

S09 T02: Apprententage Supervisat - Regressions - part 1

```
import pandas as pd
import matplotlib.pyplot as plt

# create a column to categorize ArrDelay, explained in previous exercise
data_sample['ArrDelay'] = data_sample['ActualElapsedTime'] - data_sample['CRSElapsedTime']
data_sample['ArrDelay'] = data_sample['ArrDelay'].astype('category')
data_sample.head()
```

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample

data_sample


```
In [92]: ridge_tuned_test_prediction = ridge_tuned.predict(X_test)
ridge_tuned_train_prediction = ridge_tuned.predict(X_train)

print_results_test(y_test, ridge_tuned_test_prediction)
print_results_train(y_train, ridge_tuned_train_prediction)
```

```
R2 for test set : 0.928971632080551
MSE (mean squared error) for test set: 87.54340999204696
-----
R2 for train set : 0.972976916178329
MSE (mean squared error) for train set: 87.6634510133845
```

R2 is slightly higher in this case, and MSE is a little lower compared to the case when we used alpha =0.1. It seems changing this parameter has worked well for our model. The results we got before were:

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[93]:	# see results with initial parameters				
	results_df[1:2]				

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[93]:	0 Ridge Regression 0.972895 87.54992 0.972978 87.658762				

```
In [94]: # we will save this results
results_df_6 = pd.DataFrame(data=[('Tuned Ridge Regression', get_results(y_test, ridge_tuned_test_prediction)[1],
get_results(y_train, ridge_tuned_train_prediction)[1],
get_results(y_train, ridge_tuned_train_prediction)[0]),
('R2_test', 'MSE_test', 'R2_train', 'MSE_train')])
results_df = pd.concat([results_df, results_df_6])
```

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[95]:	0 Ridge Regression 0.972895 87.54992 0.972978 87.658762				

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[95]:	0 Simple Linear Regression 0.909596 290.479789 0.909844 293.118977				
	0 Ridge Regression 0.972895 87.549920 0.972978 87.658762				
	0 Random Forest Regression 0.988119 38.377083 0.998283 5.570921				
	0 XGBoost Regression 0.988203 38.105111 0.991336 28.107768				
	0 Simple Linear Regression no intercept 0.909072 292.166020 0.909467 294.343506				
	0 Tuned Ridge Regression 0.972897 87.543450 0.972977 87.663451				

Random Forest Regression

Now we'll change parameters of our Random Forest Regression, although default parameter for n_estimators is 100, in the previous model I changed it to 50 to speed up the model, let's see if we got better results setting n_estimators to 100.

```
In [96]: # look for params of our previous model
regressor.get_params()
```

```
Out[96]: {'bootstrap': True,
'ccp_alpha': 0.0,
'criterion': 'squared_error',
'max_depth': None,
'max_features': 'auto',
'max_leaf_nodes': None,
'max_samples': None,
'min_impurity_decrease': 0.0,
'min_samples_leaf': 1,
'min_samples_split': 2,
'min_weight_fraction_leaf': 0.0,
'n_estimators': 50,
'n_jobs': None,
'oob_score': False,
'random_state': 0,
'verbosity': 0,
'warm_start': False}
```

We will now train our model with max_depth=7 and n_estimators = 100 (before it was None and 50).

- Max_depth refers to the minimum number of samples each branch must have after splitting a node.
- n_estimator refers to the number of trees in the model. Before we set it to 50 to speed the model but the default mode is 100 so we will do now 100 trees.

```
In [97]: # instantiate and fit our model with the new parameters
regressor_changed = RandomForestRegressor(n_estimators = 100, max_depth = 7, random_state=0)
regressor_changed.fit(X_train, y_train)
```

```
Out[97]: RandomForestRegressor(max_depth=7, random_state=0)
```

```
In [98]: forest_changed_test_prediction = regressor_changed.predict(X_test)
forest_changed_train_prediction = regressor_changed.predict(X_train)

print_results_test(y_test, forest_changed_test_prediction)
print_results_train(y_train, forest_changed_train_prediction)
```

```
R2 for test set : 0.9782949798548104
MSE (mean squared error) for test set: 70.10824603418195
-----
R2 for train set : 0.97965801386427
MSE (mean squared error) for train set: 66.11959075235038
```

