

CS550 - Machine Learning and Business Intelligence

End-to-End Machine Learning Project

By : Prince Rajasekaran
Instructor: Dr. Chang, Henry

Table of content

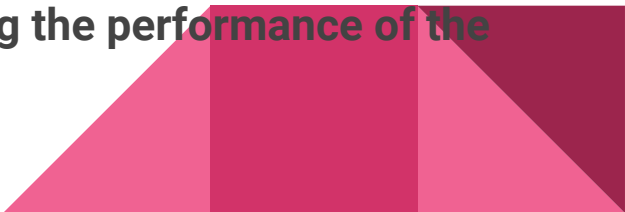
- **Introduction**
- **Theory**
- **Implementation**
- **Conclusion**
- **References**



Introduction

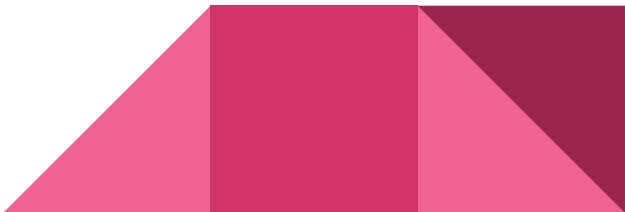
An end-to-end machine learning project entails creating a comprehensive pipeline that uses raw data as input, trains a machine learning model on it, and then applies the trained model to new, unforeseen data to make predictions. These projects typically consist of several steps, including model selection and training, feature engineering, data collection and preprocessing, and deployment of the trained model.

A project that uses end-to-end machine learning aims to produce a trustworthy and accurate model that can be applied to a particular issue, such as image classification, natural language processing, or predictive analytics. This entails not only choosing the right algorithms and methods for the job at hand but also comprehending the demands and constraints that are unique to the domain and maximizing the performance of the model.



Theory

7 steps involved in Machine learning :

- 1. Data collection**
 - 2. Data preparation**
 - 3. Choose a model**
 - 4. Train the model**
 - 5. Evaluate the model**
 - 6. Parameter tuning**
 - 7. Make predictions**
- 

1. Data collection

- The quantity & quality of your data dictate how accurate our model is
- The outcome of this step is generally a representation of data which we will use for training
- Using pre-collected data



2. Data preparation

- Wrangle data and prepare it for training
- Clean that which may require it Ex : remove duplicates, correct errors, deal with missing values, normalization, data type conversions, etc.
- Randomize data, which erases the effects of the particular order in which we collected and/or otherwise prepared our data
- Visualize data to help detect relevant relationships between variables or class imbalances (bias alert!), or perform other exploratory analysis
- Split into training and evaluation sets



3. Choose a model

Different algorithms are for different tasks; choose the right one

Example :

- 1. Linear Regression**
 - 2. Logistic Regression**
 - 3. Decision Tree**
 - 4. SVM (Support Vector Machine)**
 - 5. Naive Bayes**
 - 6. kNN (k- Nearest Neighbors)**
 - 7. K-Means**
- 

4. Train the model

- The goal of training is to answer a question or make a prediction correctly as often as possible
- Linear regression example: algorithm would need to learn values for m (or W) and b (x is input, y is output)
- Each iteration of process is a training step



5. Evaluate the model

- Uses some metric or combination of metrics to "measure" objective performance of model
- Test the model against previously unseen data
- This unseen data is meant to be somewhat representative of model performance in the real world, but still helps tune the model (as opposed to test data, which does not)
- Good train/eval split? 80/20, 70/30, or similar, depending on domain, data availability, dataset particulars, etc.



6. Parameter Tuning

- This step refers to *hyperparameter* tuning, which is an "artform" as opposed to a science

Tune model parameters for improved performance

Simple model hyperparameters may include: number of training steps, learning rate, initialization values and distribution, etc.



7. Make predictions

Using further (test set) data which have, until this point, been withheld from the model (and for which class labels are known), are used to test the model, a better approximation of how the model will perform in the real world



Implementation (using Colab)

Environment : Colab, Tensorflow 2

Programming Language: Python

Library :

```
import sys
```

```
import numpy as np
```

```
import pandas as pd
```

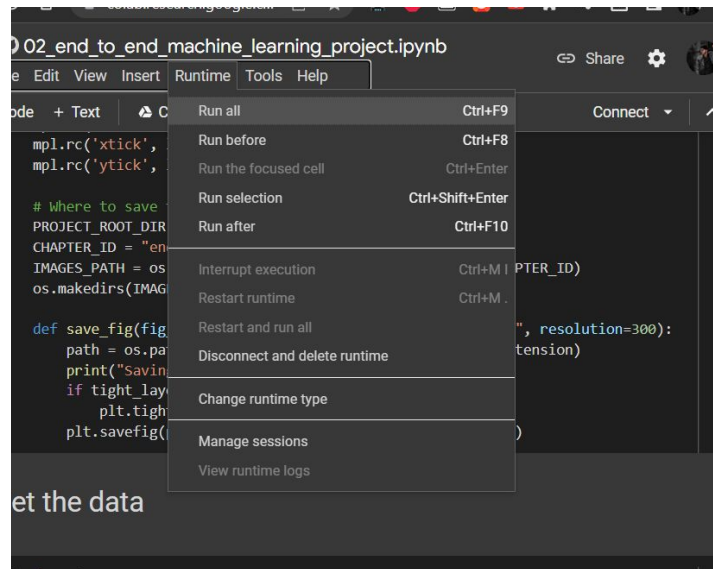
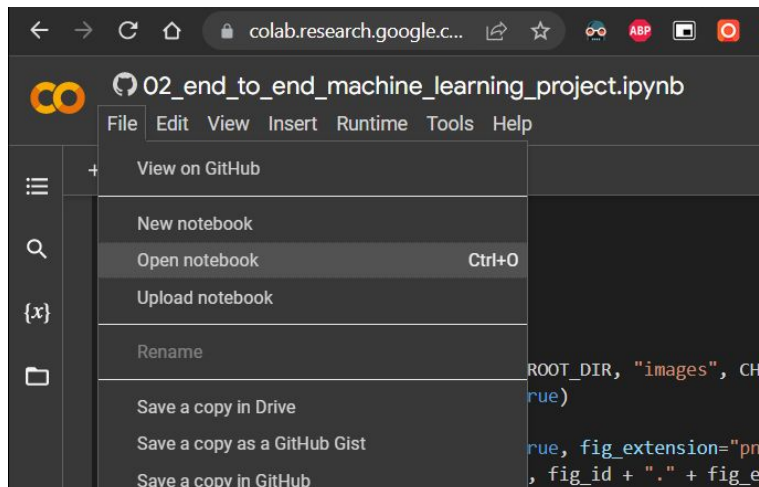
```
import os import matplotlib as mpl
```

```
import matplotlib.pyplot as plt
```



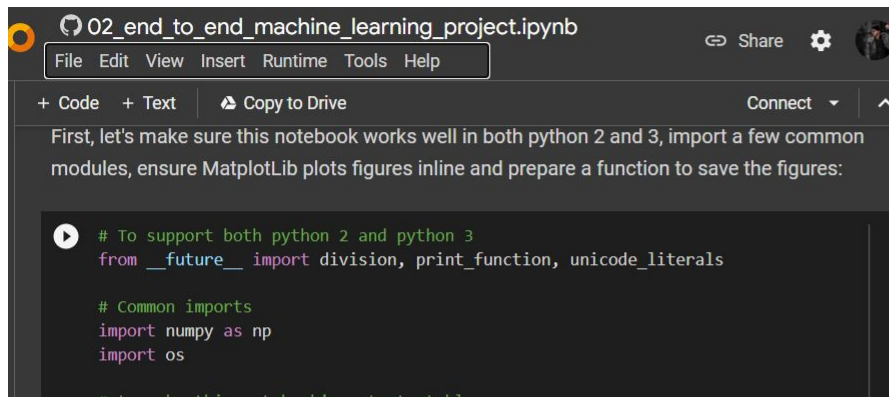
Implementation

Import and run the code



Implementation

Setting up the Colab notebook for the desired version (python 2 or 3) by importing some modules



The screenshot shows a Google Colab notebook titled "02_end_to_end_machine_learning_project.ipynb". The interface includes a menu bar with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". Below the menu bar, there are tabs for "+ Code" and "+ Text", and a "Copy to Drive" button. The notebook content starts with a text cell explaining the goal: "First, let's make sure this notebook works well in both python 2 and 3, import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures:". This is followed by a code cell containing the following Python code:

```
# To support both python 2 and python 3
from __future__ import division, print_function, unicode_literals

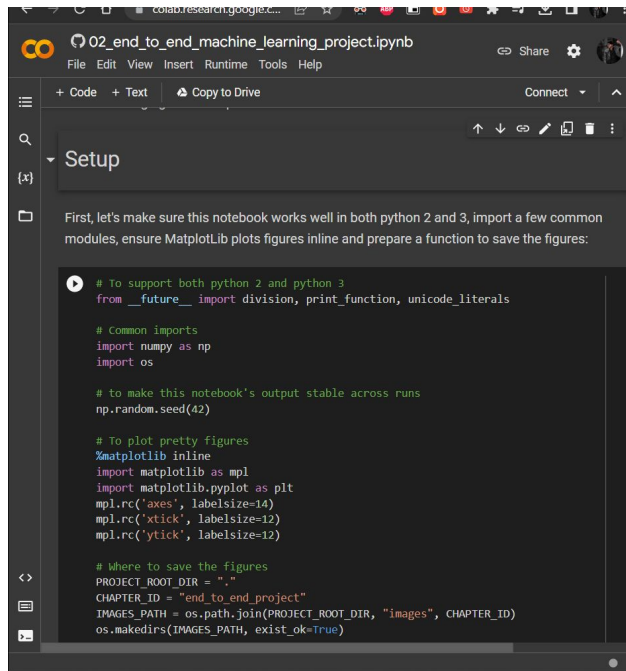
# Common imports
import numpy as np
import os
```

The code cell has a play button icon on the left, indicating it is ready to be executed. The notebook interface also shows a "Connect" button and a user profile icon in the top right corner.

Implementation

Getting the data We have housing data residing in California that contains properties of a house like

longitude,latitude,housing_median_age,total_rooms,total_bedrooms,population,households, median_income, median_house_value and ocean_proximity



```
02_end_to_end_machine_learning_project.ipynb
File Edit View Insert Runtime Tools Help
+ Code + Text Copy to Drive
Setup
First, let's make sure this notebook works well in both python 2 and 3, import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures:

# To support both python 2 and python 3
from __future__ import division, print_function, unicode_literals

# Common imports
import numpy as np
import os

# to make this notebook's output stable across runs
np.random.seed(42)

# To plot pretty figures
%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt
mpl.rc('axes', labelsz=14)
mpl.rc('xtick', labelsz=12)
mpl.rc('ytick', labelsz=12)

# Where to save the figures
PROJECT_ROOT_DIR = "."
CHAPTER_ID = "end_to_end_project"
IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
os.makedirs(IMAGES_PATH, exist_ok=True)
```

Implementation

After Fetching the data, we can have the loaded data on our Colab, as shown below.

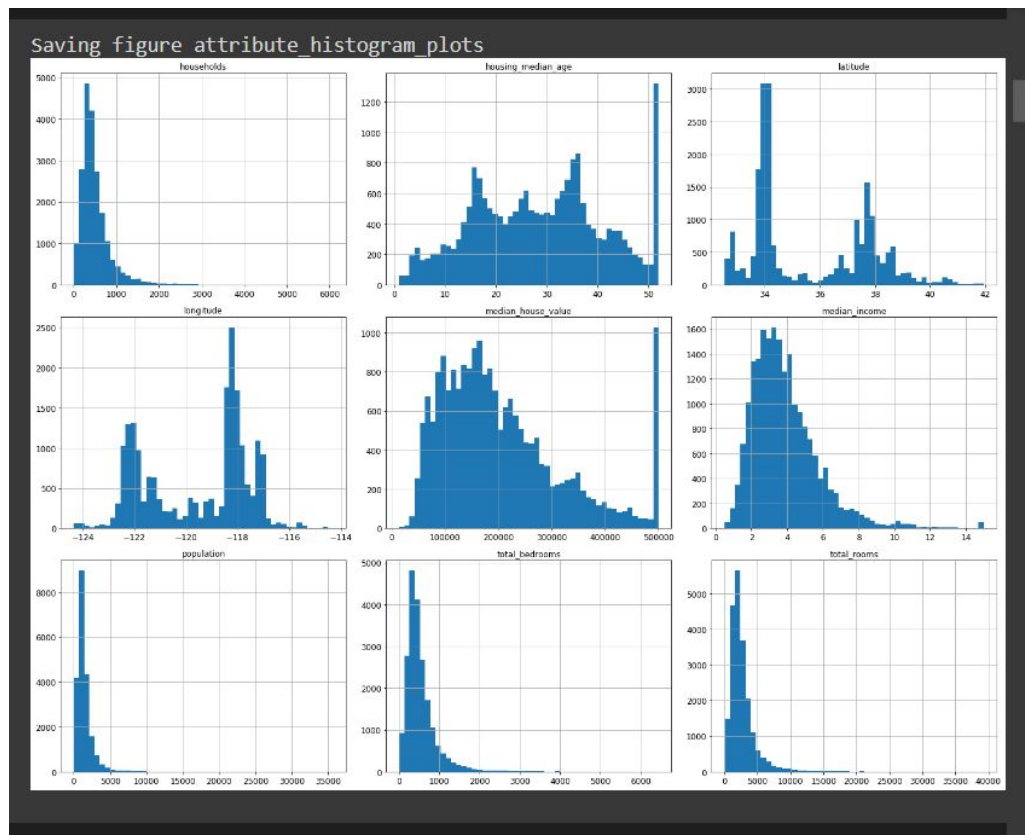
```
housing.head()
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252	452600.0	NEAR BAY
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014	358500.0	NEAR BAY
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574	352100.0	NEAR BAY
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431	341300.0	NEAR BAY
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462	342200.0	NEAR BAY

```
[1]: housing.info()
```


Implementation

We have the generated figure for the histogram information of the house.

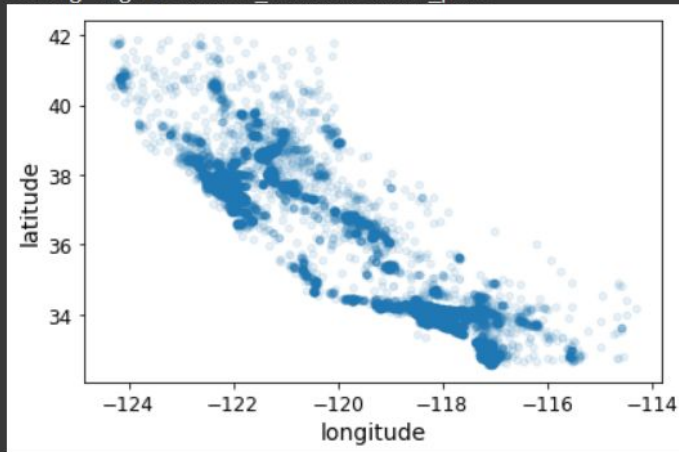


Implementation

Discover and visualize the data to gain insights

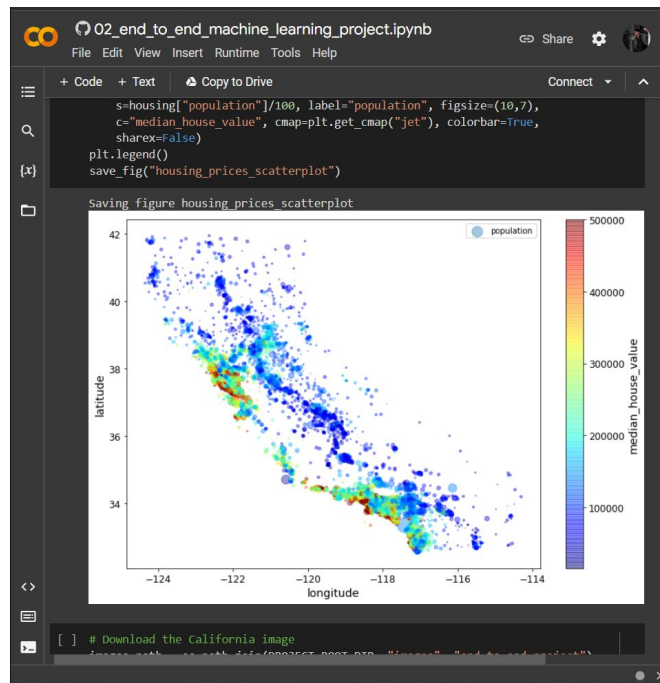
```
] housing.plot(kind="scatter", x="longitude", y="latitude", alpha=0.5)  
save_fig("better_visualization_plot")
```

Saving figure better_visualization_plot



The argument `shades=False` fixes a display bug (the x-axis values and legend

Implementation



Implementation

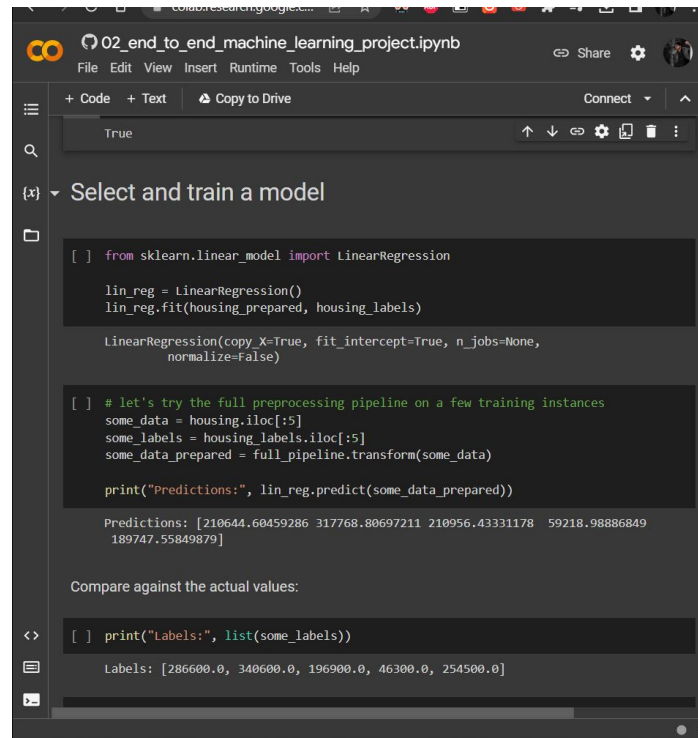
Prepare the data for Machine Learning algorithms

- While cleaning the data, we have the preprocessing phase, where we clean the data to create the machine learning model.
- The preprocessing steps include cleaning the data that is filling missing information or dropping unnecessary columns or rows that are not important later for the model to be created



Implementation

Select and train a model



```
02_end_to_end_machine_learning_project.ipynb
File Edit View Insert Runtime Tools Help
+ Code + Text Copy to Drive Connect
True
Select and train a model
[ ] from sklearn.linear_model import LinearRegression

lin_reg = LinearRegression()
lin_reg.fit(housing_prepared, housing_labels)

LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None,
normalize=False)

[ ] # let's try the full preprocessing pipeline on a few training instances
some_data = housing.iloc[:5]
some_labels = housing_labels.iloc[:5]
some_data_prepared = full_pipeline.transform(some_data)

print("Predictions:", lin_reg.predict(some_data_prepared))

Predictions: [210644.60459286 317768.80697211 210956.43331178 59218.98886849
189747.55849879]

Compare against the actual values:

[ ] print("Labels:", list(some_labels))

Labels: [286600.0, 340600.0, 196900.0, 46300.0, 254500.0]
```

Implementation

Fine tune your model

The screenshot shows a Jupyter Notebook interface with the following components:

- Top Bar:** Google Colab logo, file name "02_end_to_end_machine_learning_project.ipynb", and a "Cannot save changes" warning.
- Menu Bar:** File, Edit, View, Insert, Runtime, Tools, Help.
- Left Sidebar:** Navigation icons for code, search, and file explorer.
- Code Cell:** Contains Python code for fine-tuning a model.

```
[91] from sklearn.model_selection import cross_val_score

scores = cross_val_score(tree_reg, housing_prepared, housing_labels,
                          scoring="neg_mean_squared_error", cv=10)
tree_rmse_scores = np.sqrt(-scores)

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

display_scores(tree_rmse_scores)

Scores: [72831.45749112 69973.18438322 69528.56551415 72517.78229792
66846.14089488 72528.03725385 73997.08050233 68802.33629334
66443.28836884 70139.79923956]
Mean: 71629.89009727491
Standard deviation: 2914.035468468928

[93] lin_scores = cross_val_score(lin_reg, housing_prepared, housing_labels,
                                scoring="neg_mean_squared_error", cv=10)
lin_rmse_scores = np.sqrt(-lin_scores)
display_scores(lin_rmse_scores)

Scores: [71762.76364394 64114.99166359 67771.17124356 68635.19072082
66846.14089488 72528.03725385 73997.08050233 68802.33629334
66443.28836884 70139.79923956]
Mean: 69104.07998247063
```
- Output Cell:** Displays the output of the code execution, showing a table of scores and their statistics.

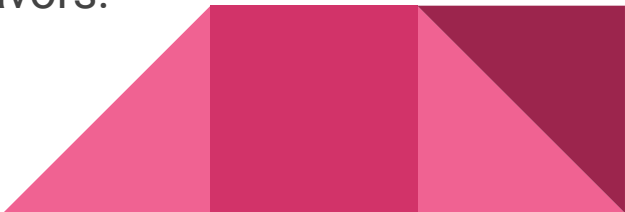
Score
72831.45749112
69973.18438322
69528.56551415
72517.78229792
66846.14089488
72528.03725385
73997.08050233
68802.33629334
66443.28836884
70139.79923956

Mean: 71629.89009727491
Standard deviation: 2914.035468468928

Conclusion

Finally, successful End-to-End Machine Learning projects necessitate a blend of technical proficiency, subject-matter knowledge, and project management abilities. Collaboration between data scientists, software engineers, subject matter experts, and stakeholders is required for them.

To make sure the project achieves its goals and adds value for the company or organization. The ability to harness the power of machine learning to gain insights, make predictions, and spur innovation is what makes End-to-End Machine Learning projects potentially difficult but rewarding endeavors.



References

Banoula, M. (2023, February 16).

Machine learning steps: A complete guide: Simplilearn. Simplilearn.com. Retrieved February 21, 2023, from <https://www.simplilearn.com/tutorials/machine-learning-tutorial/machine-learning-steps>

ML-ops.org. End-to-end Machine Learning Workflow. (2022, July 22).

Retrieved February 21, 2023, from <https://ml-ops.org/content/end-to-end-ml-workflow>

