

Author: Andreas Traut

Date: 08.05.2020 (Updates 24.07.2020)

[Download as PDF](#)

Machine Learning with Python

0. Introduction

- a) Aim of this repository: "Small Data" versus "Big Data"
- b) Motivation for IDEs
- c) Structure of this repository
 - (i) First part: "Movies Database" example
 - (ii) Second part: Scikit-Learn Example ("Small Data")
 - (iii) Third part: Spark Example ("Big Data")
 - (iv) Summary Mind-Map
 - (v) Digression (Excurs) to Big Data Visualization and K-Means Clustering Algorithm and Map-Reduce
- d) Future learnings and coding & data sources

I. "Movies Database" Example

- 1. Separate "NaN"-values
- 2. Draw a stratified sample
- 3. Create a pipeline
- 4. Fit the model
- 5. Cross-validation
- 6. Prediction
- 7. Conclusion

II. "Small Data" Machine Learning using "Scikit-Learn"

Initialize and Read the CSV File

Define Auxiliary Variables and Functions

Read the csv-file

Split Training Data and Test Data

Generate Stratified Sampling

Verify if Stratified sample is good

Discover and Visualize the Data to Gain Insights

Clean NULL-Values and Prepare for Machine Learning

Find all NULL-Values

Remove all NULL-Values

Model-Specific Preprocessing

Use "Imputer" to Clean NULL-Values

Treat "Categorial" Inputs

Pipelines and Custom Transformer

Custom Transformer

Pipelines

Select and Train Model

LinearRegression Model

DecisionTreeRegressor Model

Crossvalidation

For DecisionTreeRegressor

For LinearRegression

For RandomForestRegressor

Save Model

Optimize Model

GridSearchCV

GridSearchCV on RandomForestRegressor

GridSearchCV on LinearRegressor

Randomized Search

Analyze best models

Evaluate final model on test dataset

III. "Big Data" Machine Learning using the "Spark ML Library"

Initialize and Read the CSV File

Get Ready with Docker

0. Initialize Spark

 0.1 Create Spark Context and Spark Session

 0.2 Read CSV

 0.3 Dataset Properties and some Select, Group and Aggregate Methods

 0.4 Write as Parquet or CSV

 0.5 Read Parquet

1. Cleaning the data

 1.2 Replacing and Casting

 1.3 Null-Values

 1.4 String Values

2. Model-specific preprocessing

 2.0 Check missing entries and define userdefined scatter plot

 2.1 StringIndexer

 2.2 OntHotEncoder

 2.3 VectorAssembler

3. Aligning and numerating Features and Labels

 3.1 Aligning

 3.2 Numerating

4. Pipelines

5. Training data and Testing data

6. Apply models and evaluate

 6.1 Ordinary Least Square Regression

 6.2 Ridge Regression

 6.3 Lasso Regression

7. Minhash und Local-Sensitive-Hashing (LSH)

IV. Summary Mind-Map

V. Digression (Excurs) to Big Data Visualization and K-Means Clustering Algorithm and Map-Reduce

Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License

Machine Learning with Python

0. Introduction

a) Aim of this repository: "Small Data" versus "Big Data"

After having learnt visualization techniques in Python (which I showed in my repository ["Visualization-of-Data-with-Python"](#)), I started working on different datasets with the aim to learn and apply machine learning algorithms. I was particularly interested in better **understanding the differences and similarities of "Small Data" (Scikit-Learn) approaches versus the "Big Data" (Spark) approaches!**

Therefore I tried to focus more on this "comparison" question of "Small Data" coding vs "Big Data" coding instead of digging into too many details of each of these approaches. I haven't seen many comparisons of "Small Data" vs "Big Data" coding and I think understanding this is interesting and important.

b) Motivation for IDEs

I will use [Jupyter-Notebooks](#), which is a widespread standard today, but I will also use [Integrated Development Environments \(IDEs\)](#). The first Jupyter-Notebooks have been developed 5 years ago (in 2015). Since my first programming experience was more than 25 years ago (I started with [GW-Basic](#) then [Turbo-Pascal](#) and so on and I am also familiar with [MS-DOS](#)). I quickly learnt the advantages of using Jupyter-Notebooks. **But** I missed the comfort of an [IDE](#) from the very first days!

Why is it important for me to mention the IDEs out so early in a learning process? In my opinion Jupyter-Notebooks are good for the first examinations of data and for documenting procedures and up to a certain degree also for sophisticated data science. But it might be a good idea to learn very early how to work with an IDE. Think about how to use what has been developed so far later in a bigger environment (for example a [Lambda-Architecture](#), but you can take whatever other environment, which requires robustness&stability). I point this out here, because after having read several e-Books and having participated in seminars I see that IDEs are not in the focus.

Therefore: in my examples in this repository here I will also work with Python ".py" files. These ".py" can be executed in an IDE, like e.g. [Spyder-IDE](#), which can be downloaded for free and looks like this:

The screenshot shows the Spyder IDE interface. On the left, a code editor displays a Python script named 'Sklearn Machine-Learning_AirBnB.py'. The script contains code for fetching AirBnB data, reading CSV files, and performing data processing. On the right, there are several panes: a 'Variables manager' showing objects like 'division', 'print_function', and 'unicode_literals'; a 'Konsole 1/A' pane showing the execution of the script and output messages; and a 'IPython-Konsole' pane at the bottom. The status bar at the bottom indicates the current environment and file details.

c) Structure of this repository

(i) First part: "Movies Database" example

Therefore the *first example* uses a [Jupyter-Notebook](#) in order to learn the standard procedures (e.g. data-cleaning & preparing, model-training,...). I worked on data converting movies and their revenues.

(ii) Second part: Scikit-Learn Example ("Small Data")

The *second example* is for being used in an IDE (integrated developer environment), like the [Spyder-IDE](#) from the [Anaconda distribution](#) and apply the "[Scikit-Learn Python Machine Learning Library](#)" (you may call this example a "Small Data" example if you want). I will show you a typical structure for a machine-learning example and put it into a mind-map. The same structure will be applied on the third example.

(iii) Third part: Spark Example ("Big Data")

The *third example* is a "Big Data" example and will use a [Docker environment](#) and apply the ["Apache Machine Learning Library"](#), a scalable machine learning library. The mind-map from the second part will be extended and aligned to the second example.

In this example I also show some *Big Data Visualizations techniques*, show how the *K-Means Clustering Algorithm in Apache Spark ML* works and explain the *Map-Reduce* programming model on a Word-Count example.

(iv) Summary Mind-Map

I provide a summary mind-map, which possibly helps you to structure your code. There are lots of similarities between "Small Data" and "Big Data".

(v) Digression (Excurs) to Big Data Visualization and K-Means Clustering Algorithm and Map-Reduce

In this Digression (Excurs) I will provide some examples for Big Data Visualization, K-Means Clustering and Map-Reduce.

d) Future learnings and coding & data sources

For all of these topics various tutorials, documentation, coding examples and guidelines can be found in the internet **for free!** The Open Source Community is an incredible treasure trove and enrichment that positively drives many digital developments: [Scikit-Learn](#), [Apache Spark](#), [Spyder](#), [GitHub](#), [Tensorflow](#) and also [Firefox](#), [Signal](#), [Threema](#), [Corona-Warnapp](#)... to be mentioned. There are many positive examples of sharing code and data "for free".

Coding:

If you Google for example "*how to prepare and clean the data with spark*", you will find tons of documents around "*removing null values*" or "*encoders*" (like the OneHotEncoder for treating categorical inputs) or "*pipelines*" (for putting all the steps in an efficient, customizable order) so on. You will be overwhelmed of all this. Some resources to mention are the [official documentation](#) and a few more Github repositories like e.g. [tirthajyoti/Spark-with-Python](#) (MIT licence), [Apress/learn-pyspark](#) (Freeware License), [mahmoudparsian/pyspark-tutorial](#) (Apache License v2.0). What I will do here in my repository is nothing more than putting it together so that it works for my problem (which can be challenging as well sometimes). Adapting it for your needs should be easier from this point on.

Data:

If you would like to do further analysis or produce alternate visualizations of the Airbnb-data, you can download them from [here](#). It is available below under a [Creative Commons 1.0 Universal Public Domain Dedication](#) license. The data for the Vermont-Vendor-Payments can be downloaded from [here](#) and are available under the [Open Data Commons Open Database License](#). The movies database doesn't even mention a license and is from [Kaggle](#). There you find a lot of more datasets and also coding examples for your studies.

I. "Movies Database" Example

A good starting point for finding useful datasets is "Kaggle" ([www.kaggle.com](#)). I downloaded the movies dataset from [here](#). The dataset from Kaggle contains the following columns:

Rank | Title | Year | Score | Metascore | Genre | Vote | Director | Runtime | **Revenue** | Description | RevCat

In this example I want to predict the "**Revenue**" based on the other information, which I have for each movie (e.g. every movie has a year, a scoring, a title ...). There are some "NaN"-values in the column "Revenue" and instead of filling them with an assumption (e.g. median-value) as I did in another Jupiter-Notebook (see [here](#)), I wanted to predict these values. You might guess the conclusion already: predicting the revenue based on the available information as shown above (the columns) might not work. But essential to me is more to follow a well established standard-process of data-cleaning, data-preparing, model-training and error-calculation in this example in order to learn how to apply this process to better datasets, than the movies-dataset, later.

Therefore, here is how I approached the problem step-by-step:

1. Separate "NaN"-values

I separated the rows with "NaN"-values in column "Revenue"

These are the datarows, where column "Revenue" is null:

	Rank	Title	Year	Score	Metascore	Genre	Vote	Director	Runtime	Revenue	Description
	82	A Clockwork Orange	1971	8.3	80.0	Crime, Drama, Sci-Fi	662768	Stanley Kubrick	136	NaN	In the future, a sadistic gang leader is impr...
	513	To Kill a Mockingbird	1962	8.3	87.0	Crime, Drama	262064	Robert Mulligan	129	NaN	Atticus Finch, a lawyer in the Depression-era ...
	581	Death Proof	2007	7.0	NaN	Action, Thriller	236539	Quentin Tarantino	113	NaN	Two separate sets of voluptuous women are stal...
	620	My Neighbour Totoro	1988	8.2	86.0	Animation, Family, Fantasy	226126	Hayao Miyazaki	86	NaN	When two girls move to the country to be near ...
	685	Hachi: A Dog's Tale	2009	8.1	NaN	Drama, Family	212349	Lasse Hallström	93	NaN	A college professor's bond with the abandoned ...

```
In [10]: movies_RevenueNaN = movies[movies["Revenue"].isnull()]
movies_RevenueNaN.head()

Out[10]:
Rank      Title   Year  Score  Metascore    Genre   Vote  Director Runtime  Revenue
82  A Clockwork Orange  1971  8.3  80.0  Crime, Drama, Sci-Fi  662768  Stanley Kubrick  136       NaN  In the future, a sadistic gang leader is impr...
513  To Kill a Mockingbird  1962  8.3  87.0  Crime, Drama  262064  Robert Mulligan  129       NaN  Atticus Finch, a lawyer in the Depression-era ...
581  Death Proof  2007  7.0       NaN  Action, Thriller  236539  Quentin Tarantino  113       NaN  Two separate sets of voluptuous women are stal...
620  My Neighbour Totoro  1988  8.2  86.0  Animation, Family, Fantasy  226126  Hayao Miyazaki  86       NaN  When two girls move to the country to be near ...
685  Hachi: A Dog's Tale  2009  8.1       NaN  Drama, Family  212349  Lasse Hallström  93       NaN  A college professor's bond with the abandoned ...

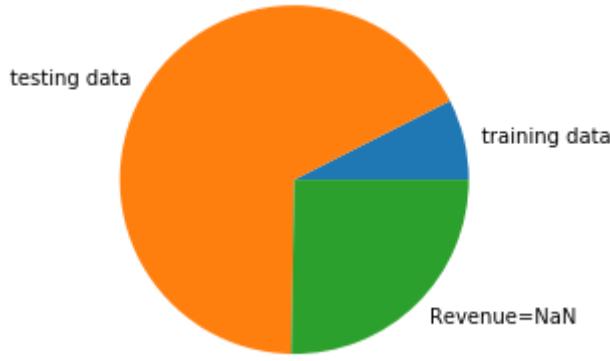
In [11]: len_movies_RevenueNaN = len(movies_RevenueNaN)
len_movies_RevenueNaN

Out[11]: 2527
```

