

Data Camp TP4

Instructions

The goal is to make submissions to [To be or not to be](#) challenge and to explore various machine learning techniques. The best way to know if a technique will improve your final score is to **try it**. Below, we give you some ideas of methods to try above - **you don't need to implement them all!**. Simply try some that seem sound in your opinion. You are free to put any code, experiment, plot, explanation, or analysis in your notebook.

This TP will be graded. You have until December the 3rd to send your notebook. You can work alone or in a team of 2 or 3 persons; we recommend teams of 2.

Please structure your notebook with three visible parts:

- **Pre-processing**: include here any code, text, or plot related to the data pre-processing such as dimensionality reduction, encoding, etc.
- **Model**: include here any code, text, or plot related to the classification algorithm(s), the hyperparameters, the validation method, etc.
- **Visualization**: use this part to add some plots. It can be exploratory analysis (e.g. 2D plots of the data) or results analysis (e.g. confusion matrix)

Grading (/20)

- Pre-processing (/5)
- Model (/5)
- Visualization (/5)
- Originality of your work (/5)
- *Bonus if your final score is better than the baseline (+1)*
- *Bonus if you're ranked 1st, 2nd, or 3rd on the leaderboard (+1)*
- *Bonus if you answered the 3 questions in the notebook correctly (+1)*

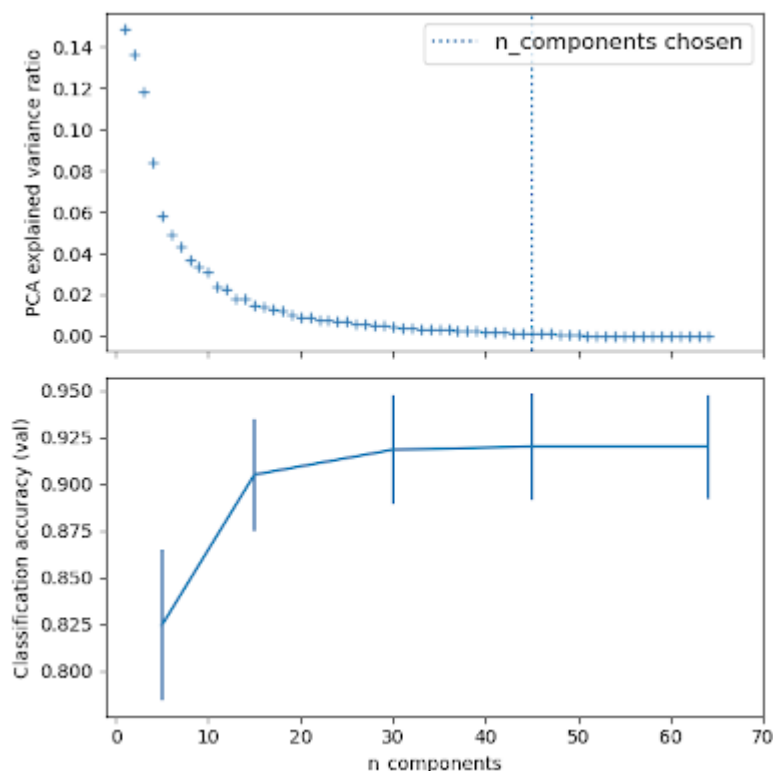
Techniques

Here you can find various ideas of methods to try to solve the problem. Most of them can be combined and can increase - or decrease - your final score. Think of this document as a kind of general machine learning cheat sheet.

To evaluate the benefit of those methods it is recommended to:

- Use a hold-out set or a cross-validation method and report the variance,
- Display the scores with Visualization (plots),
- Compare scores with only one component varying (e.g. different models but same pre-processing)

Here is an example of the kind of experiments you can put in place to improve your model.



In this plot, you can see the accuracy obtained by an algorithm on the validation set for different numbers of dimensions (dimensionality reduction method used: PCA).

1. Missing values

If you have missing values (NaN) in your data, you must remove them or replace them. The main methods are:

- Remove rows or remove columns that contain NaN
- Replace NaN with a value. It can be a specific value such as 0, or a computed value from the other values in the same row or column (mean, median, most occurring value, etc.)
- Using a machine learning model to infer the value given neighbor values

Here is a simple code to replace the missing values from your `pd.DataFrame` object:

```
x = x.fillna(method='pad')
x_test = x_test.fillna(method='pad')
```

2. Encoding

Encodings are useful to be able to give categorical variables as an input of your machine learning algorithm. Their main purpose is to transform **strings into integers**.

Some ideas of encoding:

- Ordinal encoding (or label encoding)
- Frequency encoding
- Likelihood encoding
- Target encoding
- One-hot encoding (*WARNING: it creates a new variable for each category and therefore can be very memory consuming!*)
- Deep learning embedding

- ...

Here is a piece of code useful to apply **ordinal encoding** to your data. As you can see, it is important to store the `str -> int` mapping with a dictionary to apply it to another subset of the dataset.

```
class OrdinalEncoder():
    """ Ordinal encoding.
        ['fish', 'cat', 'fish', 'dog'] -> [0, 1, 0, 2]
    """
    def __init__(self):
        # Dictionnary containing the saved mapping from String to Integers
        self.mapping = {}

    def fit(self, variable):
        """ Fit the encoding from a categorical variable.
            :param variable: 1D vector (list, np.array, pd.Series)
        """
        i = 0
        for e in variable:
            if e not in self.mapping: # This category is not encoded
                self.mapping[e] = i
                i += 1

    def transform(self, variable):
        """ Apply the encoding to a categorical variable.
            :param variable: 1D vector (list, np.array, pd.Series)
        """
        i = len(self.mapping)
        encoded = variable.copy()
        for j, e in enumerate(variable):
            if e not in self.mapping: # Unseen category
                self.mapping[e] = i
                i += 1
            encoded[j] = self.mapping[e]
        return encoded
```

To use it on your dataset:

```
X_encoded = X.copy()
X_test_encoded = X_test.copy()

str_columns = X.select_dtypes(include='object')

for column in str_columns:
    encoder = OrdinalEncoder()
    encoder.fit(X[column])
    X_encoded[column] = encoder.transform(X[column])
    X_test_encoded[column] = encoder.transform(X_test[column])
```

HINT: You can treat separately the variables representing timestamps (e.g. convert them into seconds).

3. Normalization

You can try to normalize your features.

Some ideas of normalization:

- Min-max normalization
- Standard normalization

[Here is more information on mix-max normalization](#)

4. Dimensionality reduction and feature selection

You can reduce the number of features of your data. There exist mainly two ways of doing it:

1. Combining your features (dimensionality reduction)
2. Removing some features (feature selection)

Some ideas:

- Dimensionality reduction methods: PCA, LDA, autoencoder, T-SNE, feature hashing
- Feature selection methods: Remove by hand, remove the variables with the worst Pearson's correlation with target, Lasso, tree-based method, recursive feature elimination.

[Here is the PCA from the scikit-learn package.](#)

//\ **Don't forget to vary the number of dimensions to see what works best!**

//\ **Don't forget to apply the same pre-processing to your train and test sets!**

5. Up-sampling and down-sampling

The label we try to predict is imbalanced: there are, fortunately, many more patients who leave the hospital alive.

To reduce or suppress the imbalance and **help your model in its training**, you can do:

- Up-sampling: re-sample (with replacement) new points from the least represented class
- Down-sampling: keep only a subset of most represented classes

This simple pre-processing can be very effective.

6. Model selection and hyperparameter selection

Last but not least, you need to choose a machine learning model. Each model has a set of hyperparameters that influence its behavior and therefore its performance.

Ideas of models:

- [Logistic Regression](#)
- [KNN](#)
- [Random Forest](#)
- [SVM](#)
- [Multi-layer perceptron](#)
- ... *(there are tons of models)*
- Combine several with an [ensemble of models](#)

//\ **Don't forget to read about the models' hyperparameters and to make their values vary!**

7. Visualization

Let's divide the visualizations into two groups:

- **Exploratory analysis:** to have a better understanding of the data. Some examples:
 - Show a specific variable with a simple graph or bar plot
 - Show the relationship between two variables with a scatter plot
 - Show your data in 2D or 3D using a dimensionality reduction algorithm
 - Show the correlation between your features and the target variable

- Show the distribution of the target variable
- A heatmap of the dataset
- Try to think of your own visualizations that could be interesting for your specific dataset
- **Results analysis:** to understand your models' behaviors. Some examples:
 - [Confusion matrix](#)
 - Compare several models with a box plot
 - Show the variation of the score according to a hyperparameter
 - Visualize your models' output
 - Try to think of your own interesting visualizations

The most common Python library for visualizations is [Matplotlib](#). You can also try [Seaborn](#) which is based on Matplotlib. Have a look at [Seaborn's gallery](#), it's beautiful!

Submit to Codalab

- **You must submit to Codalab**
- Create your submission file using the Jupyter notebook. You need to have a file named "mimic_synthetic_test.csv" inside a ZIP archive.
- Go to `Participate > Submit / View Results` and click on `Submit` button to upload your submission.
- If your submission fails, have a look at the scoring output and scoring error logs! Click on the cross "+" on the right of the submission interface (see the image below).

14	0.5181035075	starting_kit_submission.zip	11/19/2020 15:47:12	Finished		+
15	0.5181035075	starting_kit_submission.zip	11/19/2020 15:47:14	Finished	✓	-

Description:

update description

[Download your submission](#)
[View scoring output log](#)
[View scoring error log](#)
[Download output from scoring step](#)

Remove from Leaderboard

- If the error states `Bad prediction shape: (20002, 1) != (20001, 1)` it means that your `mimic_synthetic_test.csv` file is one line too long. Each line has to be one prediction (but the notebook's code should do that for you). Maybe it's saving it with a header, **check the first line** of your file, and add `header=False` to the `to_csv` method.

Computing power and Jupyter Hub

You can try out [LAL's Jupyter Cloud](#) if you need more computing power.