

## S09 T02: Apprenentage Supervisat - Regressions part 2

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
pd.set_option('display.max_columns', None)
import warnings
warnings.filterwarnings('ignore')
pd.set_option('display.float_format', '{:.10f}'.format)
```

I will load again the dataset from the previous exercise called data and I will apply again dummy transformation here because it takes less memory space if I do it this way (because of types of the dummy columns). I could also change dtypes from here but it takes a lot of time so it is faster doing it this way.

```
# open dataset from the previous exercise
data = pd.read_csv('data.csv')
data.head()
```

	ArrDelay	ActualElapsedTime	CRSElapsedTime	AirTime	DepDelay	Distance	Speed	Taxin	TaxiOut
0	14.0000000000	128.0000000000	150.0000000000	116.0000000000	8.0000000000	8000000000	810	418	9700000000
1	2.0000000000	128.0000000000	145.0000000000	113.0000000000	19.0000000000	810	430.9500000000	5.0000000000	10.0000000000
2	14.0000000000	96.0000000000	90.0000000000	77.0000000000	8.0000000000	515	406.5800000000	3.0000000000	17.0000000000
3	34.0000000000	90.0000000000	76.0000000000	34.0000000000	0.0000000000	515	401.3000000000	3.0000000000	10.0000000000
4	11.0000000000	101.0000000000	115.0000000000	87.0000000000	25.0000000000	688	474.4800000000	4.0000000000	10.0000000000

```
# do the sample explained in previous exercise
data_sample = data.groupby(['UniqueCarrier'], group_keys=False).apply(lambda x : x.sample(frac=0.1, random_state=1))
```

```
# now i will generate the dummy columns that are already explained in the previous exercise
# delay interval
def delay_interval(x):
    if x < 15:
        return 'Less than 15"
    elif x < 30:
        return 'Between 15 and 30"
    elif x < 60:
        return 'Between 30 and 60"
    else:
        return 'More than 60"

data_sample['DelayInterval'] = data_sample['ArrDelay'].apply(delay_interval)
data_sample.head()
```

	ArrDelay	ActualElapsedTime	CRSElapsedTime	AirTime	DepDelay	Distance	Speed	Taxin	TaxiOut
143019	21.0000000000	95.0000000000	88.0000000000	63.0000000000	14.0000000000	364	346.6700000000	11.0000000000	21.0000000000
1886001	24.0000000000	139.0000000000	126.0000000000	105.0000000000	11.0000000000	813	464.5700000000	18.0000000000	16.0000000000
1881989	2.0000000000	86.0000000000	102.0000000000	65.0000000000	18.0000000000	490	452.3100000000	9.0000000000	12.0000000000
1030789	37.0000000000	114.0000000000	90.0000000000	87.0000000000	13.0000000000	442	304.8300000000	4.0000000000	23.0000000000
687997	10.0000000000	90.0000000000	88.0000000000	61.0000000000	8.0000000000	382	375.7400000000	7.0000000000	22.0000000000

```
# creating the dummy columns
# we'll use the param drop_first = True so we don't have duplicated information
data_sample = pd.get_dummies(data = data_sample, columns = ['UniqueCarrier', 'Origin', 'Dest'], drop_first=True)
# map dummies for DelayInterval
data_sample['DelayInterval'] = data_sample['DelayInterval'].map({'Less than 15": 0, 'Between 15 and 30": 1, 'Between 30 and 60": 2, 'More than 60": 3})
data_sample.info()
```

<class 'pandas.core.frame.DataFrame'>  
Int64Index: 192825 entries, 143019 to 1541485  
Columns: 637 entries, ArrDelay to Test\_VMF  
dtypes: float64(15), int64(7), uint8(815)  
memory usage: 146.9 MB

Now I will delete the columns that show multicollinearity (explained in the previous Notebook on the VIF process)

```
data_cleaned = data_sample.copy()
```

```
data_cleaned.drop(['ActualElapsedTime', 'CRSElapsedTime', 'Taxin', 'LateAircraftDelay', 'AirTime', 'CRSDepTime', 'Speed', 'CRSAirTime', 'DepTime', 'axis = 1, inplace = True])
```

```
data_cleaned.shape
```

(192825, 628)

Now we have the same dataset that we had at the end of our previous Notebook.

## Exercici 5

Realitza algun procés d'enginyeria de variables per millorar-ne la predicció

To improve our predictions, we could scale our data, let's check the data:

```
data_scaled = data_cleaned.copy()
```

```
data_scaled.head()
```

	ArrDelay	DepDelay	Distance	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	AirTime
143019	21.0000000000	14.0000000000	364	16.0000000000	14.0000000000	0.0000000000	7.0000000000	0.0000000000	1746.0000000000
1886001	24.0000000000	13.0000000000	813	21.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	1479.0000000000
1881989	2.0000000000	18.0000000000	490	12.0000000000	-0.3600000000	0.0000000000	0.0000000000	0.0000000000	1209.0000000000
1030789	37.0000000000	13.0000000000	442	23.0000000000	13.0000000000	0.0000000000	24.0000000000	0.0000000000	2102.0000000000
687997	10.0000000000	8.0000000000	382	22.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	1208.0000000000

```
data_scaled.columns[1:6]
```

Index(['ArrDelay', 'DepDelay', 'Distance', 'TaxiOut', 'CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay', 'AirTime', 'DayOfMonth', 'DayOfWeek', 'DelayInterval', 'UniqueCarrier\_AQ', 'UniqueCarrier\_AB', 'dtype='object']

We can scale the data in different ways:

- We can use **RobustScaler** on those features that refer to delay and have heavy outliers: 'ArrDelay', 'DepDelay', 'Distance', 'TaxiOut', 'CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay'.
- We can use **MinMaxScaler** with those features that refer to the time of the flight, these features don't have outliers: 'AirTime', 'Month', 'DayOfMonth', 'DayOfWeek', they refer to the dates and hours of the flights.

We won't scale dummy columns.

## RobustScaler

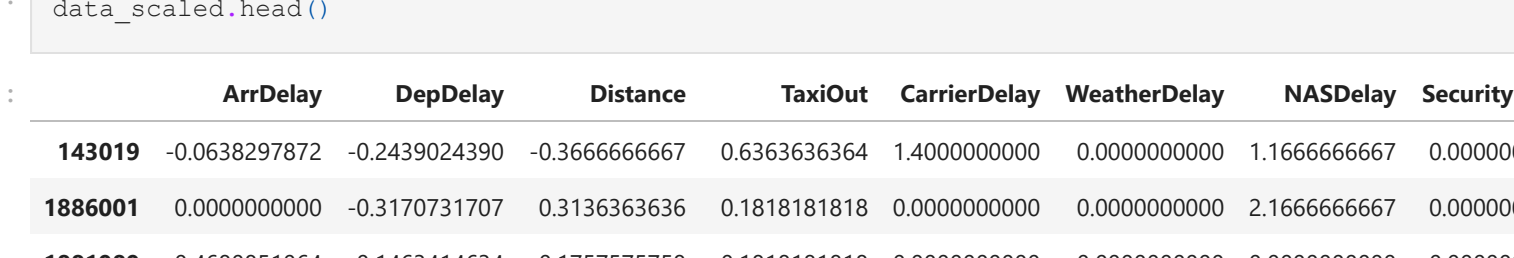
```
cols_robust = ['ArrDelay', 'DepDelay', 'Distance', 'TaxiOut', 'CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay']
```

```
data_scaled[cols_robust].describe().round(2)
```

	ArrDelay	DepDelay	Distance	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay
count	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000
mean	42.2200000000	43.1300000000	636.636364	18.2300000000	0.0000000000	0.0000000000	0.0000000000	9.7100000000
std	56.9200000000	53.5600000000	573.9400000000	14.2300000000	35.8100000000	17.9000000000	0.0000000000	28.2100000000
min	-59.0000000000	6.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000
25%	9.0000000000	12.0000000000	338.0000000000	10.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000
50%	24.0000000000	24.0000000000	606.0000000000	14.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000
75%	56.0000000000	53.0000000000	498.0000000000	21.0000000000	10.0000000000	10.0000000000	0.0000000000	6.0000000000
max	1655.0000000000	1597.0000000000	996.0000000000	356.0000000000	1312.0000000000	1148.0000000000	0.0000000000	1289.0000000000

Comparing max values and quartiles we can see we have outliers in this columns, let's draw some boxplots to confirm this.

