

Data Science | Lab 2: Build Your First Classifier

Learning Goals

libraries: mainly sklearn

- implement a basic ML process using a MinDist classifier
 - correct sampling; difference between training and test data
 - reproducible experiments
 - basic feature generation
-

Minimum Distance Classifier

Training

Compute the class means $\mu^{(c)}$ from the training data and store those. Those vectors are the representatives of each class.

Testing

For each data \mathbf{x} in `test`, predict the class by the class of the closest representative w.r.t. Euclidean distance,

$$\pi(\mathbf{x}) = \arg \min_{c \in \mathcal{L}} \|\mathbf{x} - \mu^{(c)}\|,$$

where $\|\cdot\|$ denotes the Euclidean (standard) norm.

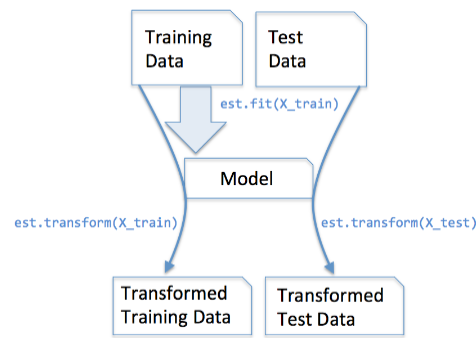
Recall the principles of the Minimum Distance (MinDist) classifier from the lecture. Take a look at scikit-learn's [git repo](#) to verify the classifier's implementation. Note that in scikit-learn, the MinDist classifier is called **Nearest Centroid** classifier.

Sklearn Essentials

A lot of sklearn objects such as preprocessors, vectorizers or scalers implement the following three methods:

- `fit`: fit the object to the current data, i.e. estimate statistics, model distributions, generate parameters etc.
- `transform`: apply model to the data using the extracted parameters from the fit-step
- `fit_transform`: short-hand method that calls and transform on the same data with just one call

Take, for example, the `StandardScaler`, which has two scaling parameters: mean and standard deviation. The scaling parameters should be estimated from the training data and **reused** on the test data. We want to test whether the assumptions we made from the training data **generalize**, i.e. are valid also on the test data.



Source: [stackoverflow](#)

TLDR: Don't use `fit` on new data!

Tasks

Basic MinDist Classification

1. Load the wine data again, reusing the code from Lab 1. Only `X` and `y` are needed.
2. Implement the ML cycle using MinDist as classifier directly applied to the (unscaled) wine data.

```

1 from sklearn.neighbors import NearestCentroid
2 from sklearn.metrics import accuracy_score
3 from sklearn.model_selection import train_test_split

```

Note that the Minimum Distance classifier is called `Nearest Centroid classifier` in sklearn.

```

1 # TODO: apply sampling
2 mindist = NearestCentroid()
3 mindist.fit(X_train, y_train)
4 y_hat = mindist.predict(X_test)
5 # TODO: evaluate the performance

```

3. What accuracy score does your classifier achieve? Compare with your colleagues!
4. Fix this issue by setting a random state.

The following steps are intended to be implemented on top of each other, i.e. add sampling, add feature extraction, add scaling.

Sampling

1. Take a look into the sklearn documentation of the `train_test_split` method to find out the default split ratio, i.e. how much data is used for the training and test set, respectively.
2. Change the default setting to a ratio of 7:3 for training and test set.
3. Rerun the classification.

Feature Extraction

1. Remove highly correlated features.

Identify the feature with highest correlation values throughout the dataset. Remove this feature from your dataset and rerun the classification. Does the accuracy score change?

```

1 X_new = X.drop('column_name', axis=1)

```

2. Can you think of other feature extraction or generation methods?

Scaling

1. Apply standard scaling to the data and rerun the classification.
2. Does the accuracy change and if yes, why?

Homework

Reusing the same `random_state` throughout, track your experiments and note down the accuracies:


Short description of training setting (*)	accuracy
most basic setting (see above)	
highest correlating feature removed	
...	

(*) This could include: the classifier used, train-test split ratio, the preprocessing methods applied, any type of feature extraction or generation, ...

In the next labs, you will update this table with new figures, so make sure to save it (either directly in the notebook, as Excel file or on a sheet of paper).

Quiz

Take the quiz in Moodle. Make sure to have your Python notebook open and your code up and running.

 **Deadline: 8.11.2022 (6:45 pm)**