R2 is lower and mse is higher than before!!! Maybe the max_depth parameter was not worth changing, we'll repeat leaving the default and repeat the model only changing number of estimators (number of trees)

```
In [99]: regressor_changed2 = RandomForestRegressor(n_estimators= 100, random_state= 0)
regressor_changed2.fit(X_train, y_train)
```

```
Out[99]: RandomForestRegressor(random_state=0)
```

```
In [100]: forest_changed2_test_prediction = regressor_changed2.predict(X_test)
forest_changed2_train_prediction = regressor_changed2.predict(X_train)
```

```
In [101]: print_results_test(y_test, forest_changed2_test_prediction)
print_results_train(y_train, forest_changed2_train_prediction)
```

```
R2 for test set : 0.9882930131220383
MSE (mean squared error) for test set: 37.81412359805006
-----
R2 for train set : 0.9983593767169089
MSE (mean squared error) for train set: 5.3222163598242656
```

Now we have improved the results of our model, R2 both for test and train is higher than before and MSE is a little lower (let's look at the results we got before):

```
In [102]: # looking at initial results
results_df[2:3]
```

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[102]:	0 Random Forest Regression 0.988119 38.377083 0.998283 5.570921				

In any case, the improvement is very little (we already have a very high R2 so it may be difficult to increase it more). Maybe it's not worth it considering than before we did the regression with half of estimators (much less computational power needed, the execution time in this model is much higher than before!) and got very close results. We will keep this results in our dataframe:

```
In [103]: results_df_7= pd.DataFrame(data = [('Random Forest Doubled Trees', get_results(y_test, forest_changed2_test_prediction)[1],
get_results(y_train, forest_changed2_train_prediction)[1],
get_results(y_train, forest_changed2_train_prediction)[0]),
('R2_test', 'MSE_test', 'R2_train', 'MSE_train')])
results_df = pd.concat([results_df, results_df_7])
results_df
```

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[103]:	0 Simple Linear Regression 0.909596 290.479789 0.909844 293.118977				
	0 Ridge Regression 0.972895 87.549920 0.972978 87.658762				
	0 Random Forest Regression 0.988119 38.377083 0.998283 5.570921				
	0 XGBoost Regression 0.988203 38.105111 0.991336 28.107768				
	0 Simple Linear Regression no intercept 0.909072 292.166020 0.909467 294.343506				
	0 Tuned Ridge Regression 0.972897 87.543450 0.972977 87.663451				
	0 Random Forest Doubled Trees 0.988293 37.814124 0.998359 5.322216				

XGBoost Regression

Let's check the params used in our previous model in XGBoost Regression:

```
In [104]: XGB_regressor.get_params()
```

```
Out[104]: {'objective': 'reg:squarederror',
'base_score': 0.5,
'booster': 'gbtree',
'colsample_bylevel': 1,
'colsample_bynode': 1,
'colsample_bytree': 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': 0,
'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}
```

This time, to find out which parameters we can change we can use Grid Search again, we will look for the combination of parameters that optimizes MSE (before we used GridSearch to maximize R2). We have choosed 3 parameters and two values for each one and we will find out if we can have some improvement in MSE:

```
In [105]: from sklearn.model_selection import GridSearchCV
```

```
In [106]: import xgboost as xgb
# we define the parameters we want to evaluate:
params = {'max_depth': [5,8], 'learning_rate': [0.05, 0.1], 'colsample_bytree': [0.3, 0.7]}
tuned_xgb = xgb.XGBRegressor()
```

```
clf = GridSearchCV(estimator=tuned_xgb,
param_grid=params,
scoring='neg_mean_squared_error',
verbose=1)
clf.fit(X_train, y_train)
print("Best parameters: ", clf.best_params_)
print("Lowest RMSE: ", (-clf.best_score_**(1/2.0))
```

```
Fitting 5 folds for each of 8 candidates, totalling 40 fits
Best parameters: {'colsample_bytree': 0.7, 'learning_rate': 0.1, 'max_depth': 8}
Lowest RMSE: 6.216730134758261
```

```
In [107]: # Output:
# Fitting 5 folds for each of 8 candidates, totalling 40 fits
# Best parameters: {'colsample_bytree': 0.7, 'learning_rate': 0.1, 'max_depth': 8}
# Lowest RMSE: 6.216730134758261
```