```
plt.figure(figsize=(12,9))
sns.boxplot(data = data_scaled[cols_robust]);
```



Some features have more outliers than others, but let's try this approach and see how it works:

```
from sklearn import preprocessing
robust_scaler = preprocessing.RobustScaler()
data_scaled[cols_robust] = robust_scaler.fit_transform(data_scaled[cols_robust])
```

```
data_scaled.head()
```

	ArrDelay	DepDelay	Distance	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	AirTime
143019	-0.636297872	-0.249302490	0.366666667	0.636363634	1.400000000	0.000000000	1.166666667	0.000000000	0.727264271
1886001	0.000000000	-0.317073107	0.313636363	0.181818181	0.000000000	0.000000000	0.166666667	0.000000000	0.591076916
1881989	-0.468051064	-0.1463414634	-0.175575758	-0.181818181	0.000000000	0.000000000	0.000000000	0.000000000	0.503543143
1030789	0.276595747	-0.2682526289	-0.248484848	0.818181818	1.300000000	0.000000000	4.000000000	0.000000000	0.875781577
687997	-0.297872304	-0.3902439024	-0.339393939	0.727272727	0.000000000	0.000000000	0.000000000	0.000000000	0.503126326

	ArrDelay	DepDelay	Distance	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay
count	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000
mean	0.3990000000	0.4700000000	0.2400000000	0.3800000000	1.2300000000	2.3900000000	1.6200000000	0.4700000000
std	1.1710000000	1.3400000000	0.8700000000	1.2700000000	3.5800000000	17.9000000000	0.0000000000	4.0000000000
25%	-0.3200000000	-0.2900000000	-0.4000000000	-0.3600000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000
50%	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000
75%	0.5800000000	0.7100000000	0.5900000000	0.6400000000	1.0000000000	0.0000000000	1.0000000000	1.0000000000
max	34.7000000000	38.3700000000	6.6000000000	31.0900000000	131.2000000000	1148.0000000000	214.8300000000	0.0000000000

We see the data has been scaled in a way that preserves the information about outliers.

## MinMaxScaler

For the columns that refer to the time of the flight ('AirTime', 'Month', 'DayOfMonth', 'DayOfWeek') I will use MinMaxScaler.

```
from sklearn.preprocessing import MinMaxScaler
cols_minmax = ['AirTime', 'Month', 'DayOfMonth', 'DayOfWeek', 'DelayInterval']
mms = MinMaxScaler()
# apply scaler on selected features
data_scaled[cols_minmax] = mms.fit_transform(data_scaled[cols_minmax])
```

```
data_scaled[cols_minmax].describe().round(2)
```

	AirTime	Month	DayOfMonth	DayOfWeek	DelayInterval
count	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000	192825.0000000000
mean	0.636297872	0.4700000000	0.4700000000	0.4700000000	0.4400000000
std	0.2700000000	0.3200000000	0.2900000000	0.3000000000	0.3900000000
min	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000
25%	0.5500000000	0.1800000000	0.2300000000	0.1700000000	0.0000000000
50%	0.7100000000	0.4500000000	0.5000000000	0.5000000000	0.3300000000
75%	0.8500000000	0.7300000000	0.7300000000	0.8300000000	0.6000000000
max	1.0000000000	1.0000000000	1.0000000000	1.0000000000	1.0000000000

Now these features have all a min value of 0 and max value of 1.

```
data_scaled.head()
```

	ArrDelay	DepDelay	Distance	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	AirTime
143019	-0.636297872	-0.249302490	0.366666667	0.636363634	1.400000000	0.000000000	1.166666667	0.000000000	0.727264271
1886001	0.000000000	-0.317073107	0.313636363	0.181818181	0.000000000	0.000000000	0.166666667	0.000000000	0.591076916
1881989	-0.468051064	-0.1463414634	-0.175575758	-0.181818181	0.000000000	0.000000000	0.000000000	0.000000000	0.503543143
1030789	0.276595747	-0.2682526289	-0.248484848	0.818181818	1.300000000	0.000000000	4.000000000	0.000000000	0.875781577
687997	-0.297872304	-0.3902439024	-0.339393939	0.727272727	0.000000000	0.000000000	0.000000000	0.000000000	0.503126326

Now let's test our models with the scaled data.

## Model 1: Simple Linear Regression

First we will apply scaled features on our Simple Linear Regression model where we are using only 2 variables:

```
# declare dependent variable
y_scaled = data_scaled['ArrDelay']
# independent variables
x_scaled = data_scaled[['Distance', 'TaxiOut']]
x_scaled.shape
```

```
# since the input has only 1 dimension, (1 feature), we reshape first
x_scaled = x_scaled.reshape(-1, 1)
x_scaled.shape
```

(192825, 1)

```
from sklearn.linear_model import LinearRegression
# we'll use a test size of 30%
# use the same random state as before
x_scaled_train_scaled, x_scaled_test_scaled, y_scaled_train_scaled, y_scaled_test_scaled = train_test_split(x_scaled, y_scaled, random_state=10)
```

