipitation* do *dataset Weather_Braga*. A remoção da coluna *precipitation* deveu-se ao facto desta apenas apresentar um único valor, o 0.

De modo a unir o resultado da preparação dos dados feita para os *datasets* anteriores, recorreu-se ao nodo *Joiner*, e uniram-se os *datasets* por *creation_date*, tendo-se efetuado, de seguida, a extração da data e do tempo.

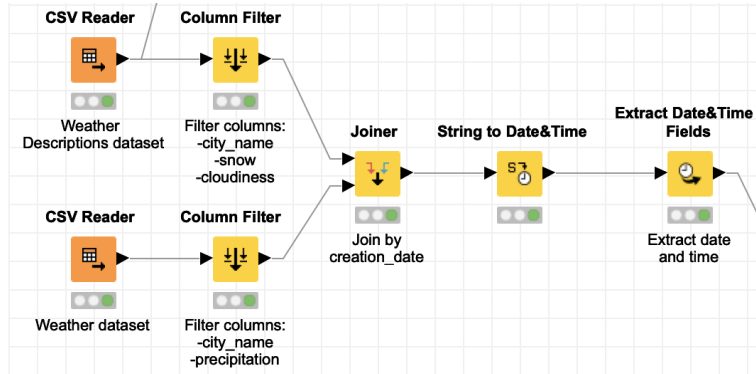


Figura 2: Preparação dos *datasets* *Weather_Descriptions_Braga* e *Weather_Braga*.

De seguida, procedeu-se à preparação do *dataset* *Traffic_Flow_Braga*, procedendo-se à remoção das colunas *city_name* e *road_name*, seguida da extração da data e hora e agrupamento dos dados por *road_num*, hora, dia do mês e mês. De modo a juntar este *dataset* ao obtido anteriormente, recorreu-se ao nodo *Joiner*, unindo-se os *datasets* por hora, dia do mês e mês, fazendo-se um *Left Outer Join*.

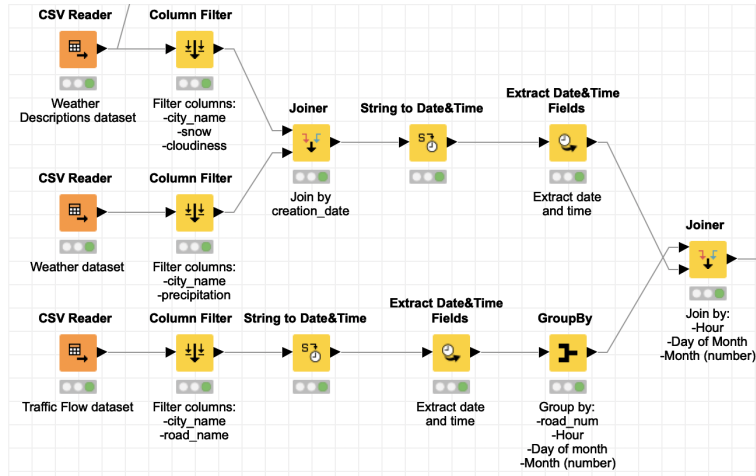


Figura 3: Preparação do *dataset* *Traffic_Flow_Braga*.

Após a junção dos *datasets*, eliminou-se a coluna *creation_date* e transformaram-se os valores "N/A", das colunas *rain*, *thunderstorm* e *atmosphere*, em *missing values*, recorrendo ao nodo *String Manipulation*. De seguida, fez-se um *merge* das colunas *rain* e *thunderstorm*, tendo-se alterado alguns dos valores ("trovoada com chuva fraca" → "chuva fraca", "trovoada com chuva forte" →

"chuva forte" e "trovoada" → "chuva"), tendo-se removido, no final, a coluna *thunderstorm*. Por fim, com os valores da coluna *clouds* construíram-se 6 intervalos:

1. céu claro:] – inf, 17[;
2. céu pouco nublado: [17, 34[;
3. nuvens dispersas: [34, 51[;
4. nuvens quebradas: [51, 68[;
5. nublado: [68, 85[;
6. muito nublado: [85, + inf[;

Tendo-se, de seguida, eliminado as colunas *sunrise*, *sunset* e *clouds*.

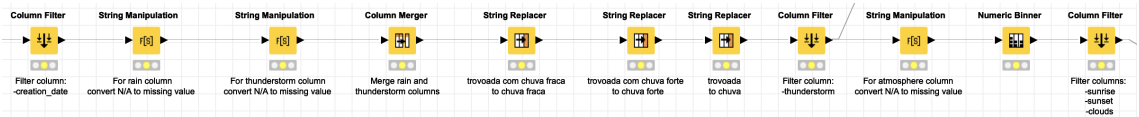


Figura 4: Preparação dos dados.

Por fim, tratou-se o *dataset* obtido após feito o tratamento do *dataset Traffic Incidents*. Procedeu-se à extração do dia e da hora e removeram-se as colunas irrelevantes. Após tratado este *dataset*, e recorrendo ao nodo *Joiner*, uniu-se este *dataset* com o obtido anteriormente por hora, dia do mês, mês e *road_num*. Deste modo, uniram-se os 4 *datasets* iniciais num único.