Now let's apply the model with the same parameters:

```
In [108]: XGB_tuned_regressor = XGBRegressor(colsample_bytree=0.7, learning_rate=0.1, max_depth=8)
# fit train set
XGB_tuned_regressor.fit(X_train, y_train)
```

```
Out[108]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=0.7, enable_categorical=False,
gamma=0, gpu_id=-1, importance_type=None,
interaction_constraints='', learning_rate=0.1, max_delta_step=0,
max_depth=8, 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 [109]: XGBtuned_test_prediction = XGB_tuned_regressor.predict(X_test)
XGBtuned_train_prediction = XGB_tuned_regressor.predict(X_train)
```

```
In [110]: print_results_test(y_test, XGBtuned_test_prediction)
print_results_train(y_train, XGBtuned_train_prediction)
```

```
R2 for test set : 0.9882761569466629
MSE (mean squared error) for test set: 37.868569845037634
-----
R2 for train set : 0.9914808862918554
MSE (mean squared error) for train set: 27.636183647999896
```

```
In [111]: # results of the model with default parameters:
results_df[3:4]
```

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[111]:	0 XGBoost Regression 0.988203 38.105111 0.991336 28.107768				

Comparing the results with our previous model we see a little improvement for both R2 and MSE for both train and test!! (The improvement is small because results were already very good).

```
In [112]: # keep results in df
results_df_8 = pd.DataFrame(data = [('XGBoost Tuned Regression', get_results(y_test, XGBtuned_test_prediction)
get_results(y_train, XGBtuned_test_prediction)[1],
get_results(y_train, XGBtuned_train_prediction)[0]),
('R2_test', 'MSE_test', 'R2_train', 'MSE_train')])
results_df = pd.concat([results_df, results_df_8])
results_df
```

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[112]:	0 Simple Linear Regression 0.909596 290.479789 0.909844 293.118977				
	0 Ridge Regression 0.972895 87.549920 0.972978 87.658762				
	0 Random Forest Regression 0.988119 38.377083 0.998283 5.570921				
	0 XGBoost Regression 0.988203 38.105111 0.991336 28.107768				
	0 Simple Linear Regression no intercept 0.909072 292.166020 0.909467 294.343506				
	0 Tuned Ridge Regression 0.972897 87.543450 0.972977 87.663451				
	0 Random Forest Doubled Trees 0.988293 37.814124 0.998359 5.322216				
	0 XGBoost Tuned Regression 0.988276 37.868570 0.991481 27.636184				

```
In [113]: # save this results df
results_df.to_csv('results_df.csv', index = False)
```

- Exercici 4

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

We will do Cross-Validation of the 4 models we have used. Since some of the models take a lot of computing time, we will use cv = 5 instead of cv = 10.

Simple Linear Regression

```
In [114]: from sklearn.model_selection import cross_val_score
```

```
In [115]: lr = LinearRegression()
# apply cross validation for r2
r2_scores_lr = cross_val_score(lr, x_sample, y_sample, cv = 10, scoring = 'r2')
r2_lr = np.mean(r2_scores_lr).round(4)
```

```
Out[115]: 0.9077
```

```
In [116]: mse_scores_lr = cross_val_score(lr, x_sample, y_sample, cv = 10, scoring = 'neg_mean_squared_error')
mse_lr = abs(np.mean(mse_scores_lr)).round(4)
```

```
Out[116]: 293.3321
```

```
In [117]: # we are going to create a dataframe for the results
results = pd.DataFrame(data = [('Simple Linear Regression', r2_lr, mse_lr]), columns = ['Model', 'R2', 'MSE',
cv_results
```

	Model	R2	MSE
Out[117]:	0 Simple Linear Regression 0.9077 293.3321		

Ridge regression

Now we'll do the same with Ridge Regression

```
In [118]: # now we apply CV to X and y with ridge model
# apply CV to X and y (not splitted)
ridge = Ridge()
ridge_r2_scores = cross_val_score(ridge, X, y, cv=5)
ridge_r2_scores
```

```
Out[118]: array([0.97113636, 0.96375803, 0.97460197, 0.96741267, 0.97066173])
```

```
In [119]: r2_cv_ridge = np.mean(ridge_r2_scores)
r2_cv_ridge
```

```
Out[119]: 0.96951015141
```

```
In [120]: ridge_mse_scores = cross_val_score(ridge, X, y, cv=5, scoring = 'neg_mean_squared_error')
```

```
In [121]: mse_cv_ridge = abs(np.mean(ridge_mse_scores))
mse_cv_ridge
```

```
Out[121]: 97.4153101417725
```

```
In [122]: ridge_results = pd.DataFrame(data = [('Ridge Regression', r2_cv_ridge, mse_cv_ridge)],
columns = ['Model', 'R2', 'MSE'])
cv_results = pd.concat([cv_results, ridge_results])
cv_results
```

	Model	R2	MSE
Out[122]:	0 Simple Linear Regression 0.90770 293.332100		
	0 Ridge Regression 0.969510 97.415301		

Random Forest Regression

```
In [123]: forest = RandomForestRegressor(n_estimators= 50, random_state= 0)
forest_r2_scores = cross_val_score(forest, X, y, cv = 5)
```

```
In [124]: r2_cv_forest = np.mean(forest_r2_scores)
r2_cv_forest
```

```
Out[124]: 0.9875035113459513
```

```
In [125]: forest_mse_scores = cross_val_score(forest, X, y, cv = 5, scoring = 'neg_mean_squared_error' )
forest_mse_scores
```

```
Out[125]: array([-51.6013453, -39.2487797, -40.37977414, -38.0480514,
-31.36632787])
```

```
In [126]: mse_cv_forest = abs(np.mean(forest_mse_scores))
mse_cv_forest
```

```
Out[126]: 40.12884643173862
```

```
In [127]: # save it to a dataframe
forest_results = pd.DataFrame(data = [('Random Forest Regression', r2_cv_forest, mse_cv_forest)],
columns = ['Model', 'R2', 'MSE'])
cv_results = pd.concat([cv_results, forest_results])
```

XGBoost Regression

```
In [128]: xgb_regressor = XGBRegressor()
xgb_r2_scores = cross_val_score(xgb_regressor, X, y, cv = 5)
xgb_r2_scores
```

```
Out[128]: array([0.98736239, 0.9869738, 0.98933323, 0.98577245, 0.98783103])
```

```
In [129]: r2_cv_xgb = np.mean(xgb_r2_scores)
```

```
Out[129]: 0.9874545798528283
```

```
In [130]: xgb_mse_scores = cross_val_score(xgb_regressor, X, y, cv = 5, scoring = 'neg_mean_squared_error')
xgb_mse_scores
```

```
Out[130]: array([-49.51210968, -40.63479463, -39.81449616, -37.4899658,
-33.27153279])
```

```
In [131]: mse_cv_xgb = abs(np.mean(xgb_mse_scores))
mse_cv_xgb
```

```
Out[131]: 40.144583804872454
```

```
In [132]: # save results to dataframe
xgb_results = pd.DataFrame(data = [('XGBoost Regression', r2_cv_xgb, mse_cv_xgb)],
columns = ['Model', 'R2', 'MSE'])
cv_results = pd.concat([cv_results, xgb_results])
```

```
In [133]: cv_results
```

	Model	R2	MSE
Out[133]:	0 Simple Linear Regression 0.907700 293.332100		
	0 Ridge Regression 0.969510 97.415301		
	0 Random Forest Regression 0.987504 40.128846		
	0 XGBoost Regression 0.987455 40.144584		

```
In [4]: # look at the df with the results of our models using train/test split
results_df = pd.read_csv('results_df.csv')
```

	Model	R2_test	MSE_test	R2_train	MSE_train
Out[4]:	0 Simple Linear Regression 0.909596 290.479789 0.909844 293.118977				
	1 Ridge Regression 0.972895 87.549920 0.972978 87.658762				
	2 Random Forest Regression 0.988119 38.377083 0.998283 5.570921				
	3 XGBoost Regression 0.988203 38.105111 0.991336 28.107768				
	4 Simple Linear Regression no intercept 0.909072 292.166020 0.909467 294.343506				
	5 Tuned Ridge Regression 0.972897 87.543450 0.972977 87.663451				
	6 Random Forest Doubled Trees 0.988293 37.814124 0.998359 5.322216				
	7 XGBoost Tuned Regression 0.988276 37.868570 0.991481 27.636184				

If we compare the df with the results of our Cross Validation and the df with our train/test results :

- MSE is a little higher for all the models (compared with train/test split) specially in the models that got better results: Random Forest and XGBoost.
- R2 is a little lower in Cross Validation but the difference is not worrying.
- We could say that our model does not suffer from overfitting, because although the results are worse than with train/test, this difference is very small and Cross Validation still shows very good results.

The exercise is continued in another Notebook in the same Github link.

```
In [134]: # save file to csv
cv_results.to_csv('cv_results', index = False)
```