2. Draw a stratified sample

I drew a stratified sample (based on "Revenue") on this remaining dataset and I received a training dataset and testing dataset:

In [24]:	split = StratifiedShuffleSplit(n_splits=1, test_size=0.9, random_state=42)
	for train_index, test_index in split.split(movies_NotNullC, movies_NotNullC["RevCat"]):
	strat_train_set = movies_NotNullC.iloc[train_index]
	strat_test_set = movies_NotNullC.iloc[test_index]
In [25]:	strat_test_set["RevCat"].value_counts() / len(strat_test_set)
Out[25]:	1 0.903657 2 0.069878 3 0.016057 4 0.010407 Name: RevCat, dtype: float64
In [26]:	movies_NotNullC["RevCat"].value_counts() / len(movies_NotNullC)
Out[26]:	1 0.903653 2 0.069851 3 0.016058 4 0.010438 Name: RevCat, dtype: float64



3. Create a pipeline

I created a pipeline to fill the "NaN"-value in other columns (e.g. "Metascore", "Score").

```
In [38]: num_pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy='median')),
])

num_attribs = ['Rank', 'Year', 'Score', 'Metascore', 'Vote', 'Runtime']

full_pipeline = ColumnTransformer([
    ('num', num_pipeline, num_attribs)
])
```

Now apply the Pipeline:

```
In [40]: movies_train_prepared = full_pipeline.fit_transform(movies_train)
```

Now let's count the "nan" values in the new prepared dataset "movies_train_prepared". I have to transform it back to a Pandas-Dataframe format first:

```
In [41]: tmp_num = movies_train.select_dtypes(include=[np.number])
tmp_prep = pd.DataFrame(movies_train_prepared, columns=tmp_num.columns, index=movies_train.index)
tmp = tmp_prep[tmp_prep["Metascore"].isnull()]
tmp
```

```
Out[41]:
```

Rank	Year	Score	Metascore	Vote	Runtime
------	------	-------	-----------	------	---------

Zero, as we wanted! All "nan"-values in "movies_train_prepared" have been removed by the "median" value (this was how the pipeline was built). Great,

4. Fit the model

I used the training dataset and fitted it with the "DecisionTreeRegressor" model

```
In [42]: tree_reg = DecisionTreeRegressor(random_state=42)
tree_reg.fit(movies_train_prepared, movies_train_labels)

movies_predictions = tree_reg.predict(movies_train_prepared)
```

5. Cross-validation

I verified with a cross-validation, how good this model/parameters are

```
In [46]: scores = cross_val_score(tree_reg,
                                 movies_train_prepared,
                                 movies_train_labels,
                                 scoring="neg_mean_squared_error",
                                 cv=10)

rmse_scores = np.sqrt(-scores)

def display_scores(scores):
    print("Scores:", scores)
    print("Mean:", scores.mean())
    print("Standard deviation:", scores.std())

display_scores(rmse_scores)

Scores: [ 69.66166984  62.59977724  62.6450061  112.98391756  47.99491086
         52.16787811  42.65104005  83.21059338  41.09199207  63.84158945]
Mean: 63.88483746577791
Standard deviation: 20.447326452253154
```

6. Prediction

I did a prediction on a subset of the testing dataset and did a side-by-side comparison of prediction and true value

```
In [53]: side_by_side = [(true, pred) for true, pred in zip(list(some_data_label), list(some_data_predictions))]
side_by_side

Out[53]: [(5.78, 10.91),
(0.05, 0.44),
(0.3, 2.68),
(159.6, 11.99),
(33.63, 2.19),
(44.9, 26.83),
(38.52, 22.52),
(33.04, 2.19),
(3.2, 11.99),
(37.49, 16.38),
(17.88, 64.19),
(41.19, 191.45),
(16.19, 12.19),
(0.05, 3.02),
(16.68, 64.19),
(35.11, 11.54),
(36.0, 40.22),
(3.61, 19.64),
(0.59, 0.05),
(0.99, 5.48)]
```

I performed a prediction on the testing dataset and calculated the mean-squared error

```
In [55]: movies_test_prepared = full_pipeline.fit_transform(movies_test)
movies_test_predictions = tree_reg.predict(movies_test_prepared)
lin_mse = mean_squared_error(movies_test_labels, movies_test_predictions)
lin_rmse = np.sqrt(lin_mse)
lin_rmse

Out[55]: 54.68522284077671
```

7. Conclusion

The conclusion of this machine learning example is obvious: it is rather not possible to predict the "Revenue" based on the available information (the most useful numerical features were "year", "score", ... and the other categorical like "genre" don't seem to have much more added value in my opinion).

Please find the complete Jupyter Notebook here:

<https://github.com/AndreasTraut/Machine-Learning-with-Python/blob/master/Movies%20Machine%20Learning%20-%20Predict%20NaNs.ipynb>

If you want to run the code immediately without installing the required "Jupyter environment" then you can use this Deepnote-Link:

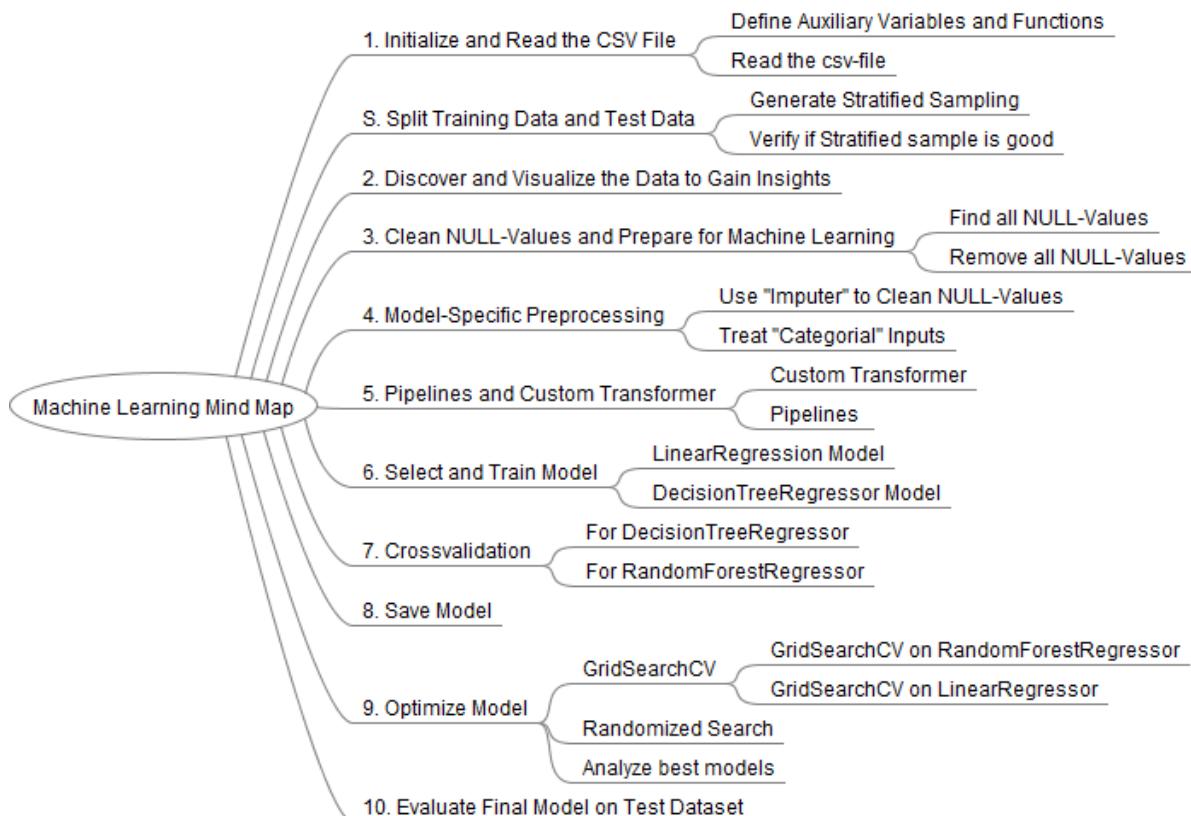
<https://beta.deepnote.com/project/754094f0-3c01-4c29-b2f3-e07f507da460>

II. "Small Data" Machine Learning using "Scikit-Learn"

In my opinion Jupyter Notebooks are **not** always the best environment for learning to code! I agree, that Jupyter Notebooks are nice for doing documentation of python code. It really looks beautiful. But I prefer debugging in an IDE instead of a Jupyter Notebook: having the possibility to set a breakpoint can be a pleasure for my nerves, specially if you have longer programs. Some of my longer Jupyter Notebooks feel from the hundreds line of code onwards more like pain than like anything helpful. And I also prefer having a "help window" or a "variable explorer", which is smoothly integrated into the IDE user interface. And there are a lot more advantages why getting familiar with an IDE is a big advantage compared to the very popular Jupyter Notebooks! I am very surprised, that everyone is talking about Jupyter Notebooks but IDEs are only mentioned very seldom. But maybe my preferences are also a bit different, because I grew up in a [MS-DOS](#) environment. :-)

I choose in this *second example* the [Spyder-IDE](#) and worked on "[Scikit-Learn](#)", a very popular python machine learning library. The structure of the Python code is a bit similar to the steps, which I followed in the Movies Database example above (you will find these sections also in the ".py" file).

So let's start with the "scikit-learn" ("SmallData", if you want). I will align this structure to the Spark "Big Data" mind map below in order to learn from each of this two approaches.



I aligned this "Small Data" structure to the Apache Spark "Big Data" structure in order to learn from each of this two approaches. Finally I will put these two Mind Maps into one big which you can take as a guide to navigate through all of your machine-learning problems.

Initialize and Read the CSV File

Define Auxiliary Variables and Functions

These auxiliary variables and auxiliary functions are intended to make your work easier:

```
PROJECT_ROOT_DIR = "."
myDataset_NAME = "AirBnB"
IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "media")
myDataset_PATH = os.path.join("datasets", "AirBnB")

def save_fig(fig_id, prefix=myDataset_NAME,
             tight_layout=True, fig_extension="png", resolution=300):
    path = os.path.join(IMAGES_PATH, prefix + "_" + fig_id + "." +
fig_extension)
    print("Saving figure", prefix + "_" + fig_id)
    if tight_layout:
        plt.tight_layout()
    plt.savefig(path, format=fig_extension, dpi=resolution)
```

Read the csv-file

Define a function for reading the CSV file:

```
def load_myDataset_data(myDataset_path=myDataset_PATH):
    csv_path = os.path.join(myDataset_path, "listings.csv")
    return pd.read_csv(csv_path)

myDataset = load_myDataset_data()
print(myDataset.head())
```

The dataset has the following format (use `myDataset.info()`):

```
RangeIndex: 25164 entries, 0 to 25163
Data columns (total 15 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   name            25114 non-null   object 
 1   host_id          25164 non-null   int64  
 2   host_name        25142 non-null   object 
 3   neighbourhood_group  25164 non-null   object 
 4   neighbourhood     25164 non-null   object 
 5   latitude          25164 non-null   float64
 6   longitude         25164 non-null   float64
 7   room_type         25164 non-null   object 
 8   price             25164 non-null   int64  
 9   minimum_nights    25164 non-null   int64  
 10  number_of_reviews 25164 non-null   int64  
 11  last_review       20636 non-null   object 
 12  reviews_per_month 20636 non-null   float64
 13  calculated_host_listings_count 25164 non-null   int64  
 14  availability_365  25164 non-null   int64 
```

```
dtypes: float64(3), int64(6), object(6)
```

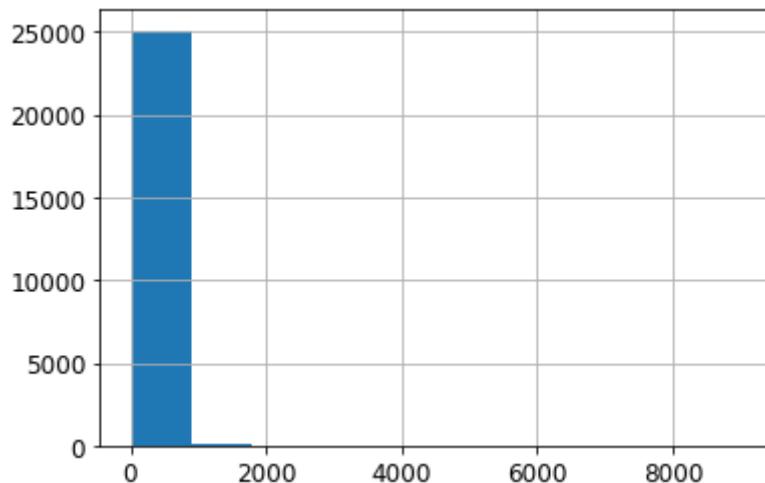
Split Training Data and Test Data

The aim is to predict the `price` and use the other columns (or some of them) as features.

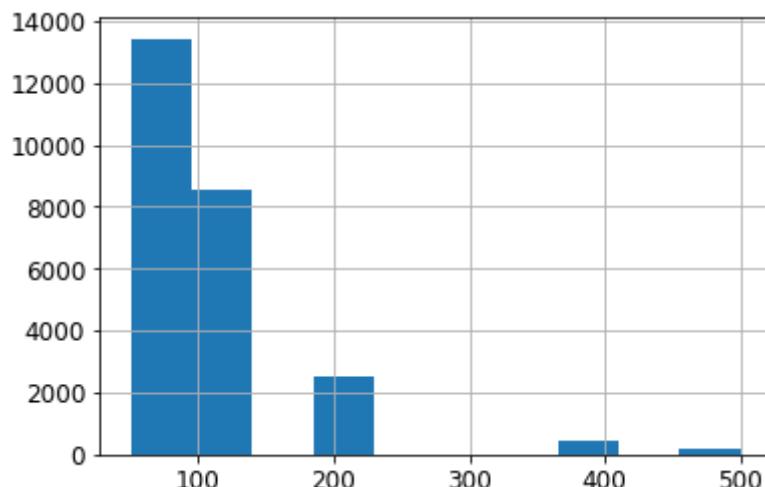
Generate Stratified Sampling

I want to create the training dataset and the test dataset by applying stratified sampling:

```
myDataset["price"].hist()  
myDataset["price_cat"] = pd.cut(myDataset["price"],  
                               bins=[-1, 50, 100, 200, 400, np.inf],  
                               labels=[50, 100, 200, 400, 500])  
print("\nvalue_counts\n", myDataset["price_cat"].value_counts())  
myDataset["price_cat"].hist()  
  
split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)  
for train_index, test_index in split.split(myDataset, myDataset["price_cat"]):  
    strat_train_set = myDataset.loc[train_index]  
    strat_test_set = myDataset.loc[test_index]
```



There is a very long tail in the price (going to about 9000).



Verify if Stratified sample is good

Next step is to verify if the stratified sample is good and compare it to a random sampling

```
def price_cat_proportions(data):
    return data["price_cat"].value_counts() / len(data)

train_set, test_set = train_test_split(myDataset, test_size=0.2,
random_state=42)

compare_props = pd.DataFrame({
    "Overall": price_cat_proportions(myDataset),
    "Stratified": price_cat_proportions(strat_test_set),
    "Random": price_cat_proportions(test_set),
}).sort_index()

compare_props["Rand. %error"] = 100 * compare_props["Random"] /
                                compare_props["Overall"] - 100
compare_props["Strat. %error"] = 100 * compare_props["Stratified"] /
                                compare_props["Overall"] - 100
```

From the `%error` columns I can see, that the stratified sample is always better, than the random sampling.

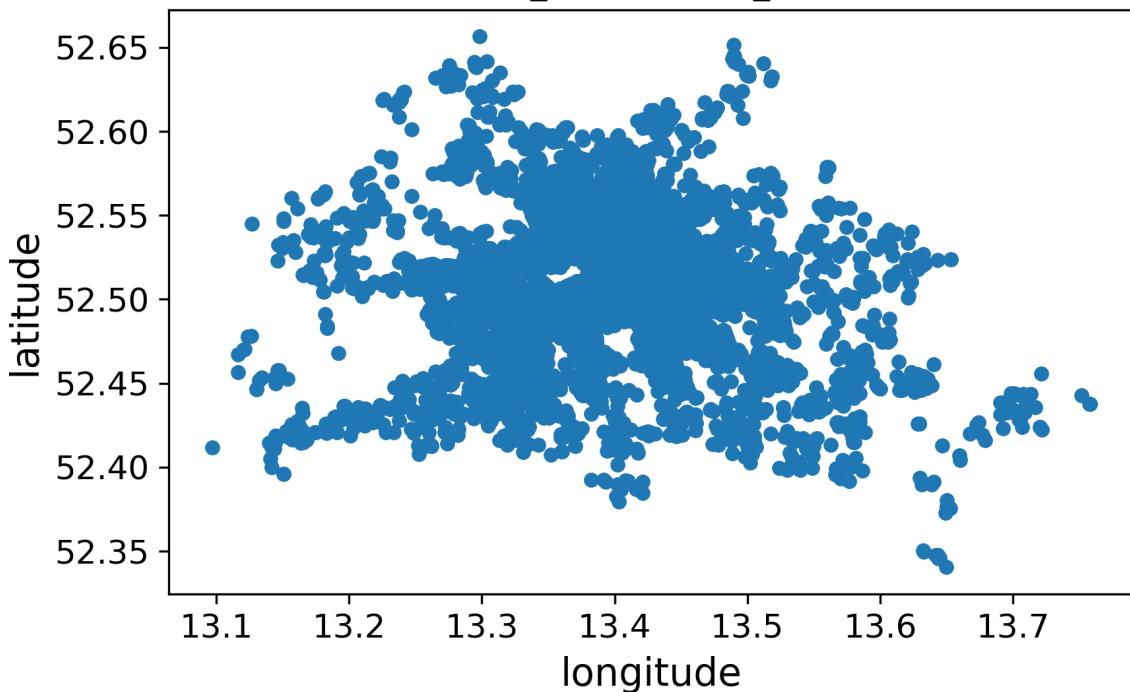
	Overall	Stratified	Random	Rand. %error	Strat. %error
50	0.533659	0.533678	0.536062	0.450249	0.003473
100	0.340168	0.340155	0.343731	1.047386	-0.003974
200	0.101256	0.101331	0.097755	-3.457526	0.074516
400	0.017326	0.017286	0.015498	-10.554013	-0.233322
500	0.007590	0.007550	0.006954	-8.380604	-0.527513

Discover and Visualize the Data to Gain Insights

Now it is time to create some visualizations as I described in [@sec:Visualization]. The first plot is a bad example, because the big blue dots don't reveal any interesting information about the `price` (which was our prediction variable).

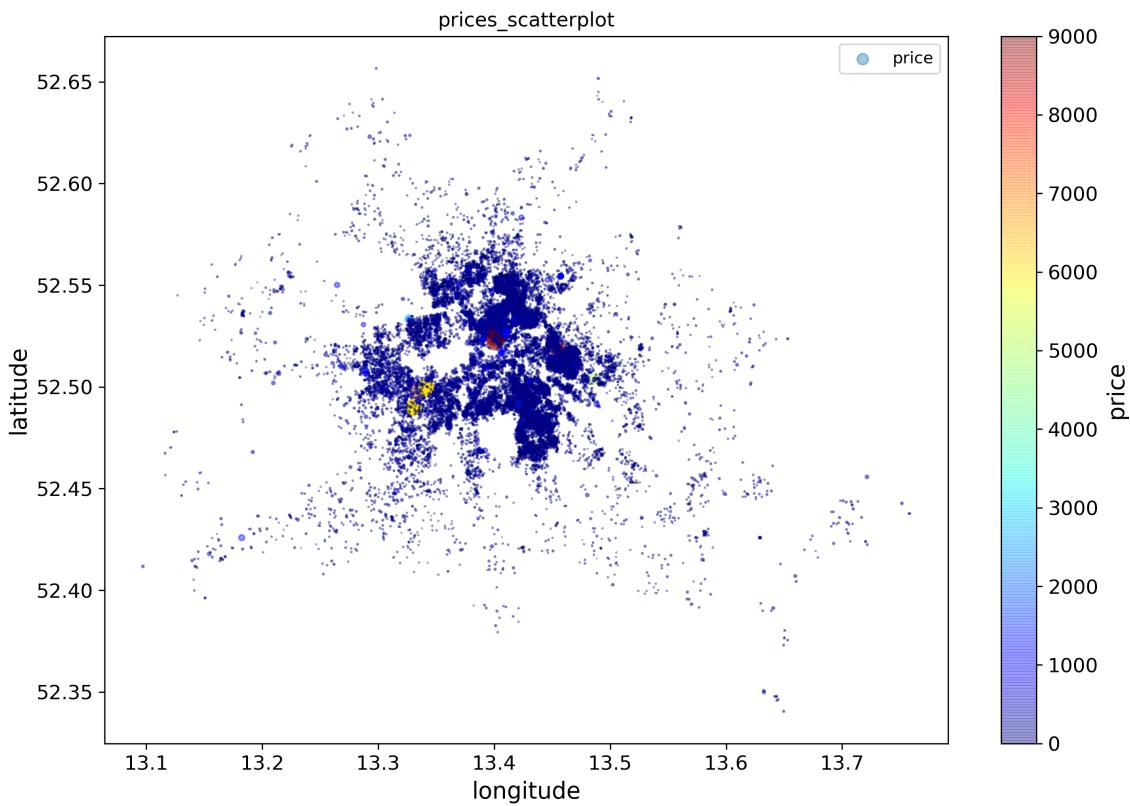
```
myDataset.plot(kind="scatter",
               x="longitude", y="latitude",
               title="bad_visualization_plot")
save_fig("bad_visualization_plot")
```

bad_visualization_plot

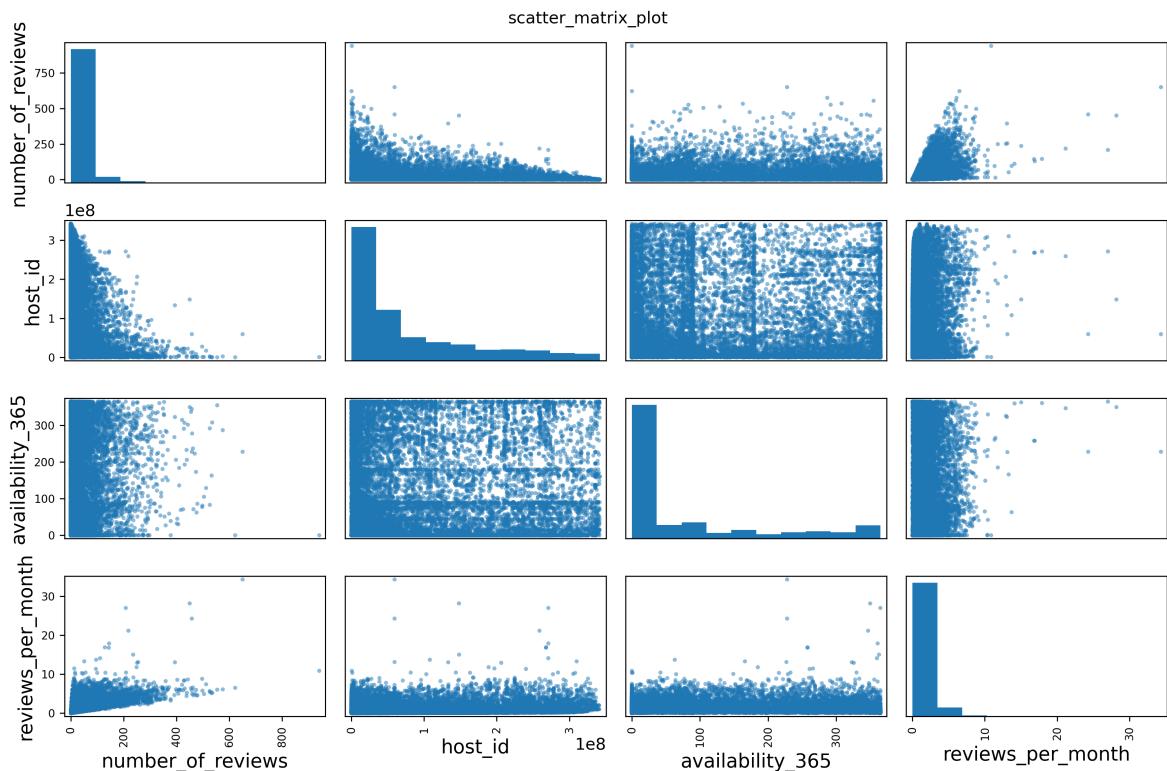


Better would be to include the price by a heat color:

```
myDataset.plot(kind="scatter",
               x="longitude", y="latitude",
               alpha=0.4, s=myDataset["price"]/100,
               label="price", figsize=(10,7),
               c="price", cmap=plt.get_cmap("jet"),
               colorbar=True, sharex=False,
               title="prices_scatterplot")
plt.legend()
save_fig("prices_scatterplot")
```



```
attributes = ["number_of_reviews", "host_id",
              "availability_365", "reviews_per_month"]
scatter_matrix(myDataset[attributes],
               figsize=(12, 8))
plt.suptitle("scatter_matrix_plot")
save_fig("scatter_matrix_plot")
```



As I am interested in predicting the `price` I will calculate the correlation matrix in order to see, which column is most important:

```
corr_matrix = myDataset.corr()
print("correlation:\n", corr_matrix["price"].sort_values(ascending=False))
```

```
correlation:
price                      1.000000
availability_365            0.096979
calculated_host_listings_count 0.077545
host_id                      0.045434
reviews_per_month             0.034332
latitude                     0.007836
number_of_reviews              0.000611
minimum_nights                 -0.006361
longitude                     -0.036490
Name: price, dtype: float64
```

Clean NULL-Values and Prepare for Machine Learning

Find all NULL-Values

```
print("\nHow many Non-NULL rows are there?\n")
print(myDataset.info())
print("\nAre there NULL values in the columns?\n",
      myDataset.isnull().any())
print("\nAre there NULLs in column reviews_per_month?\n",
      myDataset["reviews_per_month"].isnull().any())
print("\nShow some rows with NULL (head only):\n",
      myDataset[myDataset["reviews_per_month"].isnull()].head())
```

How many Non-Null rows are there?

```
Data columns (total 14 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   name             20088 non-null   object 
 1   host_id          20131 non-null   int64  
 2   host_name        20114 non-null   object 
 3   neighbourhood_group 20131 non-null   object 
 4   neighbourhood     20131 non-null   object 
 5   latitude         20131 non-null   float64
 6   longitude        20131 non-null   float64
 7   room_type        20131 non-null   object 
 8   minimum_nights   20131 non-null   int64  
 9   number_of_reviews 20131 non-null   int64  
 10  last_review       16501 non-null   object 
 11  reviews_per_month 16501 non-null   float64
 12  calculated_host_listings_count 20131 non-null   int64  
 13  availability_365 20131 non-null   int64  
dtypes: float64(3), int64(5), object(6)
```

Are there NaN values in the columns?

name	True
host_id	False
host_name	True

neighbourhood_group	False
neighbourhood	False
latitude	False
longitude	False
room_type	False
minimum_nights	False
number_of_reviews	False
last_review	True
reviews_per_month	True
calculated_host_listings_count	False
availability_365	False

Remove all NULL-Values

There are different options for handling NULL values, which occur in a column:

Option 1: we can delete the entire row

Option 2: we can delete the entire column

Option 3: we can fill these with an assumption, like for example the median (which is often a good assumption for filling NULL values, but not always, as we have seen in the movies database example).

```
sample_incomplete_rows = myDataset[myDataset.isnull().any(axis=1)]

# option 1: remove rows which contains NaNs
# sample_incomplete_rows.dropna(subset=["total_bedrooms"])

# option 2 : remove columns with contain NaNs
# sample_incomplete_rows.drop("total_bedrooms", axis=1)

# option 3 : replace Nan by median
median = myDataset["reviews_per_month"].median()
sample_incomplete_rows["reviews_per_month"].fillna(median, inplace=True)

print("sample_incomplete_rows\n",
sample_incomplete_rows['reviews_per_month'].head())
```

```
sample_incomplete_rows
24411    0.43
15700    0.43
9311     0.43
21882    0.43
24259    0.43
```

Model-Specific Preprocessing

Use "Imputer" to Clean NULL-Values

Remove all text attributes because median can only be calculated on numerical attributes

```

imputer = SimpleImputer(strategy="median")
myDataset_num = myDataset.select_dtypes(include=[np.number]) #or: myDataset_num
= myDataset.drop('ocean_proximity', axis=1)
imputer.fit(myDataset_num)
print("\n imputer.strategy\n",
      imputer.strategy)
print("\n imputer.statistics_\n",
      imputer.statistics_)
print("\n myDataset_num.median\n",
      myDataset_num.median().values)
print("\n myDataset_num.mean\n",
      myDataset_num.mean().values)

```

```

imputer.strategy
median

imputer.statistics_
[3.9804571e+07 5.2509580e+01 1.3416280e+01 3.0000000e+00 5.0000000e+00
 4.3000000e-01 1.0000000e+00 0.0000000e+00]

myDataset_num.median
[3.9804571e+07 5.2509580e+01 1.3416280e+01 3.0000000e+00 5.0000000e+00
 4.3000000e-01 1.0000000e+00 0.0000000e+00]

myDataset_num.mean
[7.74925763e+07 5.25100259e+01 1.34059193e+01 7.20863345e+00
 2.17043863e+01 1.02009696e+00 2.43967016e+00 7.33118573e+01]

```

Transform the training set:

```

X = imputer.transform(myDataset_num)
myDataset_tr = pd.DataFrame(X, columns=myDataset_num.columns,
                           index=myDataset.index)
myDataset_tr.loc[sample_incomplete_rows.index.values]

```

Treat "Categorial" Inputs

Use the OneHotEncoder for the categorial values:

```

myDataset_cat = myDataset[['room_type']]
print("myDataset_cat.head\n", myDataset_cat.head(10), "\n")

cat_encoder = OneHotEncoder()
myDataset_cat_1hot = cat_encoder.fit_transform(myDataset_cat)
print("cat_encoder.categories_:\n", cat_encoder.categories_)
print("myDataset_cat_1hot.toarray():\n", myDataset_cat_1hot.toarray())
print("myDataset_cat_1hot:\n", myDataset_cat_1hot)

```

```

myDataset_cat.head
      room_type
5177    Private room
19616   Entire home/apt
10275    Private room
2078    Entire home/apt

```

```

23131 Entire home/apt
8408 Entire home/apt
17918 Entire home/apt
24411 Entire home/apt
8872 Private room
16064 Hotel room

cat_encoder.categories_:
[array(['Entire home/apt', 'Hotel room', 'Private room', 'Shared room'],
      dtype=object)]

myDataset_cat_1hot.toarray():
[[0. 0. 1. 0.]
 [1. 0. 0. 0.]
 [0. 0. 1. 0.]
 ...
 [0. 0. 1. 0.]
 [0. 0. 1. 0.]
 [0. 0. 1. 0.]]

```

Pipelines and Custom Transformer

Custom Transformer

```

print("myDataset.columns\n", myDataset_num.columns)
number_of_reviews_ix, availability_365_ix,
calculated_host_listings_count_ix,
reviews_per_month_ix = [
    list(myDataset_num.columns).index(col)
    for col in ("number_of_reviews", "availability_365",
                "calculated_host_listings_count",
                "reviews_per_month")]

def add_extra_features(X):
    number_reviews_dot_reviews_per_month =
        X[:, number_of_reviews_ix] * X[:, reviews_per_month_ix]
    return np.c_[X, number_reviews_dot_reviews_per_month]

attr_adder = FunctionTransformer(add_extra_features, validate=False)
myDataset_extra_attribs = attr_adder.fit_transform(myDataset_num.values)

myDataset_extra_attribs = pd.DataFrame(
    myDataset_extra_attribs,
    columns=list(myDataset_num.columns) +
    ["number_reviews_dot_reviews_per_month"],
    index=myDataset_num.index)
print("myDataset_extra_attribs.head()\n", myDataset_extra_attribs.head())

```

Pipelines

```

num_pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy="median")),
    ('attribs_adder', FunctionTransformer(add_extra_features,
                                         validate=False)),
    ('std_scaler', StandardScaler())])
myDataset_num_tr = num_pipeline.fit_transform(myDataset_num)

```

```

print("myDataset_num_tr\n", myDataset_num_tr)

num_attribs = list(myDataset_num)
cat_attribs = ["room_type"]

full_pipeline = ColumnTransformer([
    ("num", num_pipeline, num_attribs),
    ("cat", OneHotEncoder(), cat_attribs),
])
myDataset_prepared = full_pipeline.fit_transform(myDataset)
print("myDataset_prepared\n", myDataset_prepared)

```

Select and Train Model

LinearRegression Model

```

lin_reg = LinearRegression()
lin_reg.fit(myDataset_prepared, myDataset_labels)
some_data = myDataset.iloc[:10]
some_labels = myDataset_labels.iloc[:10]
some_data_prepared = full_pipeline.transform(some_data)
print("Predictions:\n", lin_reg.predict(some_data_prepared))
print("Labels:\n", list(some_labels))

myDataset_predictions = lin_reg.predict(myDataset_prepared)
lin_mse = mean_squared_error(myDataset_labels, myDataset_predictions)
lin_rmse = np.sqrt(lin_mse)
print("lin_rmse\n", lin_rmse)

print("mean of labels:\n", myDataset_labels.mean())
print("std deviation of labels:\n", myDataset_labels.std())

```

```

Predictions:
[ 40.29259797  89.54082186  45.36959079 101.64252692  82.68120391
 89.57892837  76.33903757  76.42930129  38.55979179 861.9719368 ]
Labels:
[41, 190, 35, 50, 100, 60, 69, 80, 32, 140]
lin_rmse
213.51224401460684
mean of labels:
74.19313496597287
std deviation of labels:
227.66240520718222

```

DecisionTreeRegressor Model

```

tree_reg = DecisionTreeRegressor(random_state=42)
tree_reg.fit(myDataset_prepared, myDataset_labels)
myDataset_predictions = tree_reg.predict(myDataset_prepared)

tree_mse = mean_squared_error(myDataset_labels, myDataset_predictions)
tree_rmse = np.sqrt(tree_mse)
print("tree_rmse\n", tree_rmse)

```

```
tree_rmse  
2.5907022436676304
```

Crossvalidation

For DecisionTreeRegressor

```
scores = cross_val_score(tree_reg, myDataset_prepared,  
                         myDataset_labels,  
                         scoring="neg_mean_squared_error",  
                         cv=10)  
tree_rmse_scores = np.sqrt(-scores)  
def display_scores(scores):  
    print("Scores:", scores)  
    print("Mean:", scores.mean())  
    print("Standard deviation:", scores.std())  
  
display_scores(tree_rmse_scores)
```

```
Scores: [137.8762754 312.21141325 218.18522714 237.35937087 72.91732588  
       65.62821113 304.05961303 86.65346214 178.337086 207.20918128]  
Mean: 182.0437166129294  
Standard deviation: 85.6479894222897
```

For LinearRegression

```
lin_scores = cross_val_score(lin_reg, myDataset_prepared,  
                            myDataset_labels,  
                            scoring="neg_mean_squared_error",  
                            cv=10)  
lin_rmse_scores = np.sqrt(-lin_scores)  
display_scores(lin_rmse_scores)
```

```
Scores: [216.69779596 285.34262945 264.35619589 299.03183929 241.53524149  
       78.73412537 208.49809202 186.28646055 147.62364622 91.09183577]  
Mean: 201.91978620154234  
Standard deviation: 72.64261021086485
```

For RandomForestRegressor

```
forest_reg = RandomForestRegressor(n_estimators=10, random_state=42)  
forest_reg.fit(myDataset_prepared, myDataset_labels)  
  
myDataset_predictions = forest_reg.predict(myDataset_prepared)  
forest_mse = mean_squared_error(myDataset_labels, myDataset_predictions)  
forest_rmse = np.sqrt(forest_mse)  
print("forest_rmse\n", forest_rmse)  
  
forest_scores = cross_val_score(forest_reg, myDataset_prepared,  
                                myDataset_labels,  
                                scoring="neg_mean_squared_error",  
                                cv=10)  
forest_rmse_scores = np.sqrt(-forest_scores)
```

```
display_scores(forest_rmse_scores)
```

```
forest_rmse
59.973965336421536
Scores: [157.89997398 130.24045184 208.33415063 235.2803762 93.2680475
51.26602041 138.75095351 194.70999772 114.94817824 152.59500511]
Mean: 147.72931551364474
Standard deviation: 52.34950554724271
```

Save Model

```
joblib.dump(forest_reg, "forest_reg.pkl")
# and later...
my_model_loaded = joblib.load("forest_reg.pkl")
```

Optimize Model

GridSearchCV

GridSearchCV on RandomForestRegressor

```
param_grid = [
    {'n_estimators': [30, 40, 50], 'max_features': [2, 4, 6, 8, 10]},
    {'bootstrap': [False], 'n_estimators': [3, 10], 'max_features': [2, 3, 4]}]

forest_reg = RandomForestRegressor(random_state=42)
grid_search = GridSearchCV(forest_reg, param_grid, cv=5,
                           scoring='neg_mean_squared_error',
                           return_train_score=True)
grid_search.fit(myDataset_prepared, myDataset_labels)

print("Best Params: ", grid_search.best_params_)
print("Best Estimator: ", grid_search.best_estimator_)
print("\nResults (mean_test_score and params):")
cvres = grid_search.cv_results_
for mean_score, params in zip(cvres["mean_test_score"], cvres["params"]):
    print(np.sqrt(-mean_score), params)
```

```
Best Params: {'max_features': 6, 'n_estimators': 50}
Best Estimator: RandomForestRegressor(max_features=6, n_estimators=50,
random_state=42)

Results (mean_test_score and params):
141.45328245705397 {'max_features': 2, 'n_estimators': 30}
142.78195035217828 {'max_features': 2, 'n_estimators': 40}
142.37123716346338 {'max_features': 2, 'n_estimators': 50}
135.775474892488 {'max_features': 4, 'n_estimators': 30}
134.4190304245193 {'max_features': 4, 'n_estimators': 40}
135.55354286489987 {'max_features': 4, 'n_estimators': 50}
134.16110647203996 {'max_features': 6, 'n_estimators': 30}
134.1926275989663 {'max_features': 6, 'n_estimators': 40}
132.63913672411098 {'max_features': 6, 'n_estimators': 50}
136.15995027100854 {'max_features': 8, 'n_estimators': 30}
134.6124097344303 {'max_features': 8, 'n_estimators': 40}
```

```
133.2709123855618 {'max_features': 8, 'n_estimators': 50}
137.99385113493992 {'max_features': 10, 'n_estimators': 30}
137.62914189800856 {'max_features': 10, 'n_estimators': 40}
137.05278896563013 {'max_features': 10, 'n_estimators': 50}
163.96946102421438 {'bootstrap': False, 'max_features': 2, 'n_estimators': 3}
144.03648391473084 {'bootstrap': False, 'max_features': 2, 'n_estimators': 10}
149.43499632921868 {'bootstrap': False, 'max_features': 3, 'n_estimators': 3}
139.80926358472138 {'bootstrap': False, 'max_features': 3, 'n_estimators': 10}
143.39742710695754 {'bootstrap': False, 'max_features': 4, 'n_estimators': 3}
137.45096056556272 {'bootstrap': False, 'max_features': 4, 'n_estimators': 10}
```

GridSearchCV on LinearRegressor

```
param_grid = [
    {'fit_intercept': [True], 'n_jobs': [2, 4, 6, 8, 10]},
    {'normalize': [False], 'n_jobs': [3, 10]},
]

lin_reg = LinearRegression()
# train across 5 folds, that's a total of (12+6)*5=90 rounds of training
lin_grid_search = GridSearchCV(lin_reg, param_grid, cv=5,
                               scoring='neg_mean_squared_error',
                               return_train_score=True)
lin_grid_search.fit(myDataset_prepared, myDataset_labels)

print("Best Params: ", lin_grid_search.best_params_)
print("Best Estimator: ", lin_grid_search.best_estimator_)
print("\nResults (mean_test_score and params):")
cvres = lin_grid_search.cv_results_
for mean_score, params in zip(cvres["mean_test_score"], cvres["params"]):
    print(np.sqrt(-mean_score), params)
```

```
Best Params:  {'fit_intercept': True, 'n_jobs': 2}
Best Estimator: LinearRegression(n_jobs=2)

Results (mean_test_score and params):
214.25154731008828 {'fit_intercept': True, 'n_jobs': 2}
214.25154731008828 {'fit_intercept': True, 'n_jobs': 4}
214.25154731008828 {'fit_intercept': True, 'n_jobs': 6}
214.25154731008828 {'fit_intercept': True, 'n_jobs': 8}
214.25154731008828 {'fit_intercept': True, 'n_jobs': 10}
214.25154731008828 {'n_jobs': 3, 'normalize': False}
214.25154731008828 {'n_jobs': 10, 'normalize': False}
```

Randomized Search

```

        random_state=42)
rnd_search.fit(myDataset_prepared, myDataset_labels)

print("Best Params: ", rnd_search.best_params_)
print("Best Estimator: ", rnd_search.best_estimator_)
print("\nResults (mean_test_score and params):")
cvres = rnd_search.cv_results_
for mean_score, params in zip(cvres["mean_test_score"], cvres["params"]):
    print(np.sqrt(-mean_score), params)

```

```

Best Params: {'max_features': 7, 'n_estimators': 180}
Best Estimator: RandomForestRegressor(max_features=7, n_estimators=180,
random_state=42)

Results (mean_test_score and params):
133.8480285550475 {'max_features': 7, 'n_estimators': 180}
136.57169310969803 {'max_features': 5, 'n_estimators': 15}
139.51632044685252 {'max_features': 3, 'n_estimators': 72}
137.1771768583531 {'max_features': 5, 'n_estimators': 21}
134.3830342778431 {'max_features': 7, 'n_estimators': 122}
139.33370530867487 {'max_features': 3, 'n_estimators': 75}
138.9071602851763 {'max_features': 3, 'n_estimators': 88}
135.5081501949792 {'max_features': 5, 'n_estimators': 100}
138.3278914625561 {'max_features': 3, 'n_estimators': 150}
184.05585159780387 {'max_features': 5, 'n_estimators': 2}

```

Analyze best models

```

feature_importances = grid_search.best_estimator_.feature_importances_
print("feature_importances:\n", feature_importances)
extra_attribs = ["number_reviews_dot_reviews_per_month"]
cat_encoder = full_pipeline.named_transformers_["cat"]
cat_one_hot_attribs = list(cat_encoder.categories_[0])
attributes = num_attribs + extra_attribs + cat_one_hot_attribs
print("\nattributes:\n", attributes)
my_list = sorted(zip(feature_importances, attributes), reverse=True)
print("\nMost important features (think about removing features):")
print("\n".join('{}' for _ in range(len(my_list))).format(*my_list))

```

```

feature_importances:
[0.16069378 0.07746518 0.16413731 0.04363706 0.02335222 0.05394379
0.10185993 0.17949253 0.04558588 0.00626952 0.13746286 0.00526594
0.00083401]

attributes:
['host_id', 'latitude', 'longitude', 'minimum_nights', 'number_of_reviews',
'reviews_per_month', 'calculated_host_listings_count', 'availability_365',
'number_reviews_dot_reviews_per_month', 'Entire home/apt', 'Hotel room',
'Private room', 'Shared room']

Most important features (think about removing features):
(0.17949252977627858, 'availability_365')
(0.1641373056263463, 'longitude')
(0.16069378145971047, 'host_id')

```

```
(0.13746285698834157, 'Hotel room')
(0.10185992760999112, 'calculated_host_listings_count')
(0.07746518330470566, 'latitude')
(0.05394378731977807, 'reviews_per_month')
(0.04558587773390613, 'number_reviews_dot_reviews_per_month')
(0.043637064748633575, 'minimum_nights')
(0.023352218877592826, 'number_of_reviews')
(0.0062695183764503, 'Entire home/apt')
(0.0052659368947157795, 'Private room')
(0.0008340112835496431, 'Shared room')
```

Evaluate final model on test dataset

```
final_model = grid_search.best_estimator_
print("final_model:\n", final_model)

X_test = strat_test_set.drop("price", axis=1)
y_test = strat_test_set["price"].copy()

X_test_prepared = full_pipeline.transform(X_test)
final_predictions = final_model.predict(X_test_prepared)

final_mse = mean_squared_error(y_test, final_predictions)
final_rmse = np.sqrt(final_mse)

print ("final_predictions:\n", final_predictions )
print ("final_rmse:\n", final_rmse )

confidence = 0.95
squared_errors = (final_predictions - y_test) ** 2
mean = squared_errors.mean()
m = len(squared_errors)

# from scipy import stats
print("95% confidence interval: ",
      np.sqrt(stats.t.interval(confidence, m - 1,
                                loc=np.mean(squared_errors),
                                scale=stats.sem(squared_errors))))
)
```

```
final_model:
RandomForestRegressor(max_features=6, n_estimators=50, random_state=42)
final_predictions:
[132.24 78.2 96.3 ... 75.66 51.2 47.04]
final_rmse:
112.12588420446276
95% confidence interval: [ 53.7497262 149.18242105]
```

III. "Big Data" Machine Learning using the "Spark ML Library"

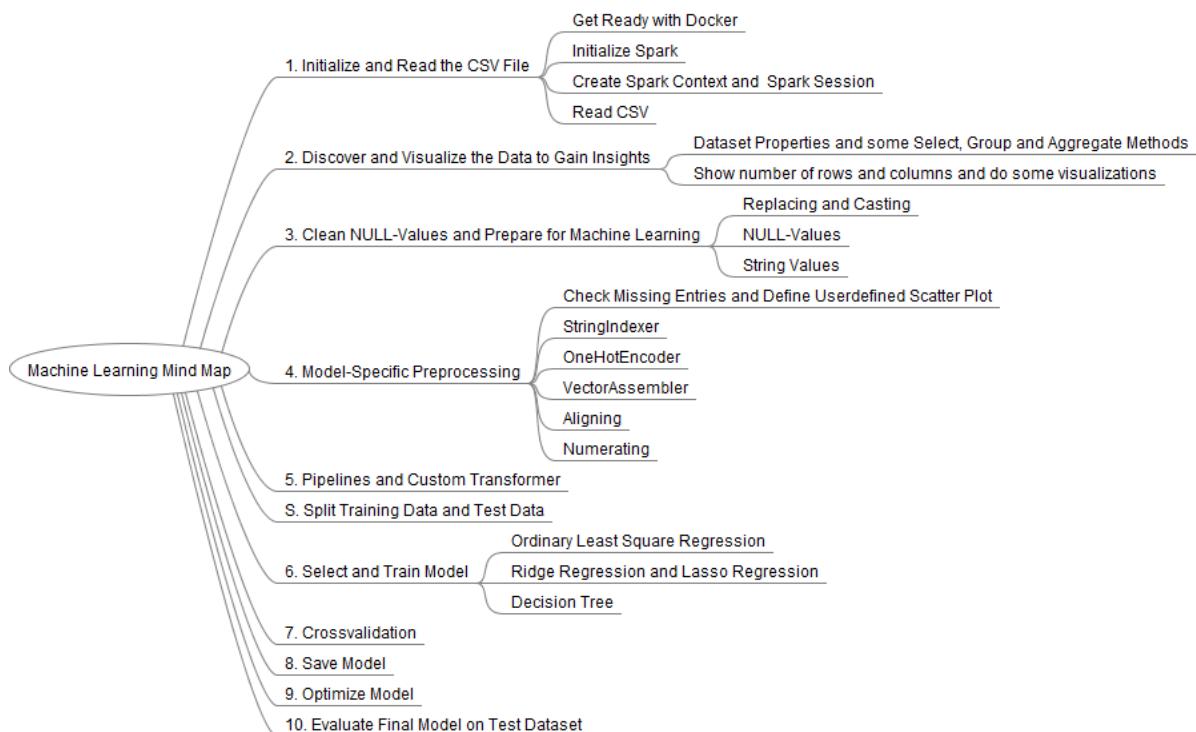
This will be an example for a "["Big-Data"](#)" environment and uses the "["Apache MLib"](#)" scalable machine learning library. Various tutorials, documentation, "code-fragments" and guidelines can be found in the internet **for free** (at least for your private use). The best is in my opinion the [official documentation](#). A few more helpful sources are the following GitHub repositories:

- [tirthajyoti/Spark-with-Python](#) (MIT license)
- [Apress/learn-pyspark](#) (Freeware License)
- [mahmoudparsian/pyspark-tutorial](#) (Apache License v2.0)

Concerning the topic "**Big Data**" I want to add the following: I passed a certification as "*Data Scientist Specialized in Big Data Analytics*". I must say: Understanding the concept of "Big-Data" and how to differentiate "standard" machine learning from a "scalable" environment is not easy. I recommend a separate training! Some steps are a bit similar to "scikit-learn" (e.g. data-cleaning, preprocessing), but the technical environment for running the code is different and also the code itself is different.

I added a "**Digression (Excuse)**" at the end of this document which covers the topics "*Big Data Visualization*", "*K-Means-Clustering in Spark*" and "*Map-Reduce*" (one of the [powerful programming models for Big Data](#)).

Let's start with the structure, which I put into a mind map (you can download it from this repository). I aligned the structure to the SkLearn mind map above in order to learn from each of this two approaches.



There are different ways to approach the Apache Spark and Hadoop environment: you can install it on your own computer (which I found rather difficult because of lack of user-friendly and easy understandable documentation). Or you can dive into a Cloud environment, like e.g. Microsoft Azure or Amazon EWS or Google Cloud and try to get a virtual machine up and running for your purposes. Have a look at my [documentation](#), where I shared my experiences, which I had with Microsoft Azure [here](#).

For the following explanation I decided to use [Docker](#). What is Docker? Docker is "*an open-source project that automates the deployment of software applications inside containers by providing an additional layer of abstraction and automation of OS-level virtualization on Linux.*" Learn from the [Docker-Curriculum](#) how it works. I found an container, which had Apache Spark Version 3.0.0 and Hadoop 3.2 installed and built my machine-learning code (using pyspark) on top of this container.

Initialize and Read the CSV File

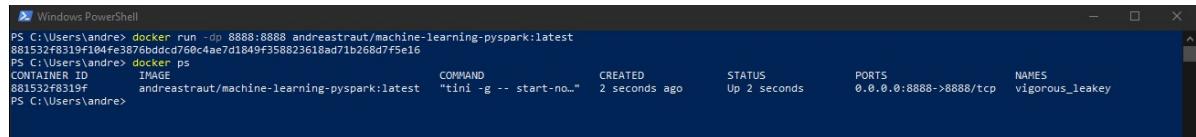
Get Ready with Docker

I shared my code and developments on Docker-Hub in the following repository [here](#). After having installed the Docker application you will need to pull my "machine-learning-pyspark" image to your computer:

```
docker pull andreastraut/machine-learning-pyspark
```

Then open Windows Powershell and type the following:

```
docker run -dp 8888:8888 andreastraut/machine-learning-pyspark:latest
```

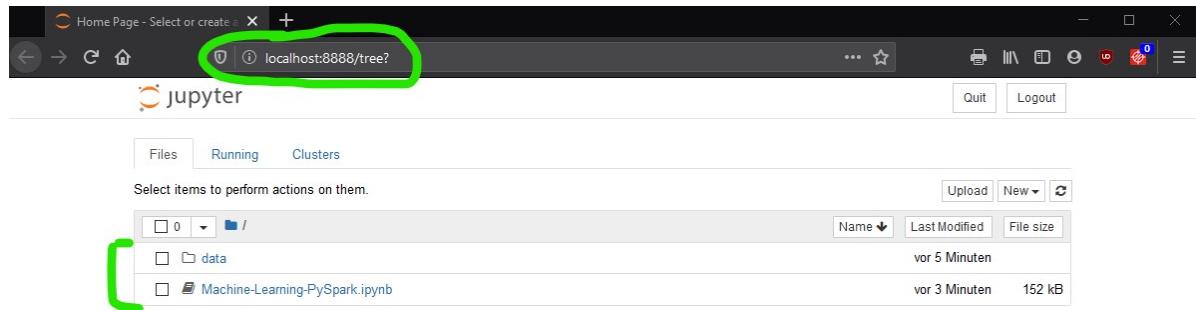


```
PS C:\Users\andre> docker run -dp 8888:8888 andreastraut/machine-learning-pyspark:latest
881532f8319f104fe3b76bddc7d80c4ae7d1849f358823618ad71b268d7f5e16
PS C:\Users\andre> docker ps
CONTAINER ID        IMAGE               COMMAND             CREATED            STATUS              PORTS               NAMES
881532f8319f        andreastraut/machine-learning-pyspark:latest   "tini -g -- start-no..."   2 seconds ago     Up 2 seconds          0.0.0.0:8888->8888/tcp   vigorous_leaky
PS C:\Users\andre>
```

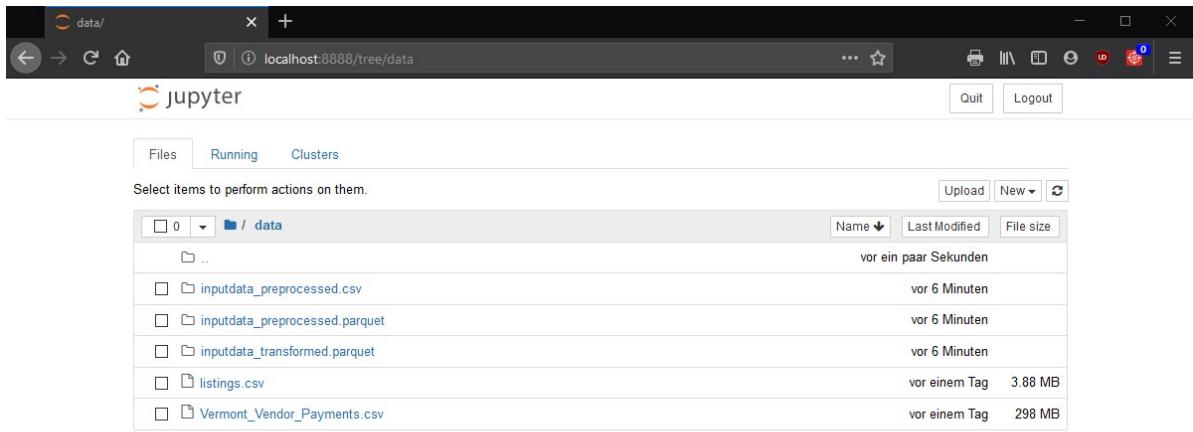
You will see in your Docker Dashboard that a container is running:



After having opened your browser (e.g. Firefox-Browser), navigate to "localhost:8888" (8888 is the port, which will be opened).



The folder "data" contains the datasets. If you would like to do further analysis or produce alternate visualizations of the Airbnb-data, you can download them from [here](#). It is available below under a [Creative Commons CC0 1.0 Universal \(CC0 1.0\) "Public Domain Dedication"](#) license. The data for the Vermont-Vendor-Payments can be downloaded from [here](#) and are available under the [Open Data Commons Open Database License](#).



When you open the Jupyter-Notebook, you will see, that Apache Spark Version 3.0.0 and Hadoop Version 3.2 is installed:

```

Author: Andreas Traut
Date: 23.07.2020

In [1]: import os
        print("APACHE_SPARK_VERSION: ", os.environ["APACHE_SPARK_VERSION"])
        print("HADOOP_VERSION: ", os.environ["HADOOP_VERSION"])
        print(os.environ)

APACHE_SPARK_VERSION: 3.0.0
HADOOP_VERSION: 3.2
{'SHELL': '/bin/bash', 'HOSTNAME': 'efcab749826e', 'LANGUAGE': 'en_US.UTF-8', 'SPARK_OPTS': '-driver-java-options=-Xms1024M --driver-java-options=-Xmx4096M --driver-java-options=-Dlog4j.logLevel=INFO', 'NB_UID': '1000', 'PWD': '/home/jovyan', 'MINICONDA_MD5': 'd63adf39f2c220950a063e0529d4ff74', 'HOME': '/home/jovyan', 'LANG': 'en_US.UTF-8', 'NB_GID': '100', 'XDG_CACHE_HOME': '/home/jovyan/.cache', 'APACHE_SPARK_VERSION': '3.0.0', 'PYTHONPATH': '/usr/local/spark/python:/usr/local/spark/python/lib/py4j-0.10.9-src.zip', 'HADOOP_VERSION': '3.2', 'SHLVL': '0', 'CONDA_DIR': '/opt/conda', 'MINICONDA_VERSION': '4.8.3', 'SPARK_HOME': '/usr/local/spark', 'CONDA_VERSION': '4.8.3', 'NB_USER': 'jovyan', 'LC_ALL': 'en_US.UTF-8', 'PATH': '/opt/conda/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/local/spark/bin', 'DEBIAN_FRONTEND': 'noninteractive', 'KERNEL_LAUNCH_TIMEOUT': '140', 'JPY_PARENT_PID': '6', 'TERM': 'xterm-color', 'CLICOLOR': '1', 'PAGE_R': 'cat', 'GIT_FINGERPRINT': 'cat', 'MPLBACKEND': 'module://ipykernel.pylab.backend_inline'}
```

```

In [2]: !conda list
# packages in environment at /opt/conda:
#
# Name          Version      Build  Channel
_libgcc_mutex   0.1          conda_forge    conda-forge
_openmp_mutex   4.5          0_gnu       conda-forge
abseil-cpp     20200225.2   he1b5a44_0  conda-forge
alembic         1.4.2        pyh9f0adid_0  conda-forge
arrow-cpp       0.17.1      py38h1234567_5_cpu  conda-forge
```

0. Initialize Spark

Initializing a Spark sessions works and reading a CSV file can be done with the following commands (see more documentation [here](#) and also have a look at a ["Get Started Guide"](#)):

```
In [3]: import pyspark
from pyspark.sql import SparkSession
from pyspark.sql import functions as F
```

0.1 Create Spark Context and Spark Session

```
In [4]: sc = pyspark.SparkContext(appName='Spark Modelling Context')
```

```
In [5]: spark = SparkSession.builder \
    .appName('Spark Modelling Session') \
    .config('spark.executor.memory','5g') \
    .config('spark.executor.cores','4') \
    .getOrCreate()
```

0.2 Read CSV

```
In [6]: import os
datapath = os.environ['PWD']
filename = datapath + "/data/listings.csv"
#read in data from csv
data = spark.read.csv(path=filename,
                      sep=',',
                      encoding='utf-8',
                      header=True,
                      inferSchema=True)
```

0.3 Dataset Properties and some Select, Group and Aggregate Methods

After then the data-cleaning and data preparation (eliminating of null values, visualization techniques) work pretty similar to the "Small data" (Sklearn) approach.

0.4 Write as Parquet or CSV

If you want to persist (=save) your intermediate you can do it as follows:

Persisting the preprocessed data

```
In [22]: data.select(*data.columns[:]).write.format("parquet") \
    .save("data/inputdata_preprocessed.parquet", mode='overwrite')

data.select(*data.columns[:]).write.csv('data/inputdata_preprocessed.csv', mode='overwrite', header=True)
<
```



```
In [23]: filename = "data/inputdata_preprocessed.parquet"
data = spark.read.parquet(filename)
data.show(5)

+-----+-----+-----+-----+-----+
| id | name | host_id | host_name | neighbourhood_group | neighbourhood | latitude | longitude |
+-----+-----+-----+-----+-----+
| 3176 | Fabulous Flat in ... | 3718 | Brittai | Pankow | Prenzlauer Berg S... | 52.5351 |
| 13.41758 | Entire home/apt | 90.0 | 62 | 145 | 2019-06-27 | 1.11 |
| 1.0 | 140 | 1 |
| 13309 | BerlinSpot Schöne... | 4108 | Jana | Tempelhof - Schön... | Schöneberg-Nord | 52.49865 |
| 13.34906 | Private room | 28.0 | 7 | 27 | 2019-05-31 | 0.34 |
| 1.0 | 320 | 1 |
| 16883 | Stylish East Side... | 16149 | Steffen | Friedrichshain-Kr... | Frankfurter Allee... | 52.51171 |
| 13.45477 | Entire home/apt | 125.0 | 3 | 133 | 2020-02-16 | 1.08 |
| 1.0 | 0 | 1 |
| 17071 | BrightRoom with s... | 17391 | BrightRoom | Pankow | Helmholtzplatz | 52.54316 |
| 13.41509 | Private room | 33.0 | 1 | 292 | 2020-03-06 | 2.27 |
| 2.0 | 45 | 1 |
| 19991 | Georgeous flat -... | 33852 | Philipp | Pankow | Prenzlauer Berg S... | 52.53303 |
| 13.41605 | Entire home/apt | 180.0 | 6 | 8 | 2020-01-04 | 0.14 |
| 1.0 | 8 | 1 |
+-----+-----+-----+-----+-----+
only showing top 5 rows
```

0.5 Read Parquet

If you want to persist (=save) your intermediate you can do it as follows:

```
data.select(*data.columns[:-1]).write.format("parquet").save("data/inputdata_preprocessed.parquet", mode='overwrite')
data.select(*data.columns[:-1]).write.csv('data/inputdata_preprocessed.csv', mode='overwrite', header=True)
```

Reading works as follows:

```
filename = "data/inputdata_preprocessed.parquet"
data = spark.read.parquet(filename)
```

1. Cleaning the data

1.2 Replacing and Casting

```
data = data.withColumn('price', F regexp_replace('price', '\$', ''))  
data = data.withColumn('price', F regexp_replace('price', ',', ''))  
data = data.withColumn('price', data['price'].cast('double'))
```

1.3 Null-Values

```
print("{} missing values for price"  
      .format(data  
              .filter(F.isnull(data['price']))  
              .count()))
```

```
104 missing values for price
```

```
data = data.fillna(0, subset='price')
```

```
print("{} missing values for  
price".format(data.filter(F.isnull(data['price'])).count()))
```

```
0 missing values for price
```

1.4 String Values

```
string_types = [x[0]  
                for x in data.dtypes  
                if x[1] == 'string']  
data.select(string_types).describe().show()
```

2. Model-specific preprocessing

2.0 Check missing entries and define userdefined scatter plot

2.1 StringIndexer

I included some examples of how features can be extracted, transformed and selected in the Jupyter-Notebook (see more documentation [here](#)). Just to mention a few here: the ["StringIndexer"](#), ["OneHotEncoder"](#) and ["VectorAssembler"](#) work as follows:

```
In [25]: data.select('neighbourhood_group').distinct().count()
Out[25]: 38

In [26]: from pyspark.ml.feature import StringIndexer
neighbourhood_indexer = StringIndexer(inputCol='neighbourhood_group', outputCol='neighbourhood_group')
neighbourhood_indexer_model = neighbourhood_indexer.fit(data)
data = neighbourhood_indexer_model.transform(data)

In [27]: data.groupby('neighbourhood_group').agg(F.collect_set('neighbourhood_group_index').alias('neighbourhood_group_index')).show()
```

neighbourhood_group	neighbourhood_group_index
Friedrichshain-Kreuzberg	[0.0]
Mitte	[1.0]
Pankow	[2.0]
Neukölln	[3.0]
Charlottenburg-Wilmersdorf	[4.0]
Tempelhof - Schöneberg	[5.0]
Lichtenberg	[6.0]
Treptow - Köpenick	[7.0]
Steglitz - Zehlendorf	[8.0]
Reinickendorf	[9.0]
Marzahn - Hellersdorf	[10.0]
Spandau	[11.0]
Downtown Apartments	[12.0]
Downtown Apartments	[13.0]
Neue Kantstraße	[14.0]
Alexanderplatz	[15.0]
Prenzlauer Berg Nord	[16.0]
Alt-Lichtenberg	[17.0]
Barstraße	[18.0]
Blankenfelde/Niederrad	[19.0]
Brunnenstr. Nord	[20.0]
Frankfurter Allee	[21.0]
Grunewald	[22.0]
Kurfürstendamm	[23.0]
Lisa	[24.0]
...	...