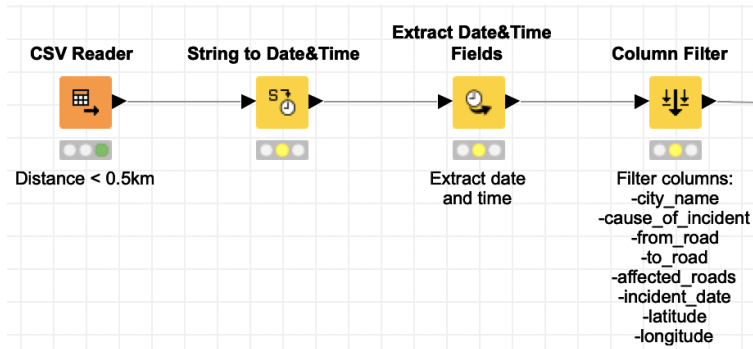


Figura 5: Preparação do *dataset* resultante do tratamento do *dataset Traffic Incidents Braga*.

Após se ter apenas um *dataset* eliminaram-se colunas que se acharam irrelevantes, servindo como apoio à decisão o nodo *Rank Correlation*, e aos valores *Undefined* da *feature descriptions* atribui-se o valor *Unknown Delay*, com o objetivo de diminuir a quantidade de atributos desta *feature*.

De seguida, e tendo em conta que as colunas *atmosphere* e *rain* apresentam muitos *missing values*, procedeu-se ao tratamento dos mesmos.

Começou-se, então, por tratar os *missing values* da coluna *atmosphere*, tendo-se separado o *dataset* em dois, recorrendo ao nodo *Rule-based Row Splitter*. Um *dataset* apresenta a coluna *atmosphere* apenas com *missing values* e o outro apresenta a coluna *atmosphere* com os vários valores. De seguida, utilizaram-se *Random Forest* para fazer a previsão dos *missing values*, para isso, particionou-se o *dataset* que apresentava os valores do atributo *atmosphere*, tendo-se usado 80% dos dados para treino. Após feita a previsão dos *missing values*, procedeu-se à previsão dos *missing values* do atributo *rain*, tendo-se utilizado o mesmo raciocínio.

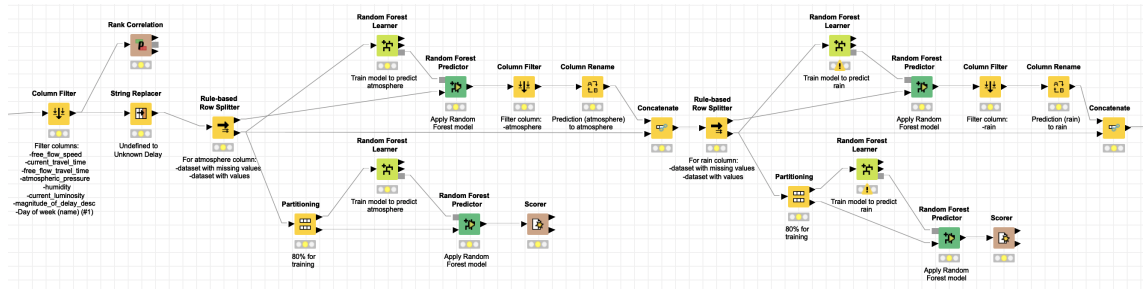


Figura 6: Previsão de *missing values*.

Por fim, recorrendo ao nodo *Duplicate Row Filter*, eliminaram-se linhas repetidas, efetuou-se o *Label Encoding* dos valores que estavam em *string* e trataram-se os *missing values*, substituindo-os por um valor *default*.

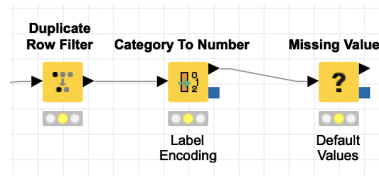


Figura 7: Tratamento final.

Para finalizar o tratamento de dados foi ainda necessário perceber quais os dias que estavam com falta de dados, procedendo-se à eliminação destes, com o intuito de se ter um *dataset* sem "buracos". Para isso, implementou-se o seguinte algoritmo:

```

1 i=0
2 for i in range(1,13):
3     for j in range(1,32):
4         L=df[(df['Month (number)']==i)&(df['Day of month']==j)].
           dropna()
5         L1=L[['Month (number)', 'Day of month', 'Hour', 'road_num']]
6         L1 = L1.drop_duplicates()
7         indexNames = df[(df['Month (number)']==i)&(df['Day of month'
           '']==j) ].index
8         if len(L1)>96:
9             if (len(L1)>=92 and len(L1)<96):
10                 i=i+1
11             if len(L1)<96:
12                 try:
13                     df.drop(indexNames, inplace=True)
14                 except:
15                     pass

```

Após feito todo o tratamento acima mencionado, o *dataset* está pronto para ser aplicado a uma rede que permita prever a *feature speed_diff*.