

# **Mini Project 1**

Machine learning (CS582), MIU <sup>1</sup>

Baraa Mousa Noufal

Ivan Krasowski Bissio

July 24, 2021

<sup>1</sup>Instructor: Anthony Sander

# Contents

<b>Abstract</b>	<b>1</b>
<b>1 About the Dataset</b>	<b>1</b>
1.1 EDA: Exploratory Data Analysis . . . . .	1
1.2 Data Preparation . . . . .	2
<b>2 Models</b>	<b>2</b>
2.1 KNN . . . . .	2
2.2 Decision Tree . . . . .	2
2.3 Support Machine Vector . . . . .	2
2.3.1 Presentation . . . . .	2
2.3.2 Defining Parameters . . . . .	2
2.3.3 Model Evaluation . . . . .	3
2.4 Neural Network . . . . .	4
2.4.1 Presentation . . . . .	4
2.4.2 Defining Hyperparameters . . . . .	5
2.4.3 Model Evaluation . . . . .	5
2.5 Model 5 . . . . .	6
2.6 Ensemble: Random Forest . . . . .	6
2.7 Voting Classifier . . . . .	6
<b>3 AUCs</b>	<b>6</b>
<b>4 AutoML</b>	<b>6</b>
4.1 Presentation . . . . .	6
4.2 Obtaining the Model . . . . .	6
<b>5 Best model: model n</b>	<b>6</b>
5.1 PCA . . . . .	6
<b>6 Conclusion</b>	<b>6</b>

## **Abstract**

# 1 About the Dataset

The chosen dataset collects information about Airline Passenger Satisfaction. We opted for this dataset because it had a lot of entries (almost 130000), many columns (including 4 categorical) and it was oriented to a binary classification problem ('satisfied' vs. 'neutral or dissatisfied').

## 1.1 EDA: Exploratory Data Analysis

- The dataset collects multiple features related to each passenger-flight, including labels for whether the passenger was satisfied with the flight or not.
- It has 24 columns (22 + id + result).
- Id column is dropped (it would add weight).
- Observation: the only column with null values is 'arrival\_delay\_in\_minutes', with 393 missing values.
- Categorical columns: 4

Features and possible values: Gender (2), customer\_type (2), type\_of\_travel (2), customer\_class (3)

Decision: use one-hot encoding in all of them (none of them is a clear candidate for being weighted)

- Looking for strong correlations: pairwise correlation function to check if two features show strong correlation.

NOTE: 'satisfaction' (the label we will try to predict) is not strongly correlated ( $>0.7$ ) with any of the other features.

*arrival\_delay\_in\_minutes* and *departure\_delay\_in\_minutes* have the highest correlation rate (0.965291), which semantically makes sense; all other variables are less than 75% correlated.

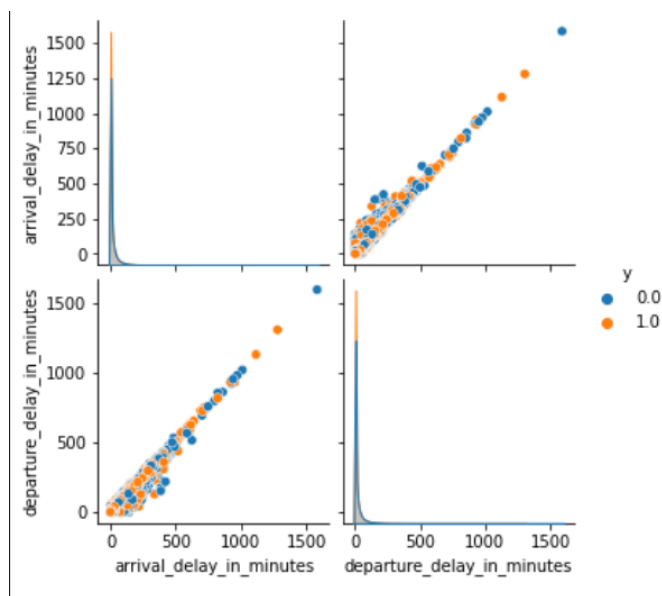


Figure 1: Correlation between most-correlated features

## 1.2 Data Preparation

1. X: all table except id and result ('satisfaction') columns; y: 'satisfaction' column
2. Since the arrival\_delay feature is highly correlated with the departure\_delay feature, and the missing values are not that many (393 out of 129879: 0.03%), we decide to remove the column.
3. We see there are 4 categorical features, with no more than 3 unique values each. So, given that none of them is a clear candidate for being weighted, we decide to use one-hot encoding in all of them.
4. Finally, we split X and y for training and validating, following a ratio of 80%/20% (the dataset is large enough).

Variables: ***X\_train, X\_val, y\_train, y\_val***

## 2 Models

### 2.1 KNN

### 2.2 Decision Tree

### 2.3 Support Machine Vector

#### 2.3.1 Presentation

Support Vector Machines are a model of supervised learning.

In summary, the model finds the hyperplane (kernel) that maximizes the Margin of Safety for classifying.

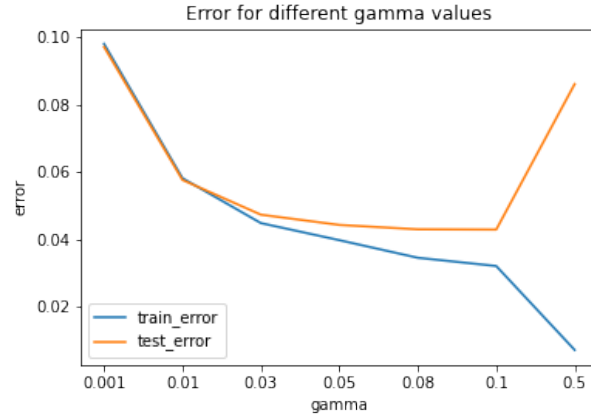
#### 2.3.2 Defining Parameters

The data for training and validating is already defined by the Data Preparation step (80% train, 20% validation).

Representative parameters for the model are *gamma* and *C*

- gamma: kernel coefficient
- C: regularization parameter (higher C, higher variance)

We train the following values for the *gamma* parameter: [0.001, 0.01, 0.03, 0.05, 0.08, 0.1, 0.5]  
To improve the predictions, we regularize the data using a StandardScaler (removes the mean and scales to unit variance) The training error and the validation error for each value allow us to plot a Complexity Curve, to select the optimal value for the parameter.

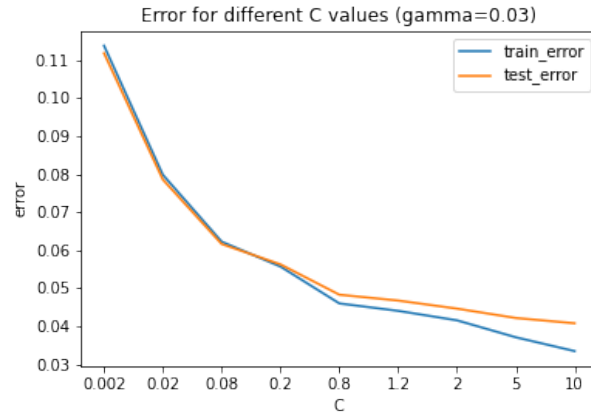


**Figure 2:** Complexity Curve:  $\gamma$  value for SVM

By observing the plot, we determine that the best value for gamma is 0.03

We train the following values for the  $C$  parameter: [0.02, 0.2, 0.8, 1.2, 2, 5, 10]

Again, we regularize the data using a StandardScaler (removes the mean and scales to unit variance)  
The training error and the validation error for each value allow us to plot a new Complexity Curve, to select the optimal value for the parameter.



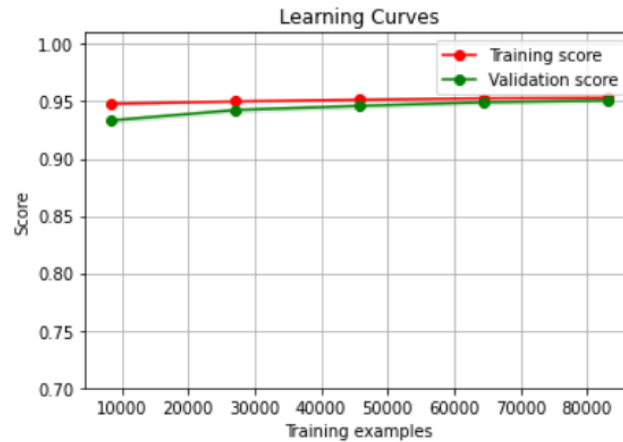
**Figure 3:** Complexity Curve:  $C$  value for SVM ( $\gamma=0.03$ )

By observing the plot, we determine that the best value for  $C$  is 0.8

### 2.3.3 Model Evaluation

Chosen the parameters: ( $\gamma=0.03$ ,  $C=0.8$ ), a learning curve shows us the training and validation scores for different data sizes.

This way, we are able to say that the score is bounded below 95%, and the model doesn't seem to continue learning after 65000/70000 training rows.

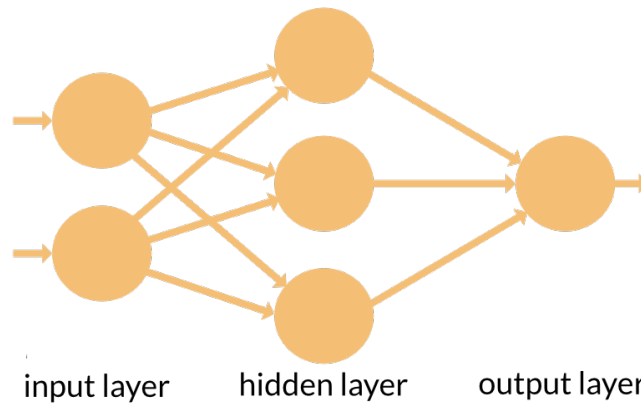


**Figure 4:** Learning Curve for SVM

## 2.4 Neural Network

### 2.4.1 Presentation

Neural Networks (or Multilayer Perceptrons) are a model of supervised learning, represented by composed non-linear functions.



**Figure 5:** Multilayer Perceptron / Neural Network

1

They consist of "neurons", each of which represents a non-linearity, divided in layers, each of which depends on the previous (the first one, on the input data) and is crucial for the next one (the last one gives the output)

They can be seen as a function, whose interface is defined by the "input layer" (the first layer) and the "output layer" (the result).

All layers between the input layer and the output layer are "hidden layers".

A Neural Network makes its predictions after having its parameters (weights) trained with labeled example data (train data)

<sup>1</sup><https://appliedgo.net/perceptron/>

### 2.4.2 Defining Hyperparameters

The data for training and validating is already defined by the Data Preparation step (80% train, 20% validation).

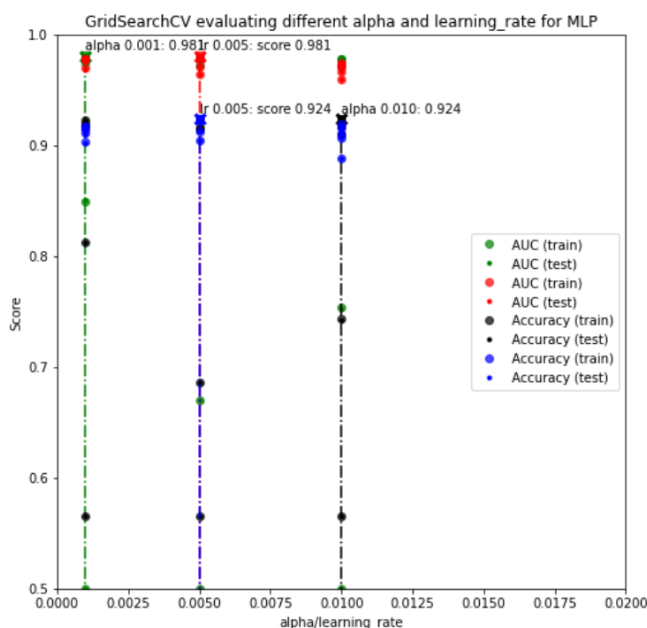
Representative hyperparameters (different from the parameters to be trained) for the model are *alpha* and *learning\_rate*

- alpha: regularization parameter (L2 penalty).
- learning\_rate: parameter for controlling the step size when updating the weights.

A Grid Search is used for getting the best combination of parameters (*alpha*, *learning\_rate*) for the input and the model.

Values used for *alpha*: [0.001, 0.005, 0.01, 0.05, 0.1]

Values used for *learning\_rate*: [0.001, 0.005, 0.01, 0.05, 0.1]



**Figure 6:** Grid Search results

Observable on the plot, and also obtainable by requesting the *grid\_search.best\_params\_*, the best values for the hyperparameters are:

- alpha: 0.001
- learning\_rate: 0.005

### 2.4.3 Model Evaluation

Chosen the parameters: (*alpha*=0.001, *learning\_rate*=0.005), a learning curve shows us the training and validation scores for different data sizes.

This way, we are able to say that the model seems to continue learning after 90000 training rows.



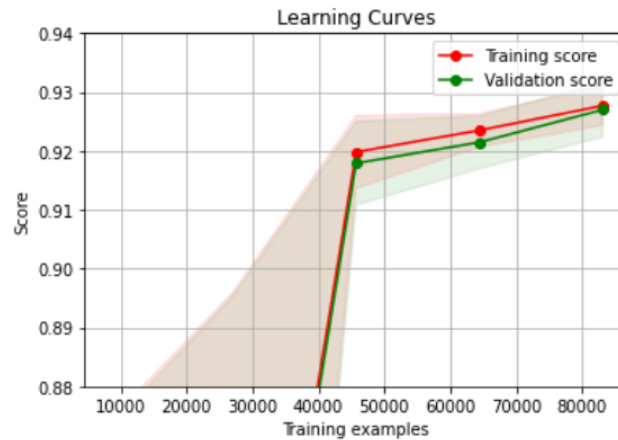


Figure 7: Learning Curve for MLP

## 2.5 Model 5

## 2.6 Ensemble: Random Forest

## 2.7 Voting Classifier

# 3 AUCs

# 4 AutoML

## 4.1 Presentation

AutoML is a tool for choosing the best model (or ensemble) for fitting and predicting a dataset. If using scikit-learn, the proper toolkit is defined as auto-sklearn.

## 4.2 Obtaining the Model

Simply fitting the training data to the `AutoSklearnClassifier` gives us an "optimal" model. In this case (shown by `running model.show_models()`), the result is an ensemble of several ensembles, involving SVMs, Random Forests, etc.

# 5 Best model: model n

## 5.1 PCA

# 6 Conclusion