

Exercici 1 (Nivell 1)

Crea almenys tres models de regressió diferents per intentar predir el millor possible l'endarreriment dels vols (ArrDelay) de DeltaFlights.csv.

```
import numpy as np #library for scientific computing
import pandas as pd #library for dataframes
import matplotlib.pyplot as plt #graphic library
import seaborn as sns #advanced graphic library based on matplotlib
import warnings #library to manage warnings
import scipy #library with algorithms for statistics and scientific computing
from sklearn.linear_model import LinearRegression #regression algorithm
from sklearn.ensemble import RandomForestRegressor #regression algorithm
from sklearn.tree import DecisionTreeRegressor #regression algorithm
from sklearn.model_selection import GridSearchCV #optimization parameter algorithm
from sklearn.model_selection import cross_val_score #cross validation algorithm
from sklearn.metrics import mean_squared_error #MSE calculation
from sklearn.metrics import mean_absolute_error #MAE calculation
from sklearn.metrics import r2_score #coefficient of determination calculation
from sklearn.preprocessing import MinMaxScaler

warnings.filterwarnings('ignore')
```

```
#Importing dataset

df_flight_imported = pd.read_csv('DelayedFlights.csv', index_col=0)

#we take a sample to speed up the algorithm
#this may affect to results
df_flight = df_flight_imported.sample(100000)
```

```
#dataset information

print(df_flight.info(nul_counts=True))
df_flight.describe().transpose()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 100000 entries, 3682056 to 3599538
Data columns (total 29 columns):
 #   Column                Non-Null Count  Dtype
---  ---
0   Year                  100000 non-null    int64
1   Month                 100000 non-null    int64
2   DayOfMonth            100000 non-null    int64
3   DayOfWeek             100000 non-null    int64
4   DepTime               100000 non-null    float64
5   CRSDepTime            100000 non-null    int64
6   ArrTime               99653 non-null     float64
7   CRSArrTime            100000 non-null    int64
8   UniqueCarrier         100000 non-null    object
9   FlightNum             100000 non-null    int64
10  TailNum               100000 non-null    object
11  ActualElapsedTime     99578 non-null     float64
12  CRSAirline            99578 non-null     float64
13  AirTime               99578 non-null     float64
14  ArrDelay              100000 non-null    float64
15  DepDelay              100000 non-null    float64
16  Origin                99578 non-null     object
17  Dest                 100000 non-null    object
18  Distance              100000 non-null    int64
19  TaxiIn                99653 non-null     float64
20  TaxiOut               99580 non-null     float64
21  Cancelled             100000 non-null    int64
22  CancellationCode      100000 non-null    object
23  Diverted              100000 non-null    int64
24  CarrierDelay          64506 non-null     float64
25  WeatherDelay          64506 non-null     float64
26  NASDelay              64506 non-null     float64
27  SecurityDelay         64506 non-null     float64
28  LateAircraftDelay     64506 non-null     float64
dtypes: float64(14), int64(14), object(1)
memory usage: 22.4+ MB
None
```

