

DATA QUALITY REPORT

HOMEWORK 1

Giuseppe Benanti –ID 07266120

INTRODUCTION

Extensive explanation of the different CRISP-DM steps for this predictive analytics project are given in the correlated Jupyter notebook (*Homework1_Comp47350_Benanti_07266120.ipynb*).

Here we summarize in tables and graphs the results from the **Data Exploration** step of the CRISP-DM process deployed in Anaconda-Jupyter environment using Python 3.6 on the dataset stored in the file *amazon-offers-10k-samples-raw.csv*.

Below, Tables 1 to 4 list descriptive statistics results from the original dataset on offers from vendors on Amazon.

The original dataset is comprised of a totality of 10000 entries (rows), one per offer, and 21 columns, with 20 descriptive features and 1 target feature, and a numeric column for row number.

Each product can be offered by many vendors, but only one of these vendors (per product) will become a 'Winner' vendor and listed at the top of the user research page. The binary categorical target feature 'IsWinner' expresses this concept, 1 for winner, 0 for loser.

The predictive model will attempt to predict the target feature based on a selection of the descriptive features data.

1 - CONTINUOUS FEATURES

Table 1. Descriptive statistics results for continuous (numerical) features. Dataset: *amazon-offers-10k-samples-raw.csv*.

[illegible]

2 - CATEGORICAL FEATURES

Table 3. Descriptive statistics results for categorical features. Dataset: *amazon-offers-10k-samples-raw.csv*. Constant features (cardinality = 1) are in red.

[illegible]

Table 4. Descriptive statistics results for categorical features. Revised table after removal of duplicates, constants and features with high missing values %. Dataset: *offers2.csv*.

[illegible]

3 – DATA VISUALIZATION

Fig 1. Histograms for continuous features.

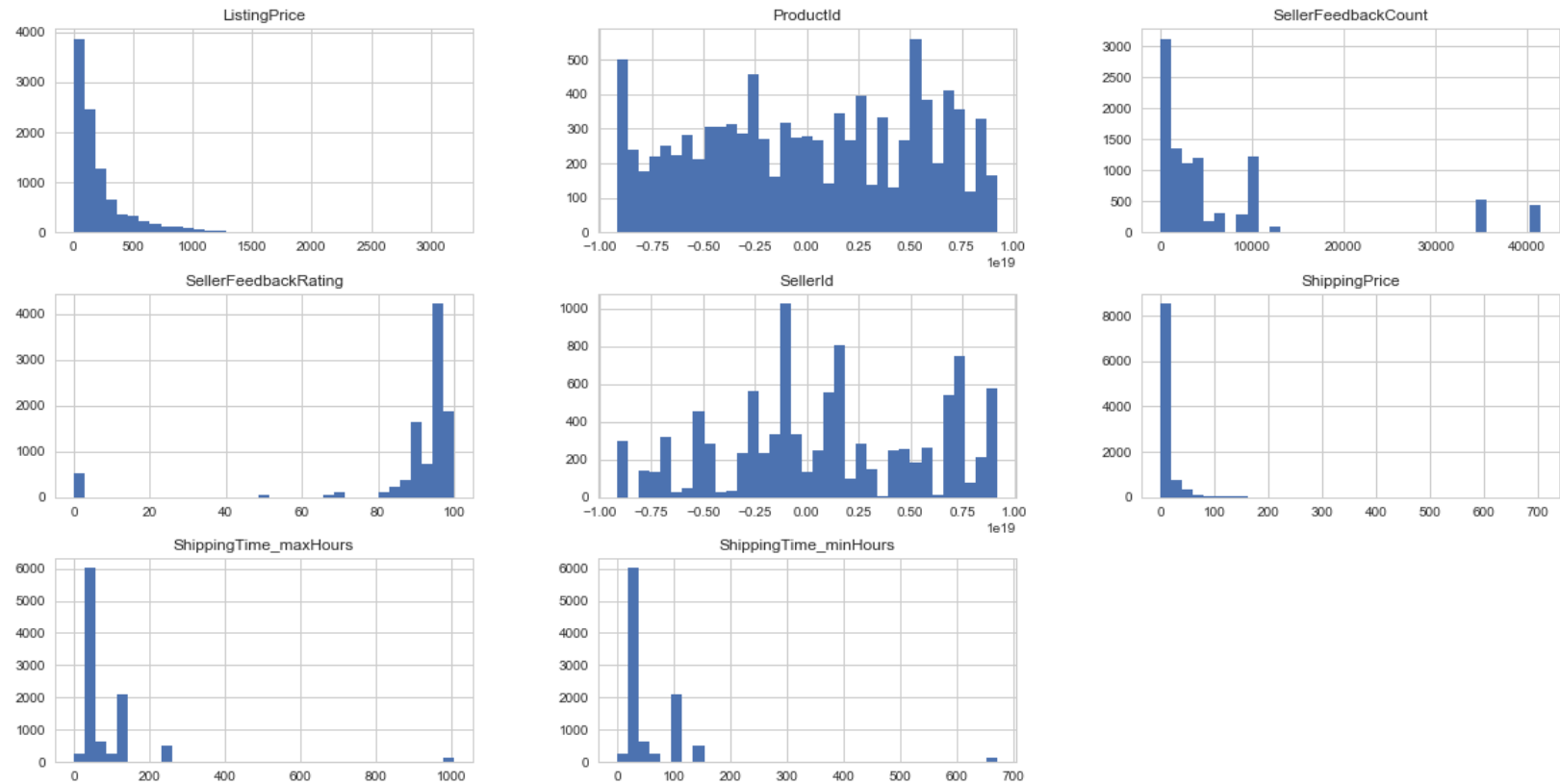


Fig 2. Box plots for continuous features.

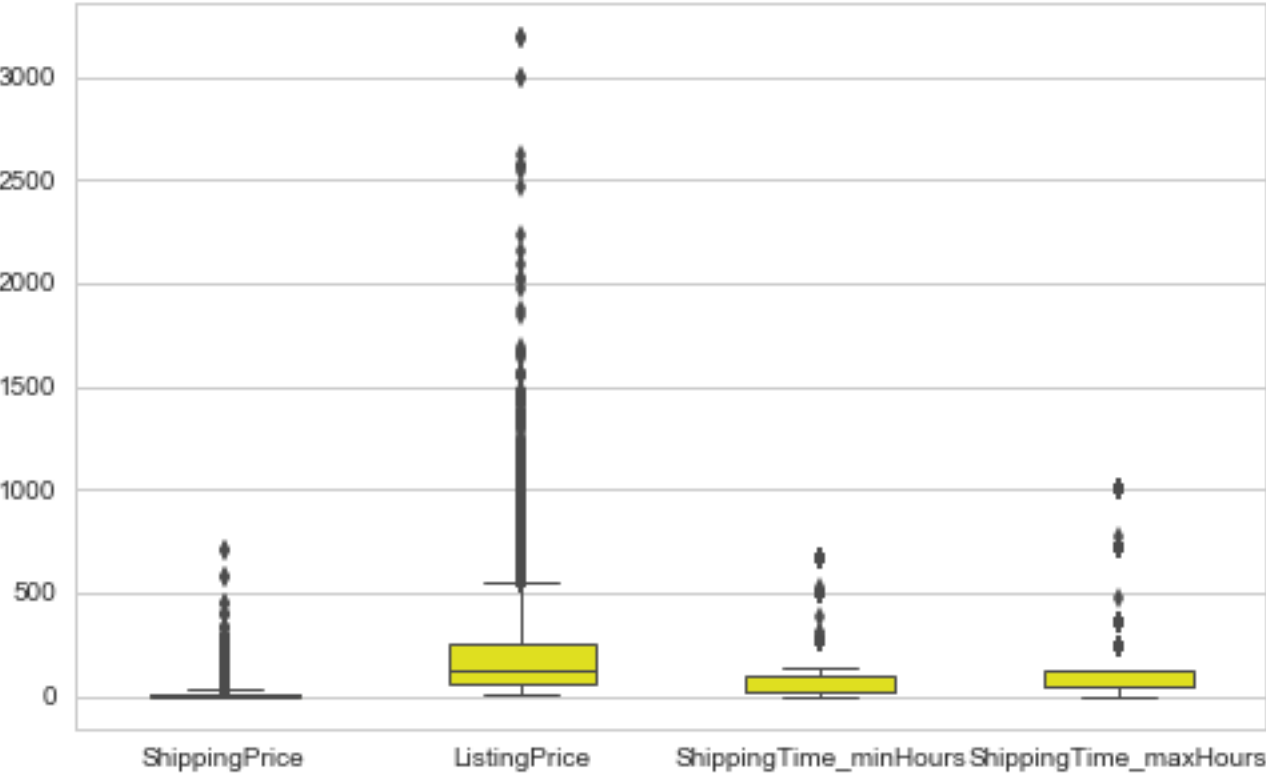


Fig 3. Box plots for continuous feature 'ShippingPrice' in more details, and box plots for the other features.

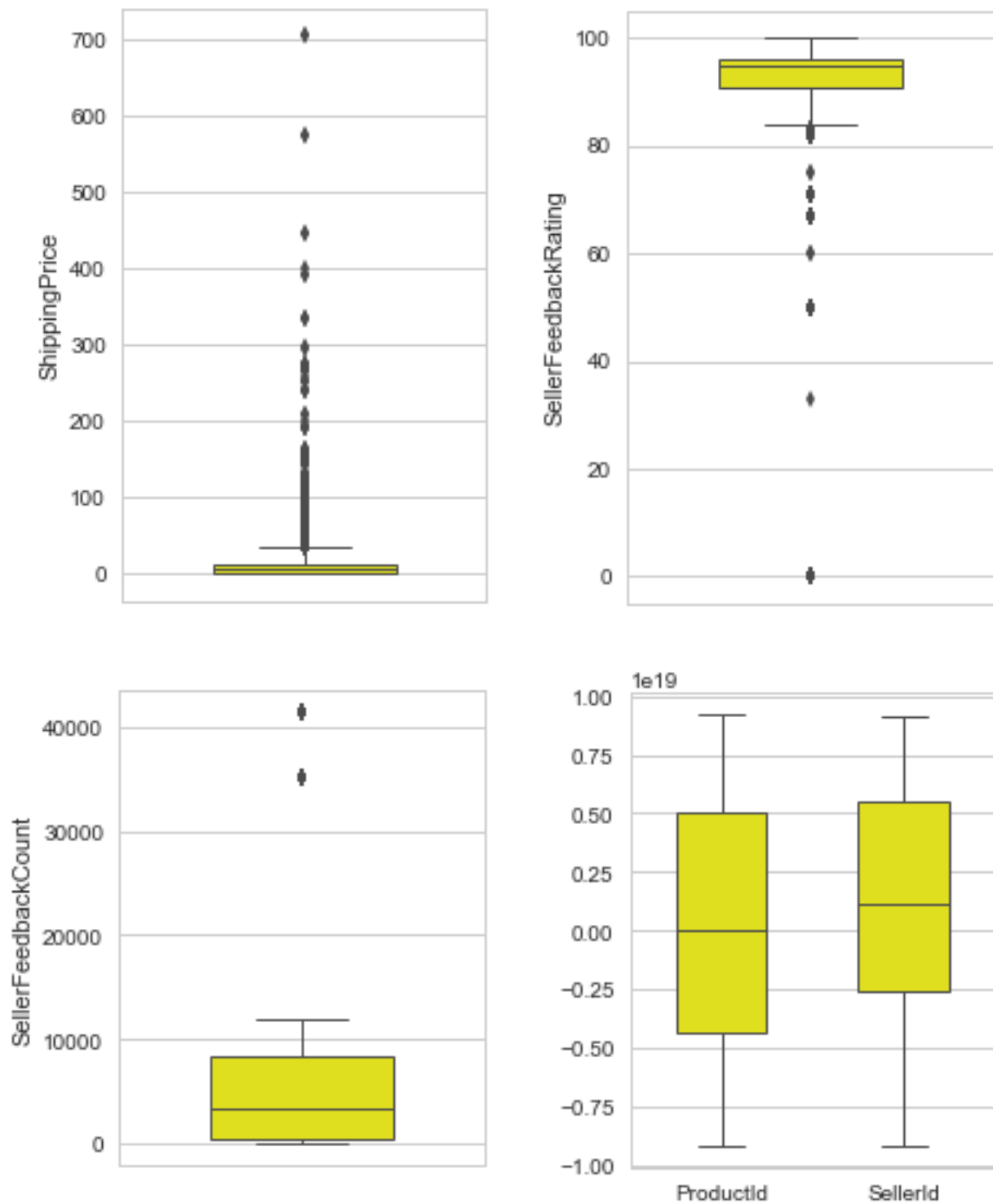


Fig 4. Bar plot for categorical feature 'IsWinner'.

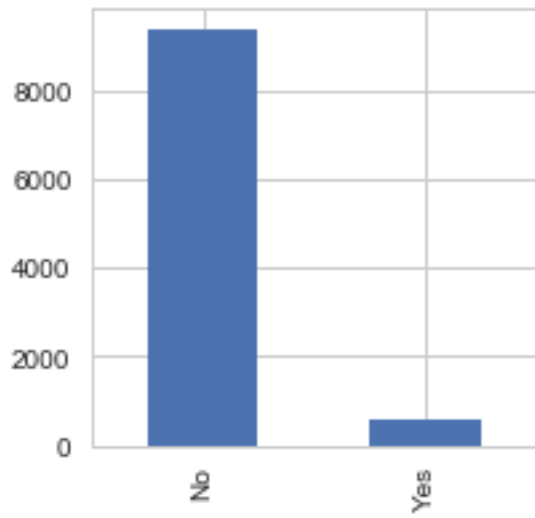


Fig 5. Bar plot for categorical feature 'IsFeaturedMerchant'.

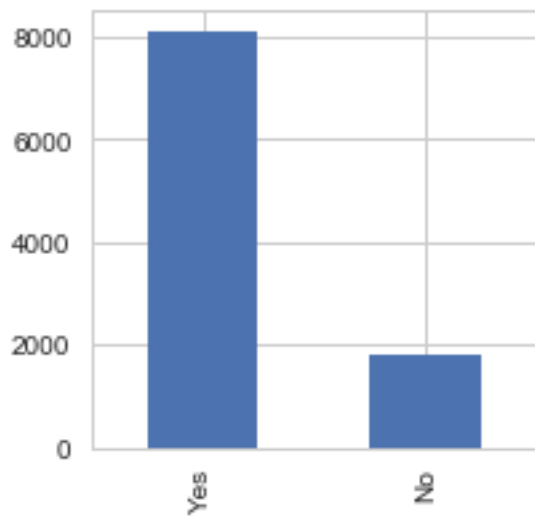


Fig 6. Bar plot for categorical feature 'IsFulfilledByAmazon'.

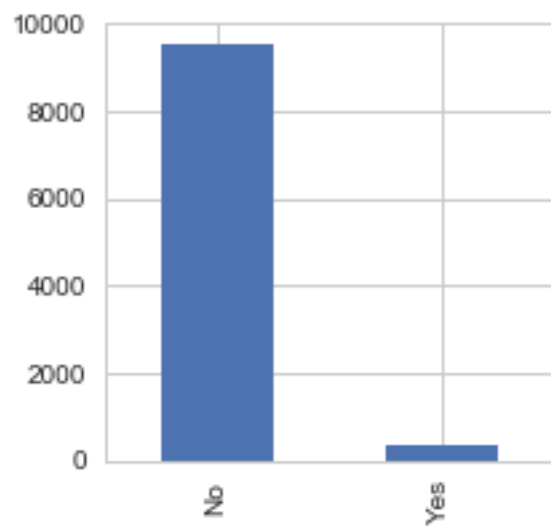
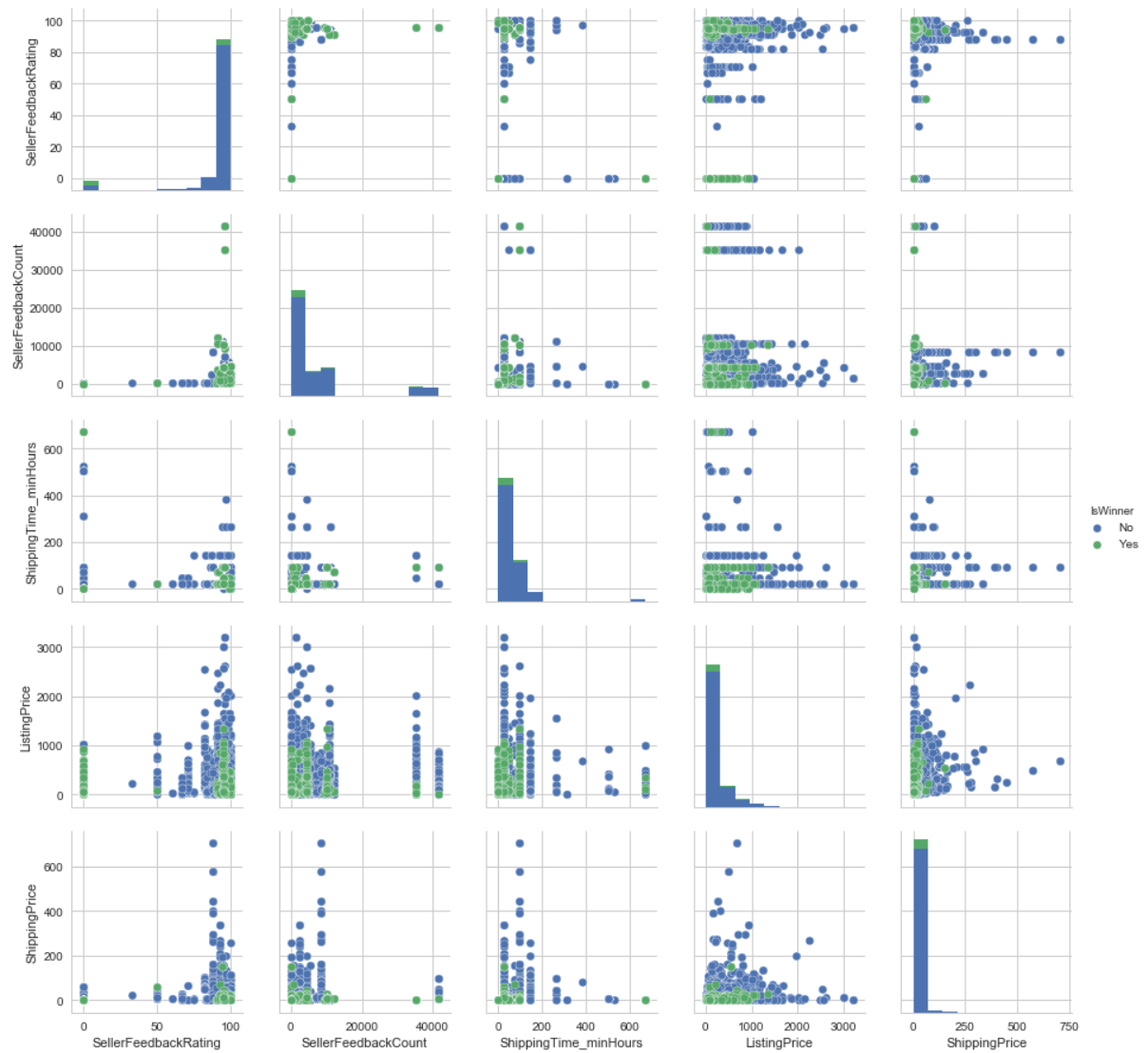


Fig 7. Scatter plot matrix for continuous features based on the target feature (Green dots for winner sellers, blue dots for non-winner sellers).



4 - CONSIDERATIONS

FROM HISTOGRAMS

- The features have different types of distributions.
- Exponential for ListingPrice and ShippingPrice.
- Skewed to the right for ShippingTime_minHours and ShippingTime_maxHours
- Skewed to the left (but almost exponential) for the SellerFeedbackRating.
- All these features could potentially have outliers, or errors, or misinterpreted missing values.
- SellerFeedbackCount could be bimodal, but not clear.
- SellerId and ProductId are both unimodal.
- SellerFeedbackRating presents a series of values at or around zero, isolated from the rest. It could indicate missing values? Or just relatively new sellers. Further inspection of the data is required. It is not clear why there is such a high difference between the data around zero and the rest of the data in the distribution, with mean at/around 90. Maybe the zero was used to indicate a missing value? If so, we could drop these values, or give them the value of the mean of the distribution, or drop altogether the feature if % of zero is too high (but from histogram seem only a small portion of the data).
- These considerations can be added to the Quality Data Plan for further evaluation later.

FROM BOX PLOTS

- We can see looking at the boxplots that many features present outliers, and we should deal with them. We take note of these in the Data Quality Plan (see file *DataQualityPlan.pdf*).
- As expected, the feature 'SellerFeedbackRating' present outliers near or at zero. We should further investigate how many of these outliers are present, then decide if to delete them, or assign to them a specific value (maybe at the lower range of the distribution, instead of using the mean). This feature should be treated with caution as it could be important in predicting the Winner sellers.
- Similar considerations should be made for the rest of the features where outliers are present, and register these on the Data Quality Plan, indicating a possible solution.

FROM BAR PLOTS

- The bar plots clearly define the two subsets in the various binary features. There are very few winner sellers in the data set, many featured merchants, and very little of the offers are fulfilled directly by Amazon.

FROM SCATTER PLOT MATRIX

- No clear correlations are seen between the descriptive features, but the division of data into winner sellers (green dots) and non-winner sellers (blue) gives some clear conclusions.
- Green dots concentrate in certain areas compared to blue dots. For instance, we can see that winner sellers have usually lower listing and shipping prices, but also higher ratings and lower minimum shipping time. No much differences were visible in rating counts.
- A good portion of values for seller ratings is at zero, but as we have seen, these are all outliers and will be dealt with in the Data Preparation Section.