```
# regression
scaled_regression = LinearRegression()
scaled_regression.fit(x_scaled_train_scaled, y_scaled_train_scaled)
scaled_regression()
```

```
# import modules for the metrics
from sklearn.metrics import r2_score
import sklearn.metrics as metrics
```

```
# functions to show the results
def print_results_test(true, predicted):
    r2_square = metrics.r2_score(true, predicted)
    mse = metrics.mean_squared_error(true, predicted)
    print('R2 for test set : ', r2_square)
    print('MSE (mean squared error) for test set: ', mse)
    print('-----')
```

```
def print_results_train(true, predicted):
    r2_square = metrics.r2_score(true, predicted)
    mse = metrics.mean_squared_error(true, predicted)
    print('R2 for train set : ', r2_square)
    print('MSE (mean squared error) for train set: ', mse)
```

```
# predictions for test sets
scaled_reg_test_pred = scaled_regression.predict(x_scaled_test_scaled)
# predictions for train
scaled_reg_train_pred = scaled_regression.predict(x_scaled_train_scaled)
```

```
print_results_test(y_scaled_test_scaled, scaled_reg_test_pred)
print_results_train(y_scaled_train_scaled, scaled_reg_train_pred)
```

```
R2 for test set : 0.395964548101837
MSE (mean squared error) for test set: 0.1314931985240448
-----
R2 for train set : 0.90984360894703
MSE (mean squared error) for train set: 0.1326930632190840
```

```
# define a function to get the results to a dataframe:
def get_results(true, predicted):
    r2_square = metrics.r2_score(true, predicted)
    mse = metrics.mean_squared_error(true, predicted)
    return r2_square, mse
```

```
# we will save the results in a dataframe
scaled_results = pd.DataFrame(data=[['Simple Linear Regression', get_results(y_scaled_test_scaled, scaled_reg_test_pred)[0], get_results(y_scaled_test_scaled, scaled_reg_test_pred)[1], get_results(y_scaled_train_scaled, scaled_reg_train_pred)[0], get_results(y_scaled_train_scaled, scaled_reg_train_pred)[1]], columns = ['Model', 'R2_test', 'MSE_test', 'R2_train', 'MSE_train'])
scaled_results
```

	Model	R2_test	MSE_test	R2_train	MSE_train
0	Simple Linear Regression	0.9095964548	0.131493199	0.9098436089	0.1326930632

If we compare this with our previous model:

```
# compare with previous model:
# open results df with previous results
results_df = pd.read_csv('results_df.csv')
```

	Model	R2_test	MSE_test	R2_train	MSE_train
0	Simple Linear Regression	0.9095964548	290.4797885540	0.9098436089	293.1189766509

- The results we got in R2 are exactly the same!!! This is logical because both ArrDelay and DepDelay are in the same model and we have scaled them equally, since in this model we only got these two features the result is the same.
- MSE cannot be compared because the scale is different

## Model 2: Ridge Regression

In this model again we will use all the features that were kept after VIF analysis.

```
# dependent variable
y_scaled = data_scaled['ArrDelay']
# independent features
X_scaled = data_scaled.drop(['ArrDelay'], axis = 1)
```

```
print('X_scaled.shape: ', X_scaled.shape)
print('y_scaled.shape: ', y_scaled.shape)
print('y_scaled.shape: ', y_scaled.shape)
```

```
data_scaled.shape: (192825, 628)
X_scaled.shape: (192825, 627)
y_scaled.shape: (192825, 1)
```

```
# use test size of 0.3
```



Out[62]:

	VIF	features
0	4.3372211651	Distance
1	21.317591452	Speed
2	3.2864654499	TaxiOut
3	1.1300343585	CarrierDelay
4	1.0325927511	WeatherDelay
5	1.4300775947	NASDelay
6	1.0011940839	SecurityDelay
7	1.2142223966	LateAircraftDelay
8	8.3143713897	ArrTime
9	3.9409143736	Month
10	4.0530068788	DayOfMonth
11	4.7030819318	DayOfWeek

- Lastly, we will drop 'Speed'.

In [63]:

```
data_final.drop(['Speed'], axis = 1, inplace = True)
data_final.shape
```

Out[63]:

```
(192825, 627)
```

In [64]:

```
cols = ['Distance', 'TaxiOut', 'CarrierDelay', 'WeatherDelay', 'NASDelay',
        'SecurityDelay', 'LateAircraftDelay', 'ArrTime', 'Month', 'DayOfMonth', 'DayOfWeek']
calculate_vif(data_final, cols)
```

Out[64]:

	VIF	features
0	2.5976295963	Distance
1	3.1706768437	TaxiOut
2	1.1263764712	CarrierDelay
3	1.0325659941	WeatherDelay
4	1.3993505497	NASDelay
5	1.0011875757	SecurityDelay
6	1.2062993968	LateAircraftDelay
7	5.9749557045	ArrTime
8	3.6347775313	Month
9	3.7528621858	DayOfMonth
10	4.1945767366	DayOfWeek

Now all of the features have a VIF under 10.

- Random Forest and XGBoost have been the models with better performance.
- I will choose for this exercise XGBoost because its prediction power is as good as Random Forest but it requires much less computational power (approx 5 times faster than Random Forest).

In [65]:

```
data_final.head()
```

Out[65]:

	ArrTime	Distance	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay	ArrTir
143019	21.0000000000	364	21.0000000000	14.0000000000	0.0000000000	7.0000000000	0.0000000000	0.0000000000	1746.00000000
1866001	24.0000000000	813	16.0000000000	0.0000000000	0.0000000000	13.0000000000	0.0000000000	11.0000000000	1419.00000000
1881989	2.0000000000	490	12.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	1209.00000000
1030789	37.0000000000	442	23.0000000000	13.0000000000	0.0000000000	24.0000000000	0.0000000000	0.0000000000	2102.00000000
687997	10.0000000000	382	22.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	1208.00000000

In [66]:

```
y = data_final['ArrDelay']
X = data_final.drop(['ArrDelay'], axis = 1)
```

In [67]:

```
print('data_final.shape: ', data_final.shape)
print('X.shape: ', X.shape)
print('y.shape: ', y.shape)
```

data\_final.shape: (192825, 627)

X.shape: (192825, 626)

y.shape: (192825,)

In [68]:

```
# use test size of 0.3
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
```

In [69]:

```
print('X_train.shape: ', X_train.shape)
print('X_test.shape: ', X_test.shape)
print('y_train.shape: ', y_train.shape)
print('y_test.shape: ', y_test.shape)
```

X\_train.shape: (134977, 626)

X\_test.shape: (57848, 626)

y\_train.shape: (134977,)

y\_test.shape: (57848,)

### XGBoost Regression model

Let's see how well it performs our model without 'DepDelay' feature.

In [70]:

```
import xgboost as xgb

regressor = XGBRegressor()
regressor.fit(X_train, y_train)
```

Out[70]:

XGBRegressor(base\_score=0.5, booster='gbtree', colsample\_bylevel=1, colsample\_bynode=1, colsample\_bynstree=1, enable\_categorical=False, gamma=0, gpu\_id=-1, importance\_type=None, interaction\_constraints='', learning\_rate=0.300000012, max\_delta\_step=0, max\_depth=6, min\_child\_weight=1, missing=nan, monotone\_constraints=(), n\_estimators=100, n\_jobs=8, num\_parallel\_tree=1, predictor='auto', random\_state=0, reg\_alpha=0, reg\_lambda=1, scale\_pos\_weight=1, subsample=1, tree\_method='exact', validate\_parameters=1, verbosity=None)

In [71]:

```
xgb_test_prediction = regressor.predict(X_test)
xgb_train_prediction = regressor.predict(X_train)

print(results_test(y_test, xgb_test_prediction))
print(results_train(y_train, xgb_train_prediction))
```

Out[71]:

R2 for test set : 0.987861778128612  
MSE (mean squared error) for test set: 39.20703110680892  
-----  
R2 for train set : 0.989655755501771  
MSE (mean squared error) for train set: 20.5807972526312

This model shows quite good results in R2, let's compare with our previous results:

In [72]:

```
# compare this results with the results we got the first time we applied XGBRegressor
results_df[3:4]
```

Out[72]:

	Model	R2 test	MSE test	R2 train	MSE train
3	XGBoost Regression	0.9882029254	38.1051112017	0.991335160	28.1077679962

- Deleting 'DepDelay' has caused the R2 to increase a little in train set but decrease in test set
- MSE is lower in train without 'DepDelay' but higher in test set.

Let's draw actual vs predicted values both in test and train set:

In [73]:

```
# df actual vs predicted in test and train
XGB_test_df = pd.DataFrame({'y_test': y_test, 'test_predictions': xgb_test_prediction})
XGB_train_df = pd.DataFrame({'y_train': y_train, 'train_predictions': xgb_train_prediction})
```

In [74]:

```
# plot actual vs predicted
plt.figure(figsize=(18,9))
plt.subplot(1,2,1)
plt.scatter(XGB_test_df['y_test'], y = XGB_test_df['test_predictions'], alpha = 0.7)
plt.plot([0, 1750], [0, 1750], color = 'green', linewidth = 1.5, linestyle = 'dashed')
plt.xlabel('Actual values', fontsize=13)
plt.ylabel('Predicted values', fontsize = 13)
plt.title('Actual vs Predicted values in test set', fontsize = 16)
```

plt.subplot(1,2,2)

plt.scatter(x=XGB\_train\_df['y\_train'], y = XGB\_train\_df['train\_predictions'], alpha = 0.7, c = 'magenta')

plt.plot([0, 1750], [0, 1750], color = 'green', linewidth = 1.5, linestyle = 'dashed')

plt.xlabel('Actual values', fontsize = 13)

plt.ylabel('Predicted values', fontsize = 13)

plt.title('Actual vs Predicted values in train set', fontsize = 16)

plt.tight\_layout()

Actual vs Predicted in XGboost Regression(without DepDelay) in test and train sets

Actual vs Predicted values in test set

Actual vs Predicted values in train set

- We observe here the dispersion in test set is much more noticeable compared with the results we got in our models where we used 'DepDelay'.
- The errors are more visible compared with the train set too: the plot is much thinner around the 45° degree line in the train set than in test set.

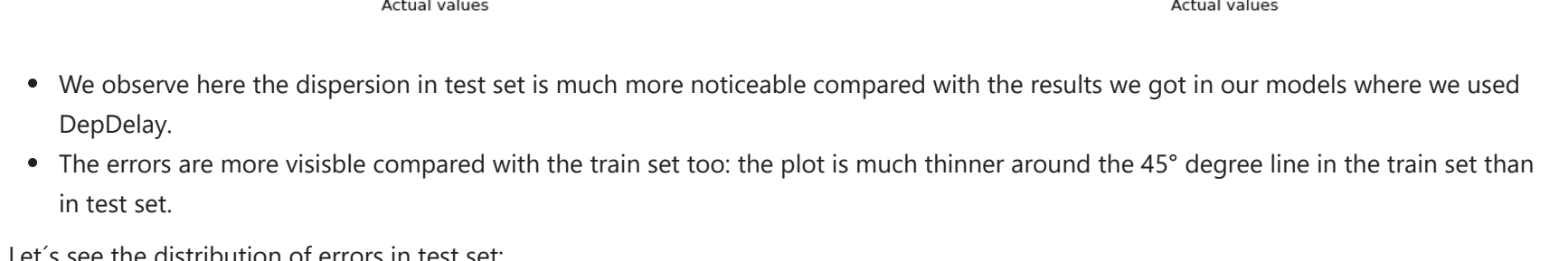
Let's see the distribution of errors in test set:

In [75]:

```
# create df of errors in test set
XGBoost_errors_df = pd.DataFrame({'Error Values': (y_test - xgb_test_prediction)})
```

In [76]:

```
# now we plot it
plt.figure(figsize=(10, 7))
plt.title('Distribution of errors in test set with XGBoost Regression', fontsize = 16)
sns.kdeplot(data = XGBoost_errors_df, fill = True)
```



In [77]:

```
XGBoost_errors_df.describe().round(2)
```

Out[77]:

	Error Values
count	57848.0000000000
mean	0.0600000000
std	6.2600000000
min	-179.5500000000
25%	-1.5100000000
50%	0.0200000000
75%	1.8000000000
max	708.8200000000

- The max value of error here is the double of the model when we used DepDelay
- On the other hand, errors seem concentrated around values around zero.

### Cross Validation

Since this model shows more differences between test and train than the model we used when we used 'depDelay' we will do cross validation now to compare the results.

In [78]:

```
from sklearn.model_selection import cross_val_score
# use cv=5, same as before
xgb_reg = XGBRegressor()
xgb_r2_scores = cross_val_score(xgb_reg, X, y, cv = 5)
xgb_mse_scores
```

Out[78]:

```
array([0.97235609, 0.99097991, 0.98266899, 0.98815041, 0.99189471])
```

In [79]:

```
# compute the mean of cv R2
r2_cv_xgb = np.mean(xgb_r2_scores)
r2_cv_xgb
```

Out[79]:

```
0.9852100210588507
```

In [81]:

```
xgb_mse_scores = cross_val_score(xgb_reg, X, y, cv = 5, scoring = 'neg_mean_squared_error')
xgb_mse_scores
```

Out[81]:

```
array([-108.3043837, -28.13787229, -64.6892613, -31.22398048,
       -22.1609237])
```

In [82]:

```
# calculate the mean of mae in cv
mse_cv_xgb = abs(np.mean(xgb_mse_scores))
mse_cv_xgb
```

Out[82]:

```
50.90328429413569
```

Now, let's compare these results with our previous Cross Validation results:

In [83]:

```
print('R2 cv without DepDelay: ', r2_cv_xgb)
print('MSE cv without DepDelay: ', mse_cv_xgb)
# open the dataframe with cross validation results :
cv_results = pd.read_csv('cv_results')
# add results from model without 'DepDelay' to our dataframe of cross validation results
cv_no_depdelay = pd.DataFrame(data= [['XGBoost Regression without DepDelay', r2_cv_xgb, mse_cv_xgb],
                                     'DayOfWeek', 'R2', 'MSE'])
cv_results = pd.concat([cv_results, cv_no_depdelay])
cv_results[3:5]
```

Out[83]:

R2 cv without DepDelay: 0.9852100210588507  
MSE cv without DepDelay: 50.90328429413569

	Model	R2	MSE
3	XGBoost Regression	0.9874545799	40.1445838049
0	XGBoost Regression without DepDelay	0.9852100211	50.9032842941

- We observe that our R2 now is lower (0.985) and MSE is higher (50.90) than before.
- The model loses predictive power when we delete DepDelay but it still has high predictive power.
- This may be explained from the fact that we have a lot of features that are highly correlated with 'ArrDelay'.

To show which features are important to predict 'ArrDelay' we can do a feature importance analysis, that shows the importance of each feature. We will order these features from more to less important and select the 5 most important features (we will exclude dummy columns):

### Feature importance

In [89]:

```
X.columns[1:11] # select the columns used in last model (we have dropped DepDelay), we will exclude dummies
```

Out[89]:

```
Index(['Distance', 'TaxiOut', 'CarrierDelay', 'WeatherDelay', 'NASDelay',
       'SecurityDelay', 'LateAircraftDelay', 'ArrTime', 'Month', 'DayOfMonth',
       'DayOfWeek'],
      dtype='object')
```

In [90]:

```
# let's do a feature importance analysis
regressor.feature_importances_[1:11]
```

Out[90]:

```
array([0.0108322, 0.0068382, 0.2974295, 0.14345698, 0.20989503,
       0.01257295, 0.27630112, 0.00107032, 0.00073453, 0.00086872,
       0.00303355], dtype=float32)
```

In [91]:

```
from sklearn.inspection import permutation_importance
# let's do a plot with the 5 most important features of our model:
perm_importance = permutation_importance(regressor, X_test, y_test)
```

Out[91]:

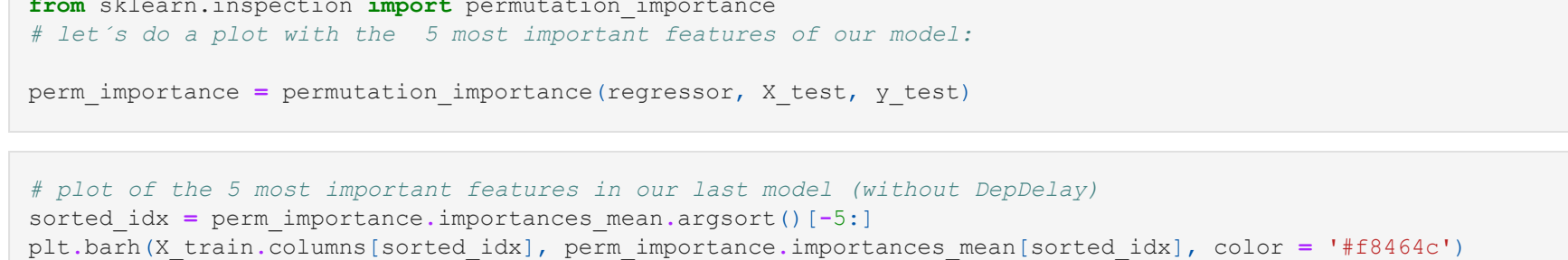
R2 cv without DepDelay: 0.9852100210588507  
MSE cv without DepDelay: 50.90328429413569

	Model	R2	MSE
3	XGBoost Regression	0.9874545799	40.1445838049
0	XGBoost Regression without DepDelay	0.9852100211	50.9032842941

- We observe that our R2 now is lower (0.985) and MSE is higher (50.90) than before.
- The model loses predictive power when we delete DepDelay but it still has high predictive power.
- This may be explained from the fact that we have a lot of features that are highly correlated with 'ArrDelay'.

To show which features are important to predict 'ArrDelay' we can do a feature importance analysis, that shows the importance of each feature. We will order these features from more to less important and select the 5 most important features (we will exclude dummy columns):

### Top 10 Variables Importance



- Out of all the features (if we exclude 'DepDelay'), the most important ones in our last model are : LateAircraftDelay, CarrierDelay, WeatherDelay and WeatherDelay. This makes sense because in the cases where ArrDelay is bigger than 15, the delay is disaggregated between these features (these features are 0 in the cases where ArrDelay is less than 15 minutes).
- Actually, DepDelay is the sum of these features in the case where ArrDelay is less than 15.

What would happen if we do our model using only these 4 features?

We will do now our model using only these 4 features and see how it performs.

In [93]:

```
# choose dependent variable
y4 = data_sample['ArrDelay']
# we choose our 4 independent features
X4 = data_sample[['WeatherDelay', 'NASDelay', 'LateAircraftDelay', 'CarrierDelay']]
```

In [95]:

```
print('X4.shape: ', X4.shape)
print('y4.shape: ', y4.shape)
```

X4.shape: (192825, 4)

y4.shape: (192825,)

In [96]:

```
# use test size of 0.3 and same random state
X4_train, X4_test, y4_train, y4_test = train_test_split(X4, y4, test_size=0.3, random_state=0)
```

In [97]:

```
xgb_4regressor = XGBRegressor()
xgb_4regressor.fit(X4_train, y4_train)
```

Out[97]:

XGBRegressor(base\_score=0.5, booster='gbtree', colsample\_bylevel=1, colsample\_bynode=1, colsample\_bynstree=1, enable\_categorical=False, gamma=0, gpu\_id=-1, importance\_type=None, interaction\_constraints='', learning\_rate=0.300000012, max\_delta\_step=0, max\_depth=6, min\_child\_weight=1, missing=nan, monotone\_constraints=(), n\_estimators=100, n\_jobs=8, num\_parallel\_tree=1, predictor='auto', random\_state=0, reg\_alpha=0, reg\_lambda=1, scale\_pos\_weight=1, subsample=1, tree\_method='exact', validate\_parameters=1, verbosity=None)

In [98]:

```
xgb4_test_prediction = xgb_4regressor.predict(X4_test)
xgb4_train_prediction = xgb_4regressor.predict(X4_train)

print(results_test(y4_test, xgb4_test_prediction))
print(results_train(y4_train, xgb4_train_prediction))
```

R2 for test set : 0.9876581318162035  
MSE (mean squared error) for test set: 39.86482036735483  
-----  
R2 for train set : 0.9916961446486076  
MSE (mean squared error) for train set: 26.37788102136856

Seeing these results, we are able to predict 'ArrDelay' quite well using only 4 features!!

### CONCLUSIONS

- We have performed models using all features available and compared the results with models using only some of the features
- 'DepDelay' alone can explain 90 % of the variability of the dependent variable 'ArrDelay'.
- 'DepDelay' is the sum of 'WeatherDelay', 'NASDelay', 'LateAircraftDelay', 'CarrierDelay' when ArrDelay > 15, so using a model with only these 4 features can return pretty good results.
- Removing 'DepDelay' does not decrease predictive power considerably because we have other features that measure delay too. In our exploratory exercise we saw that average delay varies depending on other features like Month or Carriers, it could also be interesting to build a model studying these other causes in a deeper way and exclude all the features that refer to delays (idea for another exercise).
- Although we have some outliers, the outliers in the important predictors are also outliers in the dependent variable, so the results are not so much affected by outliers (if some of the causes of delay has a big outlier, 'ArrDelay' will be an outlier too in that observation)
- Scaling the features doesn't change our model so much since the most important predictors are in the same scale as our dependent variable (scaling improves the best performing models but very little). Besides R2 without scaling is already very high, so there's actually little room for improvement.
- Regarding to the models we used, the ones with better results have been Random Forest and XGBoost Regressor, but XGBoost is much more efficient with respect to the computational time (5 times faster than Random Forest).

In [ ]: