

Projet

Le projet consiste à prévoir le nombre de vélos loués à chaque heure dans des bornes libre-services de la ville (système type Vélib'). La variable cible est ici la variable **count**.

Voici un descriptif de l'ensemble des variables :

- datetime** - date et heure du relevé
- season** - 1 = printemps, 2 = été, 3 = automne, 4 = hiver
- holiday** - indique s'il y a eu un jour de vacances ou non
- workingday** - indique si le jour est travaillé (0 si week-end ni vacances)
- weather** - 1 : Dégradé à nuageux, 2 : Brouillard, 3 : Légère pluie ou neige, 4 : Fortes averse ou neiges
- temp** - température en degrés Celsius
- atemp** - température ressentie en degrés Celsius
- humidity** - taux d'humidité
- windspeed** - vitesse du vent
- casual** - nombre de locations d'utilisateurs non abonnés
- registered** - nombre de locations d'utilisateurs abonnés
- count** - nombre total de locations de vélos

L'objectif du projet est de mener à bien la création d'un modèle qui pourrait théoriquement être déployé en production. Les étapes d'exploration des données, de traitement et de préprocessing ne sont bien entendu pas à négliger. Il ne s'agit pas d'une compétition de type Kaggle, le projet ne sera pas uniquement noté sur la performance du modèle, mais plutôt sur votre approche complète et la justification de chacun de vos choix.

Comme vu durant le cours, soyez faites attention à certains points :

- quel type de problème dois-je traiter ?
- feature engineering - est-ce que j'utilise les données correctement, toutes les données ?
- data leakage - est-ce qu'une de mes features n'est pas trop explicative ?
- ai-je bien traité toutes les données correctement ?
- est-ce que mon modèle est adapté ?
- etc, etc, etc

Soyez vigilant à expliquer et justifier votre démarche à l'aide de visualisation, de commentaires dans vos codes (pensez aux cellules markdown), etc

```
In [41]: import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib
```

```
In [42]: df = pd.read_csv("../data/input/velo.csv")

df.head()
```

```
Out[42]:
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
0	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0	3	13	16
1	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0	8	32	40
2	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0	5	27	32
3	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0	3	10	13
4	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0	0	1	1

Analyse de la structure du dataframe

```
In [43]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   datetime              10886 non-null  object
 1   season                10886 non-null  int64
 2   holiday               10886 non-null  int64
 3   workingday            10886 non-null  int64
 4   weather               10886 non-null  int64
 5   temp                 10886 non-null  float64
 6   atemp                10886 non-null  float64
 7   humidity              10886 non-null  int64
 8   windspeed            10886 non-null  float64
 9   casual                10886 non-null  int64
10   registered            10886 non-null  int64
11   count                 10886 non-null  int64
dtypes: float64(2), int64(8), object(1)
memory usage: 1020.7+ KB
```

Analyse descriptive des données

```
In [44]: df.describe()

df.describe()

Out[44]:
```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
count	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000
mean	2.506614	0.028569	0.680875	1.818427	20.23086	23.655084	61.886460	12.799395	36.021955	155.552	155.552
std	1.116174	0.166599	0.466159	0.633839	7.79159	8.474601	19.245033	8.164537	49.960477	151.039	151.039
min	1	0	0	1	0.82000	0.820000	0.000000	0.000000	0.000000	0.000000	0.0000
25%	2.000000	0.000000	0.000000	1.000000	13.94000	16.665000	47.000000	7.001500	4.000000	36.0000	36.000
50%	3.000000	0.000000	1.000000	1.000000	20.50000	24.240000	62.000000	12.998000	17.000000	118.000	118.000
75%	4.000000	0.000000	1.000000	2.000000	26.24000	31.060000	77.000000	16.997900	49.000000	222.000	222.000
max	4.000000	1.000000	1.000000	4.000000	41.00000	45.455000	100.000000	56.996900	367.000000	886.000	886.000

Nombre de valeurs uniques par colonnes

```
In [45]: df.apply(lambda x: len(x.unique()))

df.apply(lambda x: len(x.unique()))

Out[45]:
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
datetime	10886											
season	4											
holiday	2											
workingday	2											
weather	4											
temp	49											
atemp	60											
humidity	89											
windspeed	28											
casual	309											
registered	731											
count	822											
dtype:	int64											

Split de la colonne datetime en year, month, day, time afin de pouvoir analyser les variables indépendamment

```
In [46]: df["datetime"] = pd.to_datetime(df["datetime"])
df["year"] = df["datetime"].dt.year
df["month"] = df["datetime"].dt.month
df["day"] = df["datetime"].dt.day
df["hour"] = df["datetime"].dt.hour
df = df.drop(columns=["datetime"])
```

```
In [47]: df.describe()

df.describe()

Out[47]:
```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	year	month	day	hour
count	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000
mean	2.506614	0.028569	0.680875	1.818427	20.23086	23.655084	61.886460	12.799395	36.021955	155.552	155.552				
std	1.116174	0.166599	0.466159	0.633839	7.79159	8.474601	19.245033	8.164537	49.960477	151.039	151.039				
min	1	0	0	1	0.82000	0.820000	0.000000	0.000000	0.000000	0.000000	0.0000				
25%	2.000000	0.000000	0.000000	1.000000	13.94000	16.665000	47.000000	7.001500	4.000000	36.0000	36.000				
50%	3.000000	0.000000	1.000000	1.000000	20.50000	24.240000	62.000000	12.998000	17.000000	118.000	118.000				
75%	4.000000	0.000000	1.000000	2.000000	26.24000	31.060000	77.000000	16.997900	49.000000	222.000	222.000				
max	4.000000	1.000000	1.000000	4.000000	41.00000	45.455000	100.000000	56.996900	367.000000	886.000	886.000				

```
In [48]: df.head()

df.head()

Out[48]:
```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	year	month	day	hour
0	1	0	0	1	9.84	14.395	81	0.0	3	13	16	2011	1	1	0
1	1	0	0	1	9.02	13.635	80	0.0	8	32	40	2011	1	1	1
2	1	0	0	1	9.02	13.635	80	0.0	5	27	32	2011	1	1	2
3	1	0	0	1	9.84	14.395	75	0.0	3	10	13	2011	1	1	3
4	1	0	0	1	9.84	14.395	75	0.0	0	1	1	2011	1	1	4

Les colonnes weather et season sont catégorisées comme numérique / continue alors qu'il s'agit de variable qualitative / discrète

```
In [49]: df[["weather", "season"]] = df[["weather", "season"]].astype("category")

df[["weather", "season"]] = df[["weather", "season"]].astype("category")
```

Premières analyses : nous obtenons 6 variables catégorielles, 8 variables numériques et une variable cible numérique

Ajout de quelques graphiques

Analyse de la distribution de la cible afin de repérer la part des outliers

```
In [50]: sns.distplot(df["count"])

sns.distplot(df["count"])

Out[50]:
```

```
In [51]: sns.boxplot(x=df["count"])

sns.boxplot(x=df["count"])

Out[51]:
```

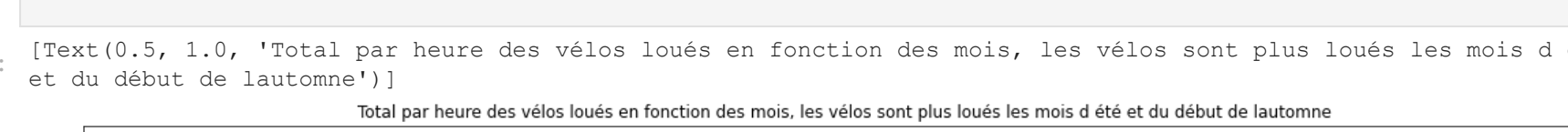
Analyse des autres variables

```
In [52]: fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='hour', y='count', ax=ax)
ax.set(title='Nombre de vélos loués toutes les heures, on note un pic à 8h et 17h/18h')

fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='hour', y='count', ax=ax)
ax.set(title='Nombre de vélos loués toutes les heures, on note un pic à 8h et 17h/18h')
```

```
Out[52]:
```

[Text(0.5, 1.0, 'Nombre de vélos loués toutes les heures, on note un pic à 8h et 17h/18h')]

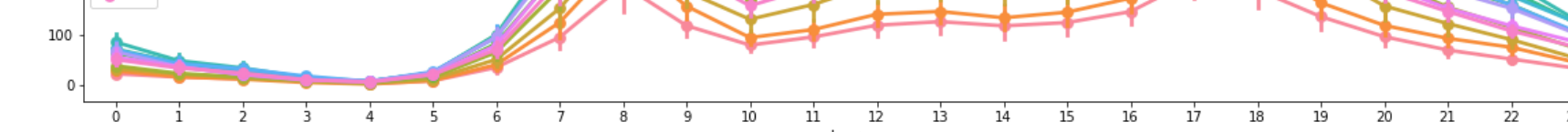


```
In [53]: fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='month', y='count', hue='month', ax=ax)
ax.set(title='Total par heure des vélos loués en fonction des mois, les vélos sont plus loués les mois d été et d'automne')

fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='month', y='count', hue='month', ax=ax)
ax.set(title='Total par heure des vélos loués en fonction des mois, les vélos sont plus loués les mois d été et d'automne')
```

```
Out[53]:
```

[Text(0.5, 1.0, 'Total par heure des vélos loués en fonction des mois, les vélos sont plus loués les mois d été et d'automne')]

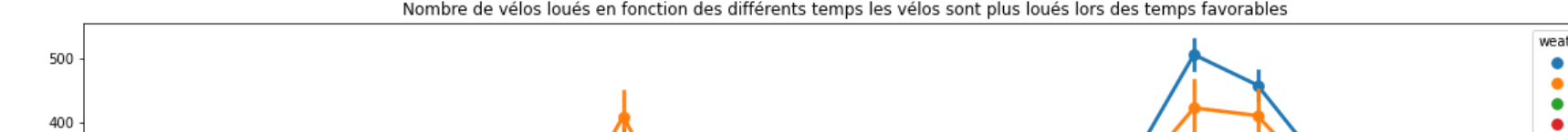


```
In [54]: fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='hour', y='count', hue='season', ax=ax)
ax.set(title='Nombre de vélos loués en fonction des différents temps les vélos sont plus loués lors des temps favorables')

fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='hour', y='count', hue='season', ax=ax)
ax.set(title='Nombre de vélos loués en fonction des différents temps les vélos sont plus loués lors des temps favorables')
```

```
Out[54]:
```

[Text(0.5, 1.0, 'Nombre de vélos loués en fonction des différents temps les vélos sont plus loués lors des temps favorables')]

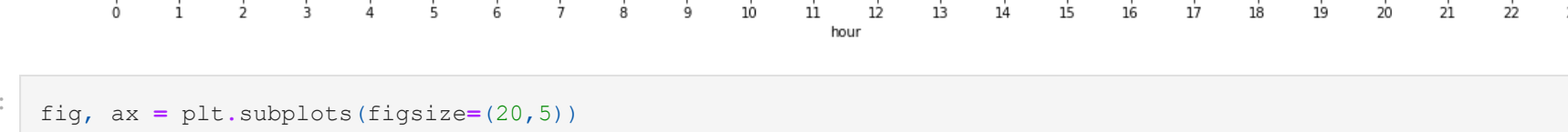


```
In [55]: fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='season', y='count', hue='season', ax=ax)
ax.set(title='Nombre de vélos loués en fonction des différentes saisons, les vélos sont moins loués en hiver')

fig, ax = plt.subplots(figsize=(20,5))
sns.pointplot(data=df, x='season', y='count', hue='season', ax=ax)
ax.set(title='Nombre de vélos loués en fonction des différentes saisons, les vélos sont moins loués en hiver')
```

```
Out[55]:
```

[Text(0.5, 1.0, 'Nombre de vélos loués en fonction des différentes saisons, les vélos sont moins loués en hiver')]



```
In [56]: sns.barplot(data=df, x='season', y='count', ci='sd', palette='dark')
ax.set(title='Nombre de vélos loués en fonction des saisons')

sns.barplot(data=df, x='season', y='count', ci='sd', palette='dark')
ax.set(title='Nombre de vélos loués en fonction des saisons')
```

```
Out[56]:
```

[Text(0.5, 1.0, 'Nombre de vélos loués en fonction des saisons')]

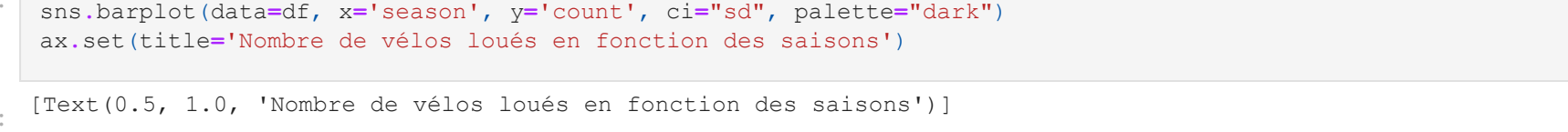


```
In [57]: sns.barplot(data=df, x='weather', y='count', ci='sd', palette='dark')
ax.set(title='Nombre de vélos loués en fonction des temps')

sns.barplot(data=df, x='weather', y='count', ci='sd', palette='dark')
ax.set(title='Nombre de vélos loués en fonction des temps')
```

```
Out[57]:
```

[Text(0.5, 1.0, 'Nombre de vélos loués en fonction des temps')]



Analyse des corrélations afin de voir s'il faudra réduire le nombre de dimension

```
In [58]: corr = df.corr()

corr = df.corr()

Out[58]:
```

Les colonnes atemp / temp sont fortement positivement corrélées et les colonnes workingday / holiday sont fortement négativement corrélées, je décide donc de supprimer les colonnes atemp et holiday. La colonne cible count correspond à la somme des colonnes registered et casual, je décide donc également de ne pas me préoccuper de ces colonnes

```
In [59]: df = df.drop(columns=["atemp", "holiday", "casual", "registered"])

df = df.drop(columns=["atemp", "holiday", "casual", "registered"])
```

Dans le premier graphique, nous observons des outliers, nous allons les supprimer grâce à la médiane et l'écart interquartile

```
In [60]: print("Shape initiale: {}".format(df.shape))
q1 = df["count"].quantile(0.25)
q3 = df["count"].quantile(0.75)
iqr = q3 - q1
df = df.loc[(df["count"] >= q1 - (1.5 * iqr)) & (df["count"] <= q3 + (1.5 * iqr))]
print("Shape après avoir supprimé les outliers: {}".format(df.shape))
sns.distplot(df["count"])

print("Shape initiale: (10886, 11)")
print("Shape après avoir supprimé les outliers: (10586, 11)")
sns.distplot(df["count"])

sns.distplot(df["count"])

Out[60]:
```

```
In [61]: df.describe()

df.describe()

Out[61]:
```

	workingday	temp	humidity	windspeed	count	year	month	day	hour
count	10586.000000	10586.000000	10586.000000	10586.000000	10586.000000	10586.000000	10586.000000	10586.000000	10586.000000
mean	0.676459	20.061494	62.165124	12.776699	175.717079	2011.487814	6.494804	9.977045	11.456641
std	0.467849	7.781496	19.231315	8.173450	156.360023	0.499875	3.464977	5.475961	6.960409
min	0.000000	0.820000	0.000000	0.000000	1.000000	2011.000000	1.000000	1.000000	0.000000
25%	0.000000	13.940000	47.000000	7.001500	40.000000	2011.000000	3.000000	5.000000	11.456641
50%	1.000000	20.500000	62.000000	12.998000	138.000000	2011.000000	6.000000	10.000000	11.000000
75%	1.000000	26.240000	78.000000	16.997900	270.000000	2012.000000	10.000000	15.000000	21.000000
max	1.000000	41.000000	100.000000	56.996900	647.000000	2012.000000	12.000000	19.000000	23.000000

Préparation du dataset pour mettre en place les algorithmes

```
In [62]: from prettytable import PrettyTable
import xgboost as xgb

# sklearn metrics for regression models
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

# sklearn ensemble models
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingRegressor

# sklearn tree model
from sklearn.tree import DecisionTreeRegressor

# sklearn linear model
from sklearn.linear_model import LinearRegression
```

Catégorisation des variables

Les colonnes de type "category" doivent être dumifiées, on utilise le paramètre drop_first=True car il est induit que si les x-1 colonnes sont à 0 alors la dernière est à 1. Cela réduit le nombre de colonnes.

Les variables weather et season sont catégorisées comme numériques : l'encodage est purement numérique, cela correspond à la factorisation. Cela peut poser problème car l'algorithme assumera que deux valeurs proches sont plus similaires que deux valeurs éloignées, ce qui n'est pas forcément le cas. C'est pourquoi les variables sont dumifiées.

Les variables year, month, day et hour n'ont pas été encodées car elles sont de nature cycliques et elles ont beaucoup de valeurs différentes. Le jour de la semaine augmenterait considérablement le nombre de variables. C'est pourquoi la conversion des valeurs en sin/cos est pertinente. En effet, de cette manière, l'algorithme considèrera les heures 0 et 23 comme étant proches.

```
In [63]: import math

def transformation(column):
    max_value = column.max()
    sin_values = [math.sin(2*math.pi*x/max_value) for x in list(column)]
    cos_values = [math.cos(2*math.pi*x/max_value) for x in list(column)]
    return sin_values, cos_values

transformation(column)

Out[63]:
```

holiday	workingday	temp	atemp	humidity	windspeed	casual	registered	count	year	month	day	hour
1	0	9.84	14.395	81	0.0000	3	13	16	2011	1	1	0
1	0	9.02	13.635	80	0.0000	8	32	40	2011	1	1	1
1	0	9.02	13.635	80	0.0000	5	27	32	2011	1	1	2
1	0	9.84	14.395	75	0.0000	3	10	13	2011	1	1	3
1	0	9.84	14.395	75	0.0000	0	1	1	2011	1	1	4

```
Out[63]:
```

[Text(0

Out[80]:

	Model	Dataset	MSE	RMSE	R ² score	MAE
XGBRegressor	training		795.43	28.20	0.97	19.10
XGBRegressor	validation		1548.69	39.35	0.94	25.31

Le modèle XGBRegressor obtient un très bon score R² sur le jeu d'entraînement, il overfit encore un peu les données.

GradientBoostingRegressor

In [81]:

```
gbr_params = {'n_estimators': 500,
              'max_depth': 6,
              'min_samples_split': 5,
              'learning_rate': 0.01}

table = PrettyTable()
table.field_names = ["Model", "Dataset", "MSE", "RMSE", "R2 score", "MAE"]

model_gbr = GradientBoostingRegressor(**gbr_params).fit(x_train, y_train)

def evaluate(x, y, dataset):
    y_pred = model_gbr.predict(x)

    mse = mean_squared_error(y, y_pred)
    score = model_gbr.score(x, y)
    rmse = np.sqrt(mse)
    mae = mean_absolute_error(y, y_pred)

    table.add_row([type(model_gbr).__name__, dataset, format(mse, '.2f'), format(rmse, '.2f'), format(score, '.2f'), format(mae, '.2f')])

evaluate(x_train, y_train, 'training')
evaluate(x_test, y_test, 'validation')

table
```

Out[81]:

	Model	Dataset	MSE	RMSE	R ² score	MAE
GradientBoostingRegressor	training		1382.72	37.19	0.94	25.61
GradientBoostingRegressor	validation		2014.46	44.88	0.92	29.81

Avec ce modèle, les résultats sur les données d'entraînement et de test sont assez similaires, cela signifie que le modèle est plus général, même si le score du modèle est légèrement moins élevé, il est plus général. Le model DecisionTree a également été testé car les résultats étaient meilleurs, mais il overfittait vraiment les données. Le model RandomForestRegressor overfit également trop les données.

Le modèle que je décide de garder est celui ci.

Conclusion

Dans le premier modèle, les variables heures, jours et températures sont celles ayant le plus d'importance dans le modèle, plus de 75% du modèle est expliqué par ces trois variables. On cherche à prédire le nombre de vélos à allouer par heure. Comme on a pu l'observer dans les graphiques, on observe des pics le matin et en fin de journée.

In [82]:

```
y_test = np.array(y_test.values.tolist())
```

In [83]:

```
x_test["hour_cos"] = x_test["hour_cos"].replace(hour_dict)
```

In [84]:

```
fig, ax = plt.subplots(figsize=(30,8))
sns.pointplot(x=x_test["hour_cos"], y=y_test, ax=ax, color='red')
sns.pointplot(x=x_test["hour_cos"], y=y_pred, ax=ax)
ax.set(title="Nombre de vélos loués toutes les heures, on note un pic à 8h et 18h, en bleu la prédiction")
```

Out[84]:

[Text(0.5, 1.0, 'Nombre de vélos loués toutes les heures, on note un pic à 8h et 18h, en bleu la prédiction')]



In [85]:

```
jupyter-nbconvert --to PDFviaHTML GAYRAUD_EVA.ipynb
```

File "C:\Users\lucaga\AppData\Local\Temp\ipykernel_5568\479933117.py", line 1
jupyter-nbconvert --to PDFviaHTML GAYRAUD_EVA.ipynb
SyntaxError: invalid syntax

In []: