

Machine Learning - SS 2021

Exercise 3: Synthetic Data

Federico Ambrogi, e1449911@student.tuwien.ac.at
Adam Höfler, e11847620@student.tuwien.ac.at
Matteo Panzeri 12039996@student.tuwien.ac.at
TU Wien

July 24, 2021

Contents

1	Data Sets Description	1
1.1	Income Data Set	1
1.2	Titanic Data Set	2
1.3	Social Data Set	2
2	Data Exploration and Pre-processing	3
2.1	Data Overview	3
3	Generation of Synthetic Data	3
3.1	GaussianCopula	4
3.2	CTGAN model	4
3.3	Copula GAN	5
4	Model Implementation	5
4.1	Holdout and Cross Validation	5
5	Performance Tests	5
6	Conclusion	7

Introduction

In this document we describe the results of the implementation of three algorithms to solve classification problems on four different dataset.

1 Data Sets Description

Here we briefly introduce our four datasets, analyze the distributions of the input data and check for significant correlations between variables.

1.1 Income Data Set

The "income" data set[?] The distribution of the relevant features, as well as a scatter plot to highlight pair-wise correlations are shown in Fig. 5, while a more general view on the correlations are shown in Fig. 6.

The dataset provides us information about people, like: age, workclass, education, marital-status, occupation, capital-gain, ecc...

Our purpose is to determine if a person income is lower or greater than a given value specifically the threshold is 50000 Dollars per year.

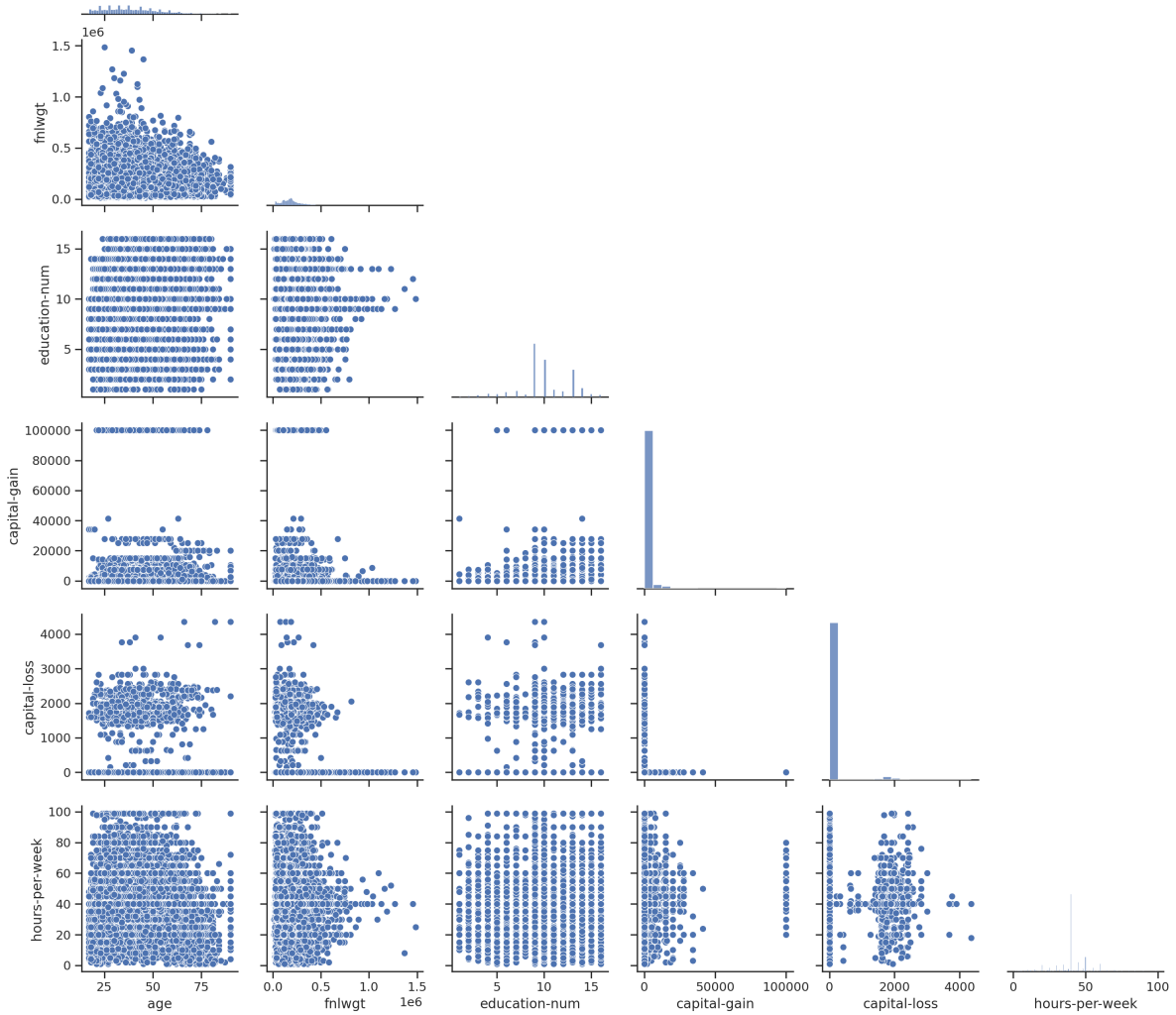


Figure 1: Distributions and pair-wise correlation for the features of the "income" data set.

1.2 Titanic Data Set

The Titanic dataset contains informations about the Titanic's passengers. Given the observations we want to determine if a given person ,represented by an observation, will survive on the Titanic. Some features of the dataset are:Name, sex, age, Passenger class, ticket fare, where the passenger embarked ecc...

1.3 Social Data Set

The Social network ads datasets collects information about social network users in particular age,sex and income. Our objective is to understand if an user will buy the advised good.

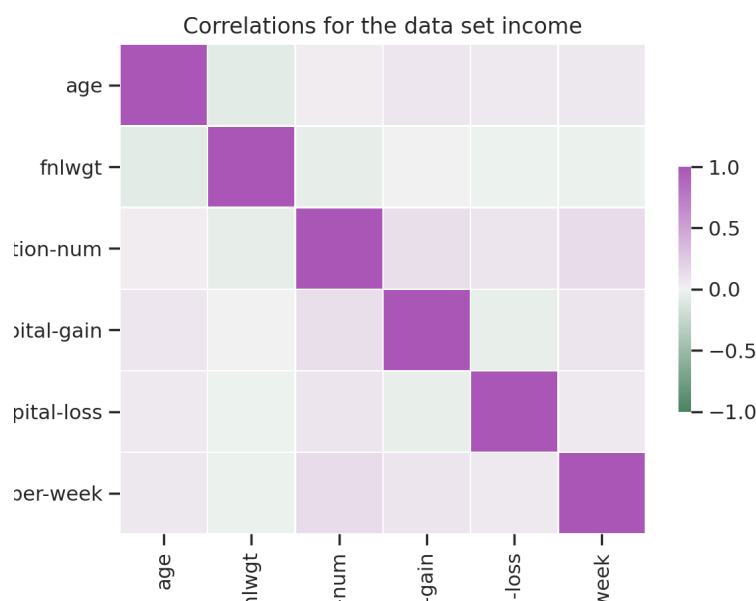


Figure 2: Correlation matrix for the "income" data set.

2 Data Exploration and Pre-processing

In this section we describe the necessary preliminary steps to import and prepare the data to make them suitable for the learning algorithms.

The steps are handled by the script *data_preparation.py* which includes dedicated function to pre-process each data set.

2.1 Data Overview

Here we provide an overview of the original data set data.

Missing Value The

Normalization The

Labelling In all datasets were present not numeric features. In order to synthesize them was necessary to change the strings into numeric values. In order to do so we used the `cat.code` method in this way for every not numeric features all the possible values were labelled using a different value. Using the labelling we had the possibility to use the three models.

Outliers The

3 Generation of Synthetic Data

In order to generate the synthetic data we relied on three different generation models. The models are defined in the SDV library. The three models aim to generate a synthetic datasets whose properties are the same as the original dataset. The three models we used are: GaussianCopula model, CTGAN model, CopulaGAN model. We now describe the idea behind these three models as reported on the SDV library user guide.

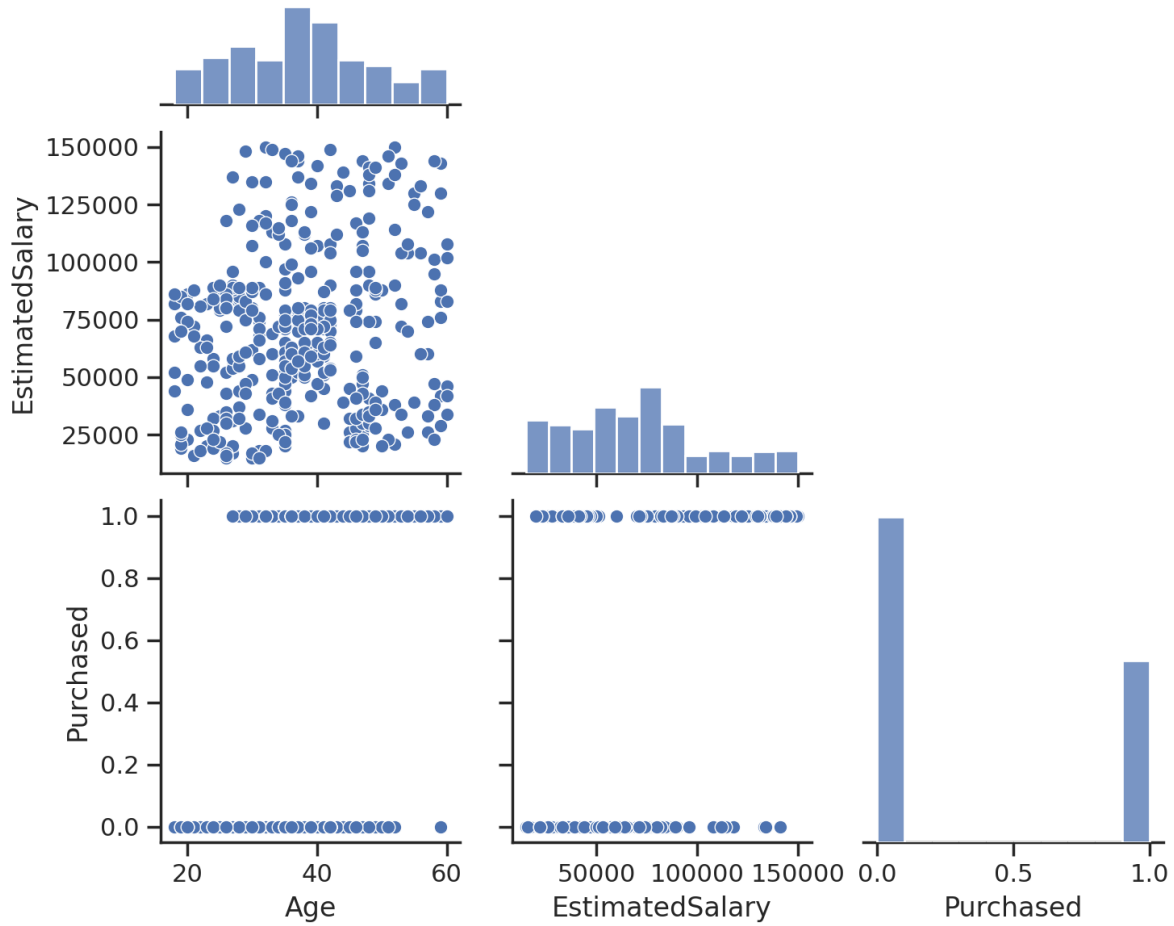


Figure 3: Distributions and pair-wise correlation for the features of the "social" data set.

GaussianCopula The Gaussian copula model is based on copulas which are distributions on $[0;1]^d$ which is constructed from a multivariate normal distribution over \mathbb{R}^d by using the probability integral transform. Intuitively, a copula is a mathematical function that allows us to describe the joint distribution of multiple random variables by analyzing the dependencies between their marginal distributions.

CTGAN model The `sdv.tabular.CTGAN` model is based on the GAN-based Deep Learning data synthesizer which was presented at the NeurIPS 2020 conference by the paper titled Modeling Tabular data using Conditional GAN.

Copula GAN The `sdv.tabular.CopulaGAN` model is a variation of the CTGAN Model which takes advantage of the CDF based transformation that the GaussianCopulas apply to make the underlying CTGAN model task of learning the data easier.

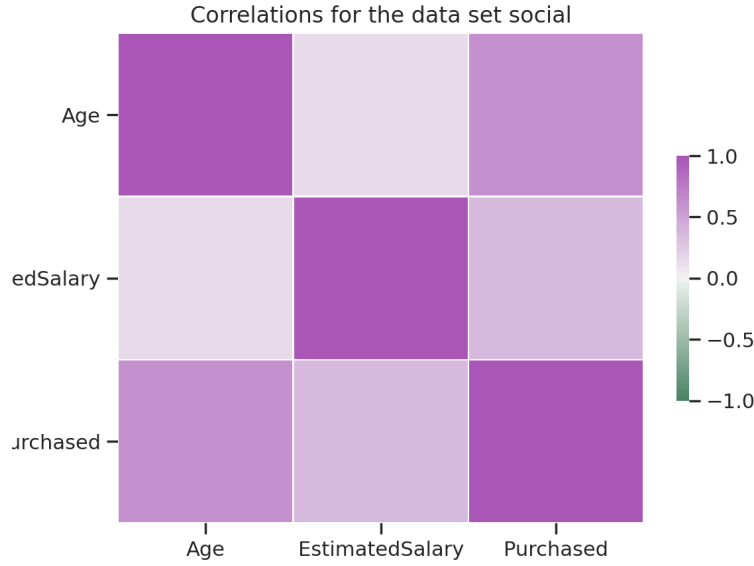


Figure 4: Correlation matrix for the "social" data set.

4 Model Implementation

4.1 Holdout and Cross Validation

Here we describe briefly the techniques of "holdout" and "cross validation" that are used to evaluate a model. The **holdout** method is essentially based on the splitting on the input data set into two subset, one used for training the model, and one used for testing the model, for example in 80% – 20% proportion, although there is not fix recipe for this split. Once the model is trained, the evaluation can be performed on the test data set, and this is therefore possible to check if the prediction of the model match the data. One big issue of this model is that the training strongly depend on the splitting of the initial data set, for example if the characteristic of the training data set are not representative of the whole data set.

A more powerful method is the **cross-validation** or "k-fold cross validation". Form the full dataset, a test is held out for final evaluation, but the validation set is no longer needed. For this, the training set is split into k sets, so that the model is trained using $k - 1$ folds as training data, and the remaining one is used for the validation as a test set (to compute the interesting metric of the model). Once we obtain such k number of metrics, the final result is the average of these parameters, obtained for each iteration of the cross-validation on each distinct fold.

5 Performance Tests

Confusion Matrix For each classifier, we produce a confusion matrix where each entry i, j corresponds to the number of observations in group i , but predicted to be in group j . We chose to normalize the entries according to the sum of each row. In case of binary classification, the matrix reduces to the number of true negatives (TN), false positives (FP), false negatives (FN) and true positives (TP).

Examples can be see in Fig. ??.

We remind here the definition of the metric parameters we will used to quantify the performance of our classifiers i.e. precision (P), recall (R), accuracy (A) and specificity S :

$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN} \quad A = \frac{TP + TN}{all} \quad S = \frac{TN}{TN + FP} \quad (1)$$

It is straightforward to calculate these parameters for binary classification tasks out of the confusion matrix. In case of multiple labels, we need to calculate these parameter for each class, given that:

1. TP s are the values in the diagonal;

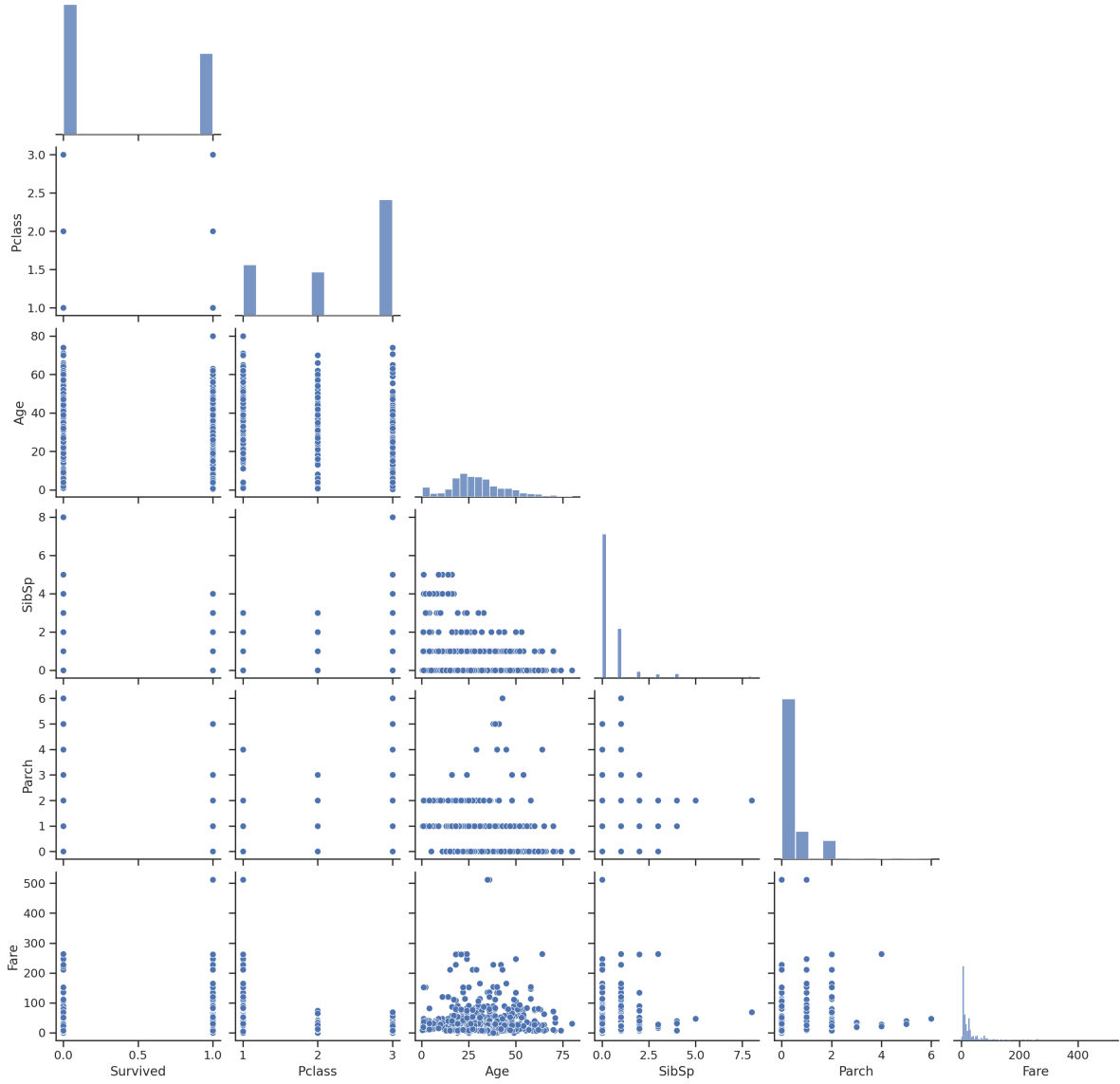


Figure 5: Distributions and pair-wise correlation for the features of the "Titanic" data set.

2. FNs for a certain class are the sum of values in the corresponding row excluding the TP ;
3. FPs for a certain class are the sum of values in the corresponding column excluding the TP ;
4. TNs for a certain class are the sum of all rows and columns, excluding the class's column and row.

Another convenient metric, particularly because it is calculated directly from the proper *scikit-learn* function, is called $f1 - score$, which is defined as the harmonic mean of the precision and recall:

$$f1 - score = 2 \times \frac{P\hat{R}}{P + R} \quad (2)$$

Describe micro averaging and macro averaging, will we report both ?

6 Conclusion

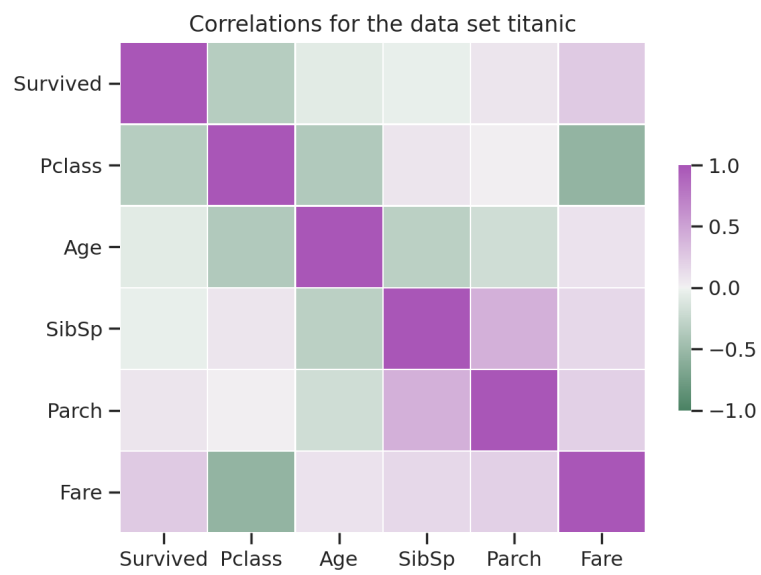


Figure 6: Correlation matrix for the "Titanic" data set.

References