	count	mean	std	min	25%	50%	75%	max
Year	100000.0	2008.000000	0.000000	2008.0	2008.0	2008.0	2008.0	2008.0
Month	100000.0	6.120900	3.485247	1.0	3.0	6.0	9.00	12.0
DayOfMonth	100000.0	15.742140	8.766615	1.0	8.0	16.0	23.00	31.0
DayOfWeek	100000.0	3.987400	1.994558	1.0	2.0	4.0	6.00	7.0
DepTime	100000.0	1517.982610	451.885171	1.0	1203.0	1545.0	1901.00	2400.0
CRSDepTime	100000.0	1466.828250	425.994317	1.0	1135.0	1510.0	1820.00	2350.0
ArrTime	99653.0	1608.228061	549.856272	1.0	1314.0	1714.0	2031.00	2400.0
CRSArrTime	100000.0	1632.634120	466.525299	0.0	1324.0	1705.0	2015.00	2400.0
FlightNum	100000.0	2187.301550	1946.763453	1.0	607.0	1540.0	3427.00	7829.0
TailNum	99578.0	133.025196	71.704471	17.0	80.0	116.0	164.75	666.0
ActualElapsedTime	99578.0	134.074345	71.156087	15.0	81.0	116.0	165.00	660.0
AirTime	99578.0	108.049398	68.383146	0.0	58.0	90.0	137.00	636.0
DepDelay	99578.0	42.309104	57.722580	-61.0	9.0	24.0	56.00	2461.0
CRSDepDelay	100000.0	43.319960	54.378297	6.0	12.0	24.0	53.00	2457.0
Distance	100000.0	764.244010	573.197992	31.0	338.0	606.0	997.00	4962.0
TaxiIn	99653.0	6.801270	5.174862	0.0	4.0	6.0	8.00	157.0
TaxiOut	99580.0	18.193189	14.209370	0.0	10.0	14.0	21.00	78.0
Cancelled	100000.0	0.000030	0.018163	0.0	0.0	0.0	0.00	1.0
Diverted	100000.0	0.003890	0.062249	0.0	0.0	0.0	0.00	1.0
CarrierDelay	64506.0	19.406133	44.592061	0.0	0.0	2.0	21.00	1552.0
WeatherDelay	64506.0	3.703392	21.968289	0.0	0.0	0.0	0.00	1090.0
NASDelay	64506.0	14.792763	33.251891	0.0	0.0	1.0	14.00	768.0
SecurityDelay	64506.0	0.091557	1.854390	0.0	0.0	0.0	0.00	131.0
LateAircraftDelay	64506.0	25.382166	42.487612	0.0	0.0	8.0	34.00	1303.0

