



# 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

<b>1</b>	<b>Introdução</b>	<b>1</b>
<b>2</b>	<b><i>Dataset</i></b>	<b>2</b>
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.

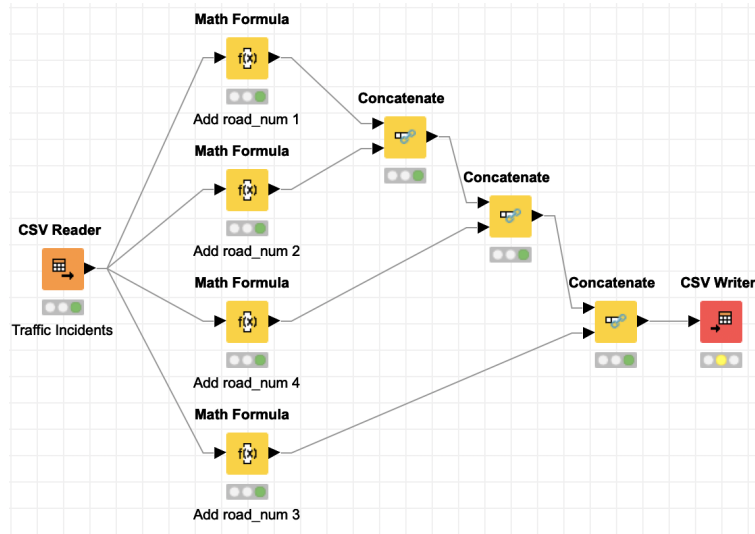


Figura 1: Preparação do *dataset Traffic\_Incidentes*.

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 dados errados, 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. Fez-se um estudo desta informação e verificou-se que estes dados dizem respeito a uma cidade que não pertence a Braga.

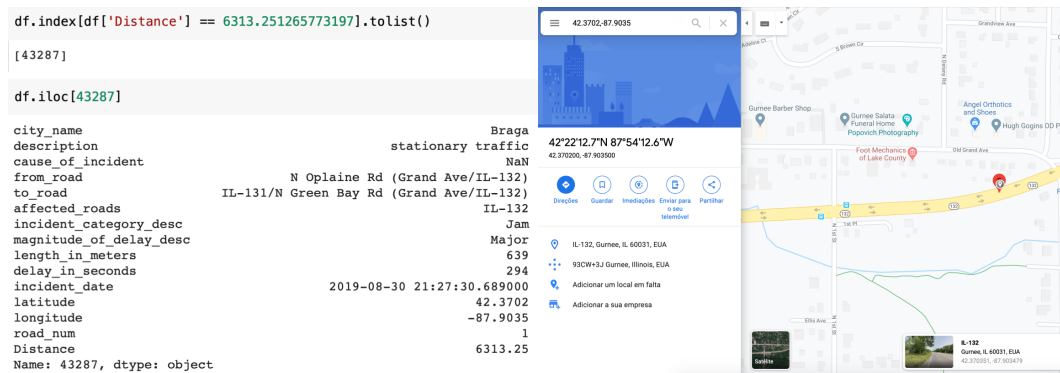


Figura 2: Dado mal classificado.

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.

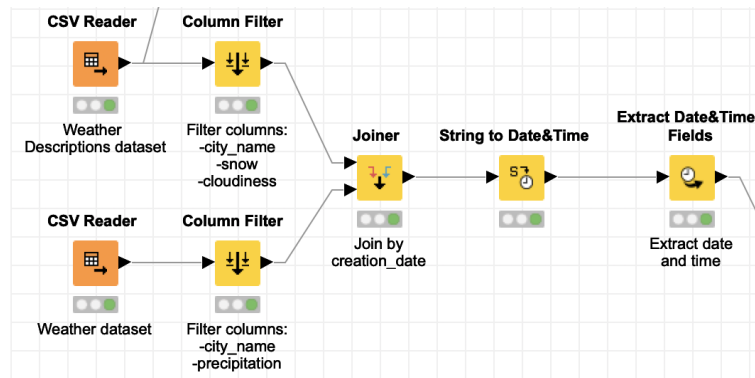
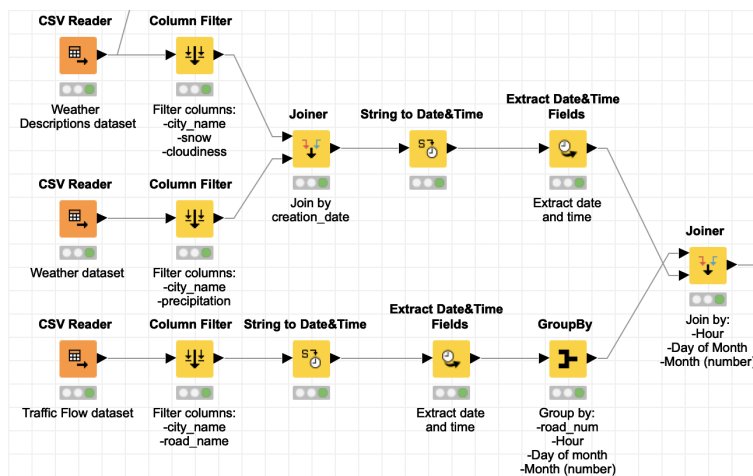


Figura 3: Preparação dos *datasets Weather\_Descriptions\_Braga* e *Weather\_Braga*.

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.

O *dataset Traffic\_Flow\_Braga* tinha registos de 20 em 20 minutos e o *dataset* obtido anteriormente tinha registos de hora em hora, assim, de modo a unir o *dataset* resultante de unir os *datasets Weather\_Descriptions\_Braga* e *Weather\_Braga*, optou-se por agrupar os registos do *dataset Traffic\_Flow\_Braga* por hora.

Assim, 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*. Optou-se por fazer um *Left Outer Join*, uma vez que não se queriam as condições atmosféricas de registos em que não havia dados de incidentes.



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, eliminaram-se as colunas *sunrise* e *sunset*.

Figura 5: Preparação dos dados.

Por fim, tratou-se o *dataset* cuja *feature Distance* tinha apenas valores inferiores a 0,5 km. 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.

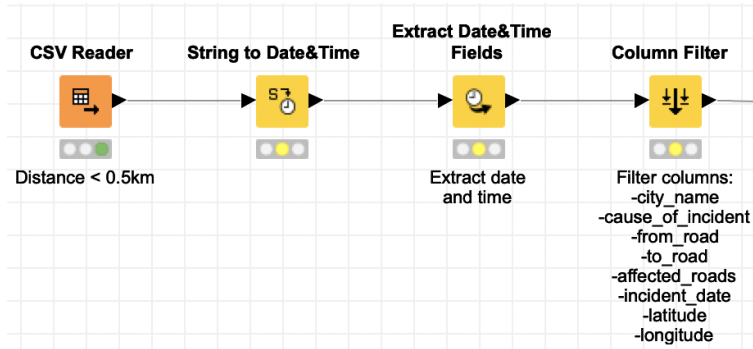


Figura 6: Preparação do *dataset* resultante do tratamento do *dataset Trafic Incidents Braga*.

Após se ter apenas um *dataset* eliminaram-se colunas que apresentavam uma correlação muito alta, de modo a tornar o *dataset* mais pequeno. Esta análise foi efetuada recorrendo ao nodo *Rank Correlation*.

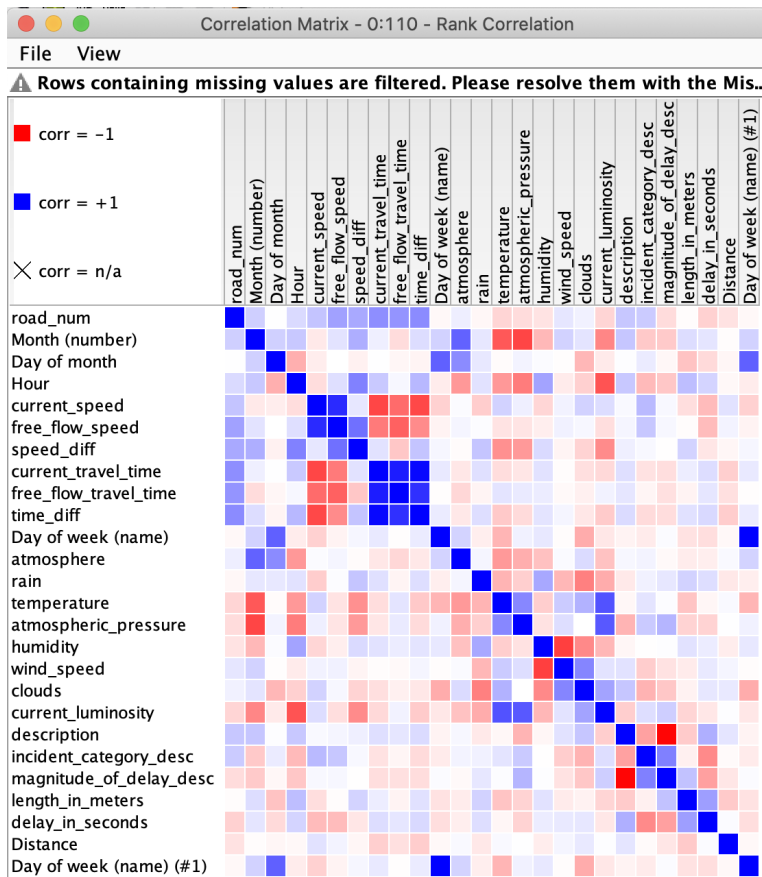


Figura 7: Análise da correlação entre as diferentes *features*.

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*, uma vez que este era o que apresentava menos *missing values*, 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*.

Com o intuito de perceber quais os melhores parâmetros a utilizar efetuou-se o *tunning* do modelo.

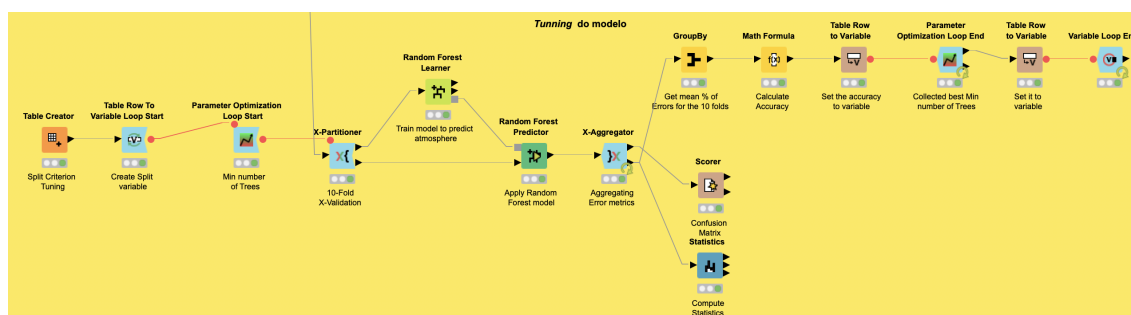


Figura 8: *Tunning* do modelo.

Após efetuado o *tunning* do modelo, concluiu-se que este apresentava melhores valores se fosse treinado com 60 árvores e usando como critério de *split* o *Information Gain*, tendo-se uma *accuracy* de cerca 99,5%

Data table of flow variables - 0:129 - Va...				Output data - 0:133 - Math Form...			
File	Hilite	Navigation	View	File	Hilite	Navigation	View
Table "default" - Rows: 3				Table "default" - Rows: 7			
Row ID	Trees	Split		Row ID	Error in %	Accuracy	
Row0	60	InformationGain		Row0	0	99.446	
Row1	160	InformationGainRatio		Row1	0.298	99.446	
Row2	160	Gini		Row2	0.299	99.446	
				Row3	0.595	99.446	
				Row4	0.597	99.446	
				Row5	0.896	99.446	
				Row6	1.194	99.446	

Figura 9: Melhores parâmetros para construir o modelo.

Por fim, fez-se a previsão dos *missing values* da *feature atmosphere*, usando 100% dos dados para treinar o modelo.

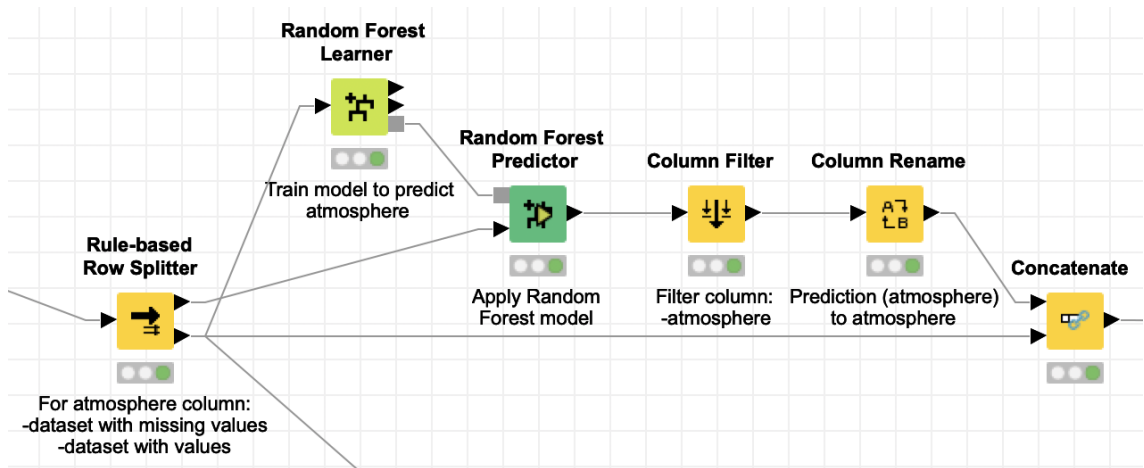


Figura 10: Previsão dos *missing values* da *feature atmosphere*.

Após feita a previsão dos *missing values* para a *feature atmosphere*, procedeu-se à previsão dos *missing values* do atributo *rain*, utilizando os mesmo parâmetros, tendo-se obtido uma *accuracy* de cerca de 97,8%.

Para finalizar o tratamento de dados, no *Knime*, recorrendo ao nó *Duplicate Row Filter*, eliminaram-se linhas repetidas e efetuou-se o *Label Encoding* dos valores correspondentes às *features*: *Day of week (name)*, *description*, *incident\_category\_desc*, *atmosphere* e *rain*.

De seguida, trataram-se os *missing values*, substituindo-os por um valor *default*. Os *missing values* existentes correspondiam a dias/horas onde não tinham ocorrido incidentes.

Por fim, observou-se que existiam dias com horas repetidas, devido ao facto de para uma mesma hora existir mais do que um incidente. Assim, de modo a que cada dia tivesse apenas 24 horas, recorreu-se ao nó *GroupBy* para agrupar os incidentes, optando-se por ficar com o incidente que estava mais próximo da rua em estudo, uma vez que seria esse que mais influenciaria o atributo *speed\_diff*.

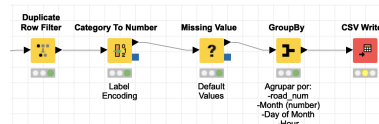


Figura 11: Tratamento final.

Após feito este tratamento, recorrendo ao nodo *Pie chart (local)* percebeu-se que o *dataset* apresentava dias e horas em falta como, por exemplo, o mês de Maio.

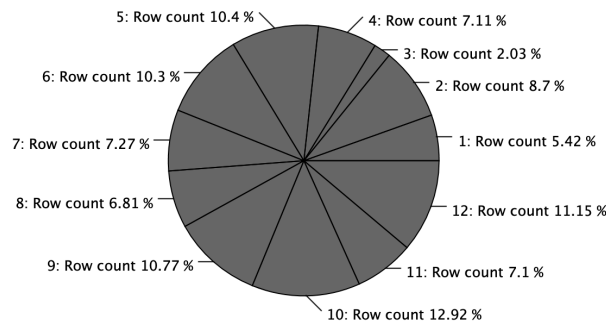


Figura 12: Dias em falta no *dataset*.

Deste modo, foi então necessário perceber quais os dias que estavam incompletos, ou seja, quais os dias que não tinham as 24 horas preenchidas, 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             try:
10                 df.drop(indexNames, inplace=True)
11             except:
12                 pass

```

Com o objetivo de ter a certeza que se treina o modelo com dias seguidos, desenvolveu-se o seguinte código *python*:

```

1 n_future = 24 # next 4 days temperature forecast
2 n_past = 24*7 # Past 30 days
3
4 x_train = []
5 y_train = []
6 label = df_1['speed_diff']
7
8

```

```

9 for i in range(0, len(df_1)-n_past-n_future+1):
10     dias = df_1.iloc[i : i + n_past+24]
11     mes = dias.iloc[0]['Month (number)']
12     dia_1 = dias.iloc[0]['Day of month']
13     dia_168 = dias.iloc[168]['Day of month']
14     if (mes == 4 or mes == 6 or mes == 9 or mes == 11) and (
15         dia_168 - dia_1 == 7 or dia_168 - dia_1 == -25):
16         x_train.append(df_1.iloc[i : i + n_past])
17         y_train.append(label.iloc[i + n_past : i + n_past +
18             n_future ])
19     elif (mes == 1 or mes == 3 or mes == 5 or mes == 7 or mes
20         == 8 or mes == 10 or mes == 12) and (dia_168 - dia_1 == 7
21         or dia_168 - dia_1 == -24):
22         x_train.append(df_1.iloc[i : i + n_past])
23         y_train.append(label.iloc[i + n_past : i + n_past +
24             n_future ])
25     elif mes == 2 and (dia_168 - dia_1 == 7 or dia_168 - dia_1
26         == -22):
27         x_train.append(df_1.iloc[i : i + n_past])
28         y_train.append(label.iloc[i + n_past : i + n_past +
29             n_future ])

```

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