2.2 OneHotEncoder

```
In [28]: from pyspark.ml.feature import OneHotEncoder
one_hot_encoder = OneHotEncoder(
    inputCol = 'neighbourhood_group_index',
    outputCol = 'one_hot_neighbourhood_group',
    dropLast=False)
one_hot_encoder_model = one_hot_encoder.fit(data)
data = one_hot_encoder_model.transform(data)
```

2.3 VectorAssembler

```
In [29]: from pyspark.ml.feature import VectorAssembler
data = data.withColumn('number_of_reviews', data['number_of_reviews'].cast('double'))
data.select('number_of_reviews').show()
```

number_of_reviews
145.0
27.0
133.0
292.0
8.0
24.0
48.0
262.0
86.0
60.0
86.0
307.0
130.0
21.0
5.0
188.0
31.0
74.0
296.0
39.0

only showing top 20 rows

3. Aligning and numerating Features and Labels

3.1 Aligning

Aligning in this context means, that you use the variable `label` for the labeled column (here it is the "price" column) and the variable `feature_cols` for the features.

```
label = 'price'
feature_cols = ['num_features',
                 'vec_label']
cols = feature_cols + [label]
data_feat = data.withColumnRenamed(label, 'label').select(cols)
```

```
+-----+-----+
| num_features | vec_label | label |
+-----+-----+
| [145.0] | [90.0] | 90.0 |
| [27.0] | [28.0] | 28.0 |
| [133.0] | [125.0] | 125.0 |
| [292.0] | [33.0] | 33.0 |
| [8.0] | [180.0] | 180.0 |
+-----+-----+
only showing top 5 rows
```

3.2 Numerating

```
feat_len = len(numeric_attributes) +
data.select('room_type', 'room_index').distinct().count()
features = numeric_attributes + [r[0] for r in
data.select('room_type', 'room_index').distinct().collect()]
feature_dict = dict(zip(range(0, feat_len), features))
```

4. Pipelines

```
from pyspark.ml import Pipeline
pipeline = Pipeline(stages=[vec_num,
                           vec_label])
pipeline_model = pipeline.fit(data)
```

If you want to work with another pipeline, then all you need is to replace the second line. For example into this:

```
pipeline = Pipeline(stages=[room_indexer,
                           vec_num,
                           vec_label])
```

5. Training data and Testing data

```
data_train, data_test = data_feat.randomSplit([0.9, 0.1],
                                              seed=42)
```

This creates a random split into a training dataset `data_train` and the testing dataset `data_test`. Remember, that in the movies data example (see [[@sec:MoviesDatabaseExample](#)]) we used the stratified sample.

6. Apply models and evaluate

6.1 Ordinary Least Square Regression

After having extracted, transformed and selected features you will want to apply some models, which are documented [here](#), for example the "[OLS Regression](#)":

```
In [38]: from pyspark.ml.regression import LinearRegression  
lr = LinearRegression(featuresCol='num_features', labelCol='label', maxIter=1000, fitIntercept=True)  
  
In [39]: lr_model = lr.fit(data_train)  
lr_model.coefficients  
  
Out[39]: DenseVector([-0.4689])  
  
In [40]: pred = lr_model.transform(data_test)
```

6.2 Ridge Regression

The coding syntax for the "Ridge Regression" and the "Lasso Regression" are pretty similar (there are also the `fit` and the `transform` methods). Please read the official documentation in order to understanding the "Ridge Regression and the Lasso Regression" [^ridgelassoregression] in detail.

6.3 Lasso Regression

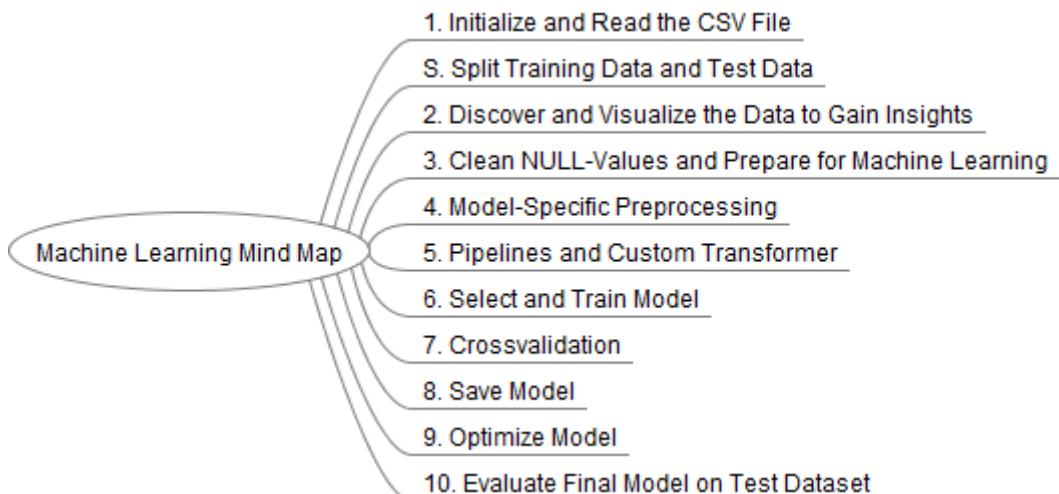
The coding syntax for the "Decision Tree" is pretty similar (there is also the `fit` and the `transform` method). Please read the official documentation in order to understanding the "Decision Tree" [^decisiontree] in detail.

7. Minhash und Local-Sensitive-Hashing (LSH)

see example: https://github.com/AndreasTraut/Deep_learning_explorations

IV. Summary Mind-Map

To summarize the whole coding structure have a look at the following and also the provided mind-maps. My mind map below may help you to structure your code:

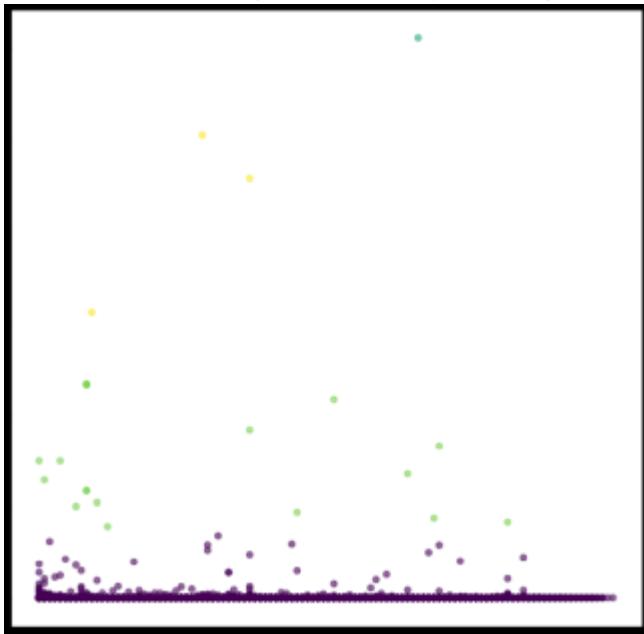


V. Digression (Excurs) to Big Data Visualization and K-Means Clustering Algorithm and Map-Reduce

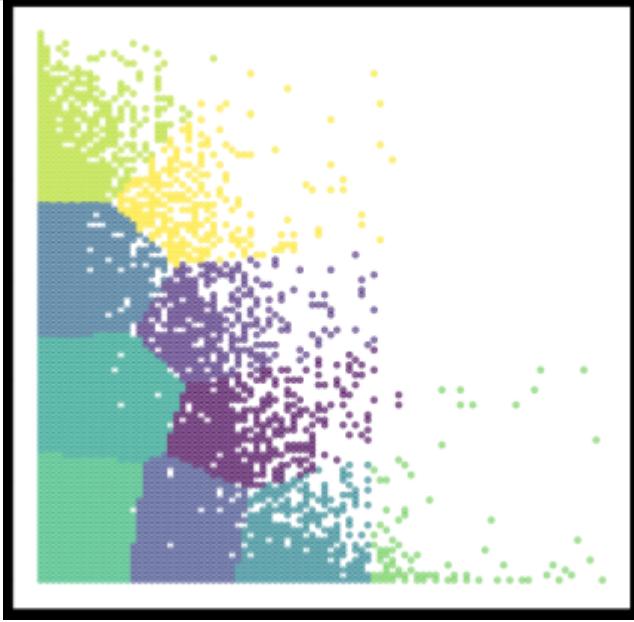
Digression (Excurs) to Big Data Visualization and K-Means Clustering Algorithm and Map-Reduce

(i) Big Data Visualization: You will see a Jupyter-Notebook (which contains the Machine-Learning Code) and a folder named "data" (which contains the raw-data and preprocessed data). As you can see: I also worked on a 298 MB big csv-file ([\("Vermont Vendor Payments.csv"\)](#)), which I couldn't open in Excel, because of the huge size. This file contains a list of all state of Vermont payments to vendors (Open Data Commons License) and has more than 1.6 million lines (exactly 1'648'466 lines). I already mentioned in my repository ["Visualization-of-Data-with-Python"](#), that the **visualization of big datasets** can be difficult when using "standard" office tools, like Excel. If you are not able to open such csv-files in Excel you have to find other solutions. One is to use PySpark which I will show you here. Another solution would have been to use the Excel built-in connection, [PowerQuery](#) or something similar, maybe Access or whatever, which is not the topic here, because we also want to be able to apply machine-learning algorithms from the [Spark Machine Learning Library](#). And there are more benefits of using PySpark instead of Excel: it can handle distributed processing, it's a lot faster, you can use pipelines, it can read many file systems (not only csv), it can process real-time data.

(ii) K-Means Clustering Algorithm: Additionally I worked on this dataset to show how the K-Means Clustering Algorithm can be applied by using the Spark Machine-Learning Library (see more documentation [here](#)). I will show how the "Vermont Vendor Payments" dataset can be clustered. In the images below every color represents a different cluster:



Digression (Excurs) to Big Data Visualization and K-Means Clustering Algorithm and Map-Reduce



(iii) Map-Reduce: This is a programming model for generating big data sets with parallel distributed algorithm on a cluster. Map-Reduce is very important for Big Data and therefore I added some Jupyter-Notebooks to better understand how it works. Learn the basis of the *Map-Reduce* programming model from [here](#) and then have a look into my jupyter notebook for details. I used the very popular "Word Count" example in order to explain Map-Reduce in detail.

In another application of Map-Reduce I found the very popular **term frequency-inverse document frequency** (short **TF-idf**) very interesting (see [Wikipedia](#)). This is a numerical statistic, which is often used in text-based recommender systems and for information retrieval. In my example I used the texts of "Moby Dick" and "Tom Sawyer". The result are two lists of most important words for each of these documents. This is what the TF-idf found:
Moby Dick: WHALE, AHAB, WHALES, SPERM, STUBB, QUEEQUEG, STRARBUCK, AYE
Tom Sawyer: HUCK, TOMS, BECKY, SID, INJUN, POLLY, POTTER, THATCHER
Applications for using TF-idf are in the [information retrieval](#) or to classify documents.

Have a look into my notebook [here](#) to learn more about Big Data Visualization, K-Means Clustering Algorithm, Map-Reduce and TF-idf.

Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License



This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.