Dataset Explanation

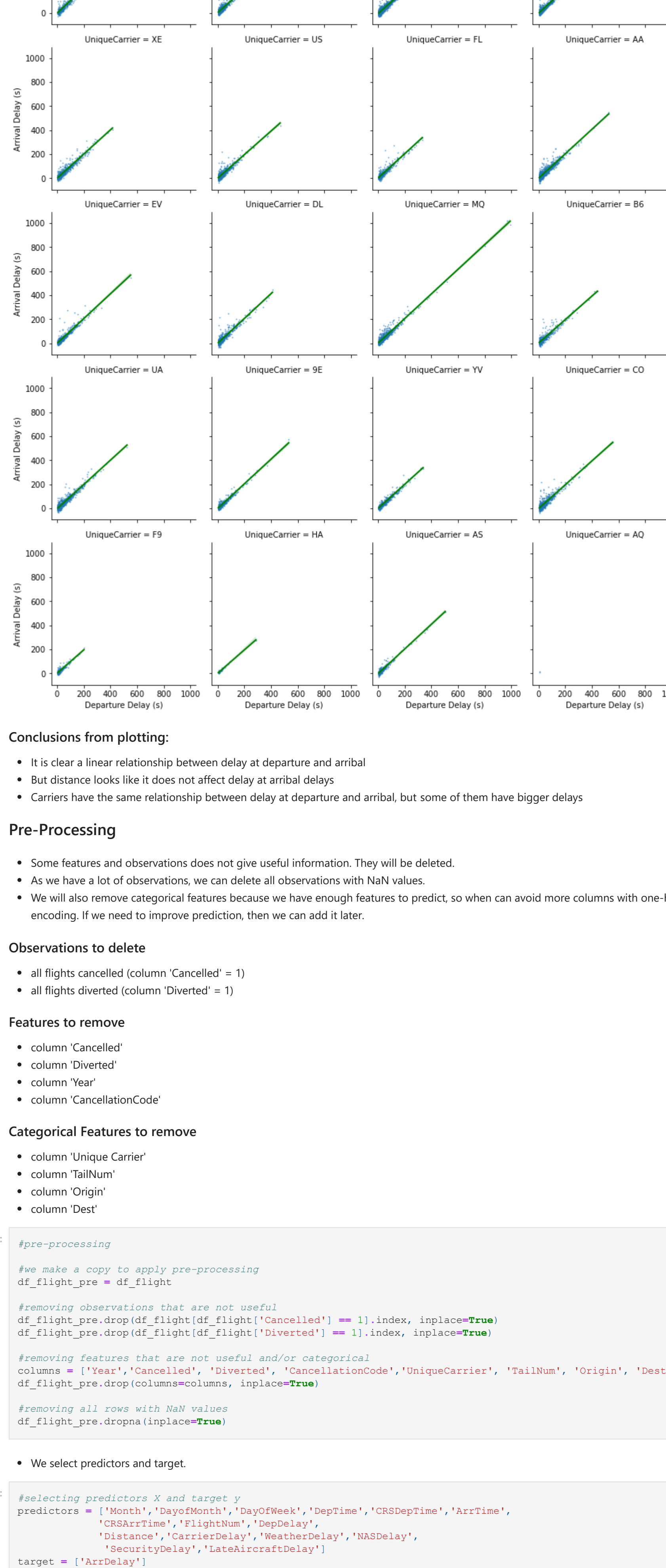
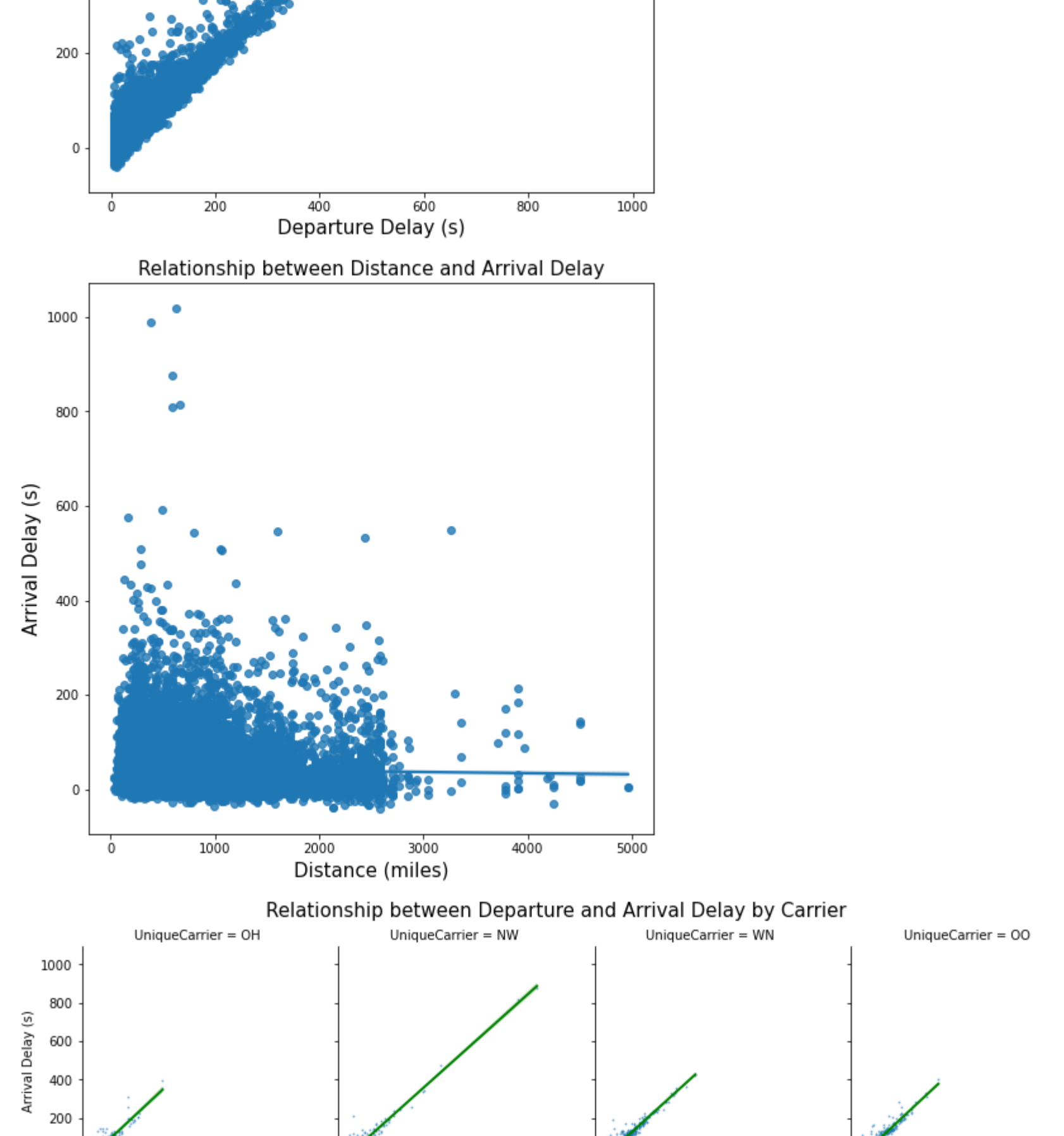
1. Year: 2008
2. Month: 1-12
3. DayOfMonth: 1-31
4. DayOfWeek: 1 (Monday) - 7 (Sunday)
5. DepTime: actual departure time (local, hh:mm)
6. CRSDepTime: scheduled departure time (local, hh:mm)
7. ArrTime: actual arrival time (local, hh:mm)
8. CRSArrTime: scheduled arrival time (local, hh:mm)
9. UniqueCarrier: unique carrier code
10. FlightNum: flight number
11. TailNum: plane tail number; aircraft registration, unique aircraft identifier
12. ActualElapsedTime: in minutes
13. CRSElapsedTime: in minutes
14. AirTime: in minutes

Dataset Explanation

- Year: 2008
- Month: 1-12
- DayOfMonth: 1-31
- DayOfWeek: 1 (Monday) - 7 (Sunday)
- DepTime: actual departure time (local, hhmm)
- CRSDepTime: scheduled departure time (local, hhmm)
- ArrTime: actual arrival time (local, hhmm)
- CRSArrTime: scheduled arrival time (local, hhmm)
- UniqueCarrier: unique carrier code
- FlightNum: flight number
- TailNum: plane tail number; aircraft registration, unique aircraft identifier
- ActualElapsedTime: in minutes
- CRSAirline: in minutes
- AirTime: in minutes
- ArrDelay: arrival delay, in minutes: A flight is counted as "on time" if it operated less than 15 minutes later the scheduled time shown in the carriers' Computerized Reservations Systems (CRS).
- DepDelay: departure delay, in minutes
- Origin: origin IATA airport code
- Dest: destination IATA airport code
- Distance: in miles
- TaxiIn: taxi in time (the movement of an aircraft on the ground after landing), in minutes
- TaxiOut: taxi out time (the movement of an aircraft on the ground before taking off), in minutes
- Cancelled: 1 = yes, 0 = no
- CancellationCode: reason for cancellation (A = carrier, B = weather, C = NAS, D = security)
- Diverted: 1 = yes, 0 = no
- CarrierDelay: in minutes: Carrier delay is within the control of the air carrier. Examples of occurrences that may determine carrier delay are: aircraft cleaning, aircraft damage, awaiting the arrival of connecting passengers or crew, baggage, bird strike, cargo loading, catering, computer, out-of-control equipment, crew legality (pilot or attendant rest), damage by hazardous goods, engineering inspection, fueling, handling disabled passengers, late crew, lavatory servicing, maintenance, oversales, potable water servicing, removal of unruly passenger, slow boarding or seating, stowing carry-on baggage, weight and balance delays.
- WeatherDelay in minutes: Weather delay is caused by one or more hazardous weather conditions that are forecasted or manifest themselves on point of departure, enroute, or on point of arrival.
- NASDelay in minutes: Delay that is within the control of the National Airspace System (NAS) may include: non-extreme weather conditions, airport operations, heavy traffic volume, air traffic control, etc.
- SecurityDelay in minutes: Security delay is caused by evacuation of a terminal or concourse, re-boarding of aircraft because of security breach, inoperative screening equipment and/or long lines in excess of 29 minutes at screening areas.
- LateAircraftDelay in minutes: Arrival delay at an airport due to the late arrival of the same aircraft at a previous airport. The ripple effect of an earlier delay at downstream airports is referred to as delay propagation.

Conclusions from dataset:

- some columns have many NaN values
- if object types (5) are going to be used as the predictors, they should be converted from categorical to numerical



Conclusions from plotting:

- It is clear a linear relationship between delay at departure and arrival
- But distance looks like it does not affect delay at arrival delays
- Carriers have the same relationship between delay at departure and arrival, but some of them have bigger delays

Pre-Processing

- Some features and observations does not give useful information. They will be deleted.
- As we have a lot of observations, we can delete all observations with NaN values.
- We will also remove categorical features because we can add enough features to predict, so when can avoid more columns with one-hot encoding. If we need to improve prediction, then we can add it later.

Observations to delete

- all flights cancelled (column 'Cancelled' = 1)
- all flights diverted (column 'Diverted' = 1)

Features to remove

- column 'Cancelled'
- column 'Diverted'
- column 'Year'
- column 'CancellationCode'

Categorical Features to remove

- column 'UniqueCarrier'
- column 'TailNum'
- column 'Origin'
- column 'Dest'

In [7]:	<pre>#pre-processing #we make a copy to apply pre-processing df_flight_pre = df_flight #removing observations that are not useful df_flight_pre.drop(df_flight[df_flight['Cancelled'] == 1].index, inplace=True) df_flight_pre.drop(df_flight[df_flight['Diverted'] == 1].index, inplace=True) #removing features that are not useful and/or categorical columns = ['Year','Cancelled', 'Diverted', 'CancellationCode', 'UniqueCarrier', 'TailNum', 'Origin', 'Dest'] df_flight_pre.drop(columns, inplace=True) #removing all rows with NaN values df_flight_pre.dropna(inplace=True)</pre>
	<ul style="list-style-type: none">• We select predictors and target.
In [8]:	<pre>#selecting predictors X and target y predictors = ['Month','DayOfMonth','DayOfWeek','DepTime','CRSDepTime','ArrTime', 'Distance','FlightNum','DepDelay', 'CarrierDelay','WeatherDelay','NASDelay', 'SecurityDelay','LateAircraftDelay'] target = ['ArrDelay'] X = df_flight_pre[predictors] y = df_flight_pre[target] #splitting in train and test dataset X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)</pre>
In [9]:	<pre>#regression models model2 = 'SLR/MLR (Single/Multiple linear Regression)' model1 = 'Decision Tree Regression' model3 = 'Random Forest Regression' model1 = LinearRegression().fit(X_train, y_train) model2 = DecisionTreeRegressor(random_state=0).fit(X_train, y_train) model3 = RandomForestRegressor(random_state=0).fit(X_train, y_train)</pre>

Evaluation Metrics:

- R2: It is the coefficient of determination. It measures the amount of variance in the predictions. Best score is 1.0 or -1.0. Worst score is 0.0.
- MSE (Mean Square Error): It is the average of the square of the errors, where error is the difference between an observed value and its predicted value. There is no correct value for MSE, the lower the value the better. 0 means the model is perfect.
- MAE (Mean Absolute Error): It is the average of all absolute errors. The absolute error is the absolute value of the difference between each actual value and its predicted value. The lower the value the better.

In [10]:	<pre>#model 1 #we obtain R2, intercept and slope from train data R2_train1 = model1.score(X_train, y_train) y_pred1ac = model1.predict(X_test) intercept_b0 = model1.intercept_ slope_b1 = model1.coef_ #we obtain predicted values from test data y_pred1 = model1.predict(X_test) #Evaluation Metrics #MSE MSE_1 = round(mean_squared_error(y_test, y_pred1),3) #MAE MAE_1 = round(mean_absolute_error(y_test, y_pred1),3) #R2 from predicted data r2_test1 = round(r2_score(y_test, y_pred1),3)</pre>
In [11]:	<pre>#model 2 #we obtain R2 from train data R2_train2 = model2.score(X_train, y_train) #we obtain predicted values from test data y_pred2 = model2.predict(X_test) #Evaluation Metrics #MSE MSE_2 = round(mean_squared_error(y_test, y_pred2),3) #MAE MAE_2 = round(mean_absolute_error(y_test, y_pred2),3) #R2 from predicted data r2_test2 = round(r2_score(y_test, y_pred2),3)</pre>
In [12]:	<pre>#model 3 #we obtain R2 from train data R2_train3 = model3.score(X_train, y_train) #we obtain predicted values from test data y_pred3 = model3.predict(X_test) #Evaluation Metrics #MSE MSE_3 = round(mean_squared_error(y_test, y_pred3),3) #MAE MAE_3 = round(mean_absolute_error(y_test, y_pred3),3) #R2 from predicted data r2_test3 = round(r2_score(y_test, y_pred3),3)</pre>

Exercici 2 (Nivell 1)

Compara's en base al MSE i al R2.

Comparative Table

```
#parameters for Model 2
print(model2)
print('\ndefault Parameters:')
print(model_2.get_params())

#training model 2

#estimator
model_2_est = DecisionTreeRegressor(random_state=0)

#selection of parameters and creating a dict for grid tool
criterion = ["squared_error", "friedman_mse", "absolute_error", "poisson"]
splitter = ["best", "random"]
max_features = ["sqrt", "log2"]

param_grid_2 = dict(criterion=criterion, splitter=splitter, max_features=max_features)

#launching grid tool
model_2_grid = GridSearchCV(estimator=model_2_est, param_grid=param_grid_2, n_jobs=-1)
model_2_result = model_2_grid.fit(X_train, y_train)
```

Conclusions:

- The Multiple linear Regression algorithm is the most simple model but it fits exactly. The ratio between computing work and results is excellent
- The other two algorithms work well, but they are more complex, they need more computing work and they does not fit as well as Linear
- Tuning parameters, applying feature engineering or taking more predictors can improve results for the Decision Tree and Random Forest algorithm.

Exercici 3 (Nivell 1)

Entrena's utilitzant els diferents paràmetres que admeten.

In [14]:	<pre>#parameters for Model 1 print(model1) print('\nDefault Parameters:') print(model1.get_params()) #training model 1 #there are no parameters to fit this model, because is quite simple SLR/MLR (Single/Multiple linear Regression) Default Parameters: {'copy_X': True, 'fit_intercept': True, 'n_jobs': None, 'normalize': False, 'positive': False}</pre>
In [15]:	<pre>#parameters for Model 2 print(model2) print('\nDefault Parameters:') print(model2.get_params()) #training model 2 #estimator model2_est = DecisionTreeRegressor(random_state=0) #selection of parameters and creating a dict for grid tool criterion = ['squared_error', 'friedman_mse', 'absolute_error', 'poisson'] splitter = ['best', 'random'] max_features = ['sqrt', 'log2'] param_grid2 = dict(criterion=criterion, splitter=splitter, max_features=max_features) #launching grid tool model2_grid = GridSearchCV(estimator=model2_est, param_grid=param_grid2, n_jobs=-1) model2_result = model2_grid.fit(X_train, y_train) #results print('\n') print('Best Score: ', model2_result.best_score_) print('Best Params: ', model2_result.best_params_) Decision Tree Regression Default Parameters: {'alpha': 0.0, 'criterion': 'mse', 'max_depth': None, 'max_features': None, 'max_leaf_nodes': None, 'min_impurity_decrease': 0.0, 'min_impurity_split': 0.0, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'random_state': 0, 'splitter': 'best'}</pre>
In [16]:	<pre>#parameters for Model 3 print(model3) print('\nDefault Parameters:') print(model3.get_params()) #training model 3 #estimator model3_est = RandomForestRegressor(random_state=0) #selection of parameters and creating a dict for grid tool n_estimators = ['squared_error', 'absolute_error', 'poisson'] n_estimators = [10, 20, 30] max_features = ['auto', 'None'] param_grid3 = dict(criterion=criterion, n_estimators=n_estimators, max_features=max_features) #launching grid tool model3_grid = GridSearchCV(estimator=model3_est, param_grid=param_grid3, n_jobs=-1) model3_result = model3_grid.fit(X_train, y_train) #results print('\n') print('Best Score: ', model3_result.best_score_) print('Best Params: ', model3_result.best_params_) Random Forest Regression Default Parameters: {'bootstrap': True, 'copy_alpha': 0.0, 'criterion': 'mse', 'max_depth': None, 'max_features': 'auto', 'max_leaf_nodes': None, 'max_samples': None, 'min_impurity_decrease': 0.0, 'min_impurity_split': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 100, 'n_jobs': None, 'oob_score': 1, 'random_state': 0, 'verbose': 0, 'warm_start': False}</pre>

Comparative Table

Best Score: 0.9775977396198756
Best Params: {'criterion': 'poisson', 'max_features': 'auto', 'n_estimators': 20}

Comparative Table

```
#comparing R2 standard and R2 with best parameters

df_EvaluationMetrics_param = pd.DataFrame(data=[(r2_test2, model2_result.best_score_), (r2_test3, model3_result.best_score_)],
                                           columns=['R2 Default', 'R2 Hypertuned'],
                                           index=[model2, model3])

df_EvaluationMetrics_param
```

	R2 Default	R2 Hypertuned
Decision Tree Regression	0.988	0.906956
Random Forest Regression	0.993	0.977598

Conclusions:

Conclusions:

- The parameters that we chose to train the algorithm with the GridSearch does not improve default parameters for R2.
- Maybe we can try with other parameters, but it takes a lot of computing power, specially the Random Forest algorithm.

Exercici 4 (Nivell 1)

Compara el seu rendiment utilitzant l'aproximació traint/test o utilitzant totes les dades (validació interna).

```
model1_cross = np.round(cross_val_score(model_1_est, X, y),3)
model2_cross = np.round(cross_val_score(model_2_est, X, y),3)
model3_cross = np.round(cross_val_score(model_3_est, X, y),3)
```

Comparative Table

```
#comparison R2 standard and R2 with cross-validation

df_EvaluationMetrics_cv = pd.DataFrame(data=[['r2_test1', model1_cross], ['r2_test2', model2_cross], ['r2_test3', model3_cross]],
                                     columns=['R2 Train-Test', 'R2 Cross-Validation'], index=[model1, model2, model3])
```

	R2 Train-Test	R2 Cross-Validation
SLR/MLR (Single/Multiple linear Regression)	1.000	[1.0, 1.0, 1.0, 1.0, 1.0]
Decision Tree Regression	0.984	[0.986, 0.989, 0.988, 0.985, 0.986]
Random Forest Regression	0.992	[0.991, 0.994, 0.993, 0.992, 0.994]

Conclusions:

- Applying feature engineering improves the algorithm performance

Exercici 5 (Nivell 2)

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

We will try to apply a feature engineering tool to improve predictions.

As we have a lot of features with different measurement units, we are going to use the MinMaxScaler to unify all features units to the scale from 0 to 1.

In [19]:	<pre>model1ac = LinearRegression().fit(X_train, y_train) model2ac = DecisionTreeRegressor(random_state=0).fit(X_train, y_train) model3ac = RandomForestRegressor(random_state=0).fit(X_train, y_train) #launching predictors X and target y predictors = ['Month','DayOfMonth','DayOfWeek','DepTime','CRSDepTime','ArrTime', 'Distance','FlightNum','DepDelay', 'CarrierDelay','WeatherDelay','NASDelay', 'SecurityDelay','LateAircraftDelay'] target = ['ArrDelay'] X = df_flight_pre[predictors] y = df_flight_pre[target] #splitting in train and test dataset X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)</pre>
In [20]:	<pre>#we obtain predicted values from test data y_pred1ac = model1ac.predict(X_test) y_pred2ac = model2ac.predict(X_test) y_pred3ac = model3ac.predict(X_test) #MSE MSE_ac1 = round(mean_squared_error(y_test, y_pred1ac),3) MSE_ac2 = round(mean_squared_error(y_test, y_pred2ac),3) MSE_ac3 = round(mean_squared_error(y_test, y_pred3ac),3) #R2 from predicted data r2_test1ac = round(r2_score(y_test, y_pred1ac),3) r2_test2ac = round(r2_score(y_test, y_pred2ac),3) r2_test3ac = round(r2_score(y_test, y_pred3ac),3)</pre>
In [21]:	<pre>#parameters for Model 1 print(model1) print('\nDefault Parameters:') print(model1.get_params()) #training model 1 #there are no parameters to fit this model, because is quite simple SLR/MLR (Single/Multiple linear Regression) Default Parameters: {'copy_X': True, 'fit_intercept': True, 'n_jobs': None, 'normalize': False, 'positive': False}</pre>
In [22]:	<pre>#parameters for Model 2 print(model2) print('\nDefault Parameters:') print(model2.get_params()) #training model 2 #estimator model2_est = DecisionTreeRegressor(random_state=0) #selection of parameters and creating a dict for grid tool criterion = ['squared_error', 'friedman_mse', 'absolute_error', 'poisson'] splitter = ['best', 'random'] max_features = ['sqrt', 'log2'] param_grid2 = dict(criterion=criterion, splitter=splitter, max_features=max_features) #launching grid tool model2_grid = GridSearchCV(estimator=model2_est, param_grid=param_grid2, n_jobs=-1) model2_result = model2_grid.fit(X_train, y_train) #results print('\n') print('Best Score: ', model2_result.best_score_) print('Best Params: ', model2_result.best_params_) Decision Tree Regression Default Parameters: {'alpha': 0.0, 'criterion': 'mse', 'max_depth': None, 'max_features': None, 'max_leaf_nodes': None, 'min_impurity_decrease': 0.0, 'min_impurity_split': 0.0, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'random_state': 0, 'splitter': 'best'}</pre>
In [23]:	<pre>#parameters for Model 3 print(model3) print('\nDefault Parameters:') print(model3.get_params()) #training model 3 #estimator model3_est = RandomForestRegressor(random_state=0) #selection of parameters and creating a dict for grid tool n_estimators = ['squared_error', 'absolute_error', 'poisson'] n_estimators = [10, 20, 30] max_features = ['auto', 'None'] param_grid3 = dict(criterion=criterion, n_estimators=n_estimators, max_features=max_features) #launching grid tool model3_grid = GridSearchCV(estimator=model3_est, param_grid=param_grid3, n_jobs=-1) model3_result = model3_grid.fit(X_train, y_train) #results print('\n') print('Best Score: ', model3_result.best_score_) print('Best Params: ', model3_result.best_params_) Random Forest Regression Default Parameters: {'bootstrap': True, 'copy_alpha': 0.0, 'criterion': 'mse', 'max_depth': None, 'max_features': 'auto', 'max_leaf_nodes': None, 'max_samples': None, 'min_impurity_decrease': 0.0, 'min_impurity_split': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 100, 'n_jobs': None, 'oob_score': 1, 'random_state': 0, 'verbose': 0, 'warm_start': False}</pre>

Comparative Table

```

#selecting predictors X and target y
predictors = ['Month','DayOfMonth','DayOfWeek','DepTime','CRSDepTime','ArrTime',
              'CRSArrTime','FlightNum',
              'Distance','CarrierDelay','WeatherDelay','NASDelay',
              'SecurityDelay','LateAircraftDelay']

target = ['ArrDelay']

X = df_flight_pre[predictors]
y = df_flight_pre[target]

#splitting in train and test dataset
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.33, random_state=1)

model_LND = LinearRegression().fit(X_train, y_train)
model_ND = DecisionTreeRegressor(random_state=0).fit(X_train, y_train)
model_ZND = RandomForestRegressor(random_state=0).fit(X_train, y_train)

```

Conclusions:

- The Linear Regression algorithm is not affected when we remove this feature.
- Decision tree and Random Forest are affected when this feature is not available, and results get worse.
- To compensate this, we could add other features that we did not use, like categorical features or try with other parameters or algorithms.