## **Practical Machine Learning 101**

An end-to-end machine learning project

Justin Looye – IIJ Research Laboratory



#### Takeaway from today's talk - 今日の話から学べること

- Overview of machine learning 機械学習の概要
- Machine learning pipeline in practice with an end-to-end project
  エンドツーエンドプロジェクトによる実践的な機械学習パイプライン



zdataset.com: A Guide to Machine Learning Pipelines and Orchest

 Popular python tools – 人気のある Python ツール Pandas, Scikit-Learn, MLflow, FastAPI



 "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow 3<sup>rd</sup> edition" by Aurélien Géron (2022)



- "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow 3<sup>rd</sup> edition" by Aurélien Géron (2022)
- Well written, easy to follow with notebooks



- "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow 3<sup>rd</sup> edition" by Aurélien Géron (2022)
- Well written, easy to follow with notebooks
- Low priors



- "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow 3<sup>rd</sup> edition" by Aurélien Géron (2022)
- Well written, easy to follow with notebooks
- Low priors
- Permissive apache license



• Let's consider a spam filter:

#### Let's consider a spam filter:





#### Let's consider a spam filter:



Software carpentry https://wvuhpc.github.io/2019-Machine-Learning/01-introduction/index.html

#### Let's consider a spam filter:



Software carpentry https://wvuhpc.github.io/2019-Machine-Learning/01-introduction/index.html

- Advantages of ML:
  - Automated decision making (easy to update)
  - No human priors (unicode abuse)

#### Machine learning problems in the real world







#### Machine learning problems in the real world





2024-07-09 5 /

#### Machine learning problems in the real world





 Insufficient quantity of training data

- Insufficient quantity of training data
- Nonrepresentative training data



Wikipedia: Martin Grandjean (vector), McGeddon (picture), Cameron Moll (concept)



- Insufficient quantity of training data
- Nonrepresentative training data
- Poor-Quality data



Wikipedia: Martin Grandjean (vector), McGeddon (picture), Cameron Moll (concept)

- Insufficient quantity of training data
- Nonrepresentative training data
- Poor-Quality data
- Irrelevant features



Wikipedia: Martin Grandjean (vector), McGeddon (picture), Cameron Moll (concept)

- Insufficient quantity of training data
- Nonrepresentative training data
- Poor-Quality data
- Irrelevant features
- Overfitting the training data
- Underfitting the training data



Jeeva Saravanan @ Medium: How to Evaluate your Machine Learning model

## **End-to-end ML project**

- We are a recently hired data scientist in a real estate company
- The company want to be able to predict district median house values





## **End-to-end ML project**

- We are a recently hired data scientist in a real estate company
- The company want to be able to predict district median house values



Company gives us a dataset of district median house values

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



22

1)Look at the big picture – frame the problem2)Get the data

- 3)Discover and visualize data to gain insights
- 4) Prepare data for ML algorithms
- 5)Select a model and train it
- 6)Launch, monitor and maintain



# 1)Supervised, unsupervised, reinforcement learning?



1)Supervised, unsupervised, reinforcement learning?

2)Task? Regression, classification, clustering,...?



- 1)Supervised, unsupervised, reinforcement learning?
- 2)Task? Regression, classification, clustering,...?3)Performance criterion?

- 1)Supervised, unsupervised, reinforcement learning?
- 2)Task? Regression, classification, clustering,...?3)Performance criterion?
- 4) Baseline? Human-level performance? Experts?

 Open-source Python library providing highperformance, easy-to-use data structures and analysis tools







- Open-source Python library providing highperformance, easy-to-use data structures and analysis tools
- Key feature: DataFrame object for tabular data manipulation with integrated indexing pandas



- Open-source Python library providing highperformance, easy-to-use data structures and analysis tools
- Key feature: DataFrame object for tabular data manipulation with integrated indexing
- Core features:
  - Data manipulation
  - Data cleaning
  - File handling
  - Time series

pandas



#### **Tool #1: Pandas for handling data**

- Open-source Python library providing high-performance, easy-touse data structures and analysis tools
- Key feature: DataFrame object for tabular data manipulation with integrated indexing
- Core features:
  - Data manipulation
  - Data cleaning
  - File handling
  - Time series
- Use case in ML:
  - Data preparation
  - Exploratory data analysis and visualization





2024-07-09 10

## **Tool #2 : scikit-learn for Machine learning**

- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on Numpy, SciPy and matplotlib
- Open source, commercially usable – BSD license

#### Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition. Algorithms: Gradient boosting, nearest neighbors, random forest, logistic regression, and more...





#### **Dimensionality reduction**

Reducing the number of random variables to consider.

Applications: Visualization, increased efficiency. Algorithms: <u>PCA</u>, feature selection, <u>non-negative</u> matrix factorization, and <u>more...</u>



#### Regression

Model selection

models

tunina.

and more ...

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, stock prices. Algorithms: Gradient boosting, nearest neighbors, random forest, ridge, and more...



Comparing, validating and choosing parameters and

Applications: Improved accuracy via parameter

Algorithms: Grid search, cross validation, metrics,

Examples

#### **Clustering**

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, grouping experiment outcomes. Algorithms: <u>k-Means</u>, <u>HDBSCAN</u>, <u>hierarchical</u>

clustering, and more...





Examples

#### Preprocessing

Feature extraction and normalization.

Applications: Transforming input data such as text for use with machine learning algorithms. Algorithms: <u>Preprocessing</u>, feature extraction, and more...



Examples



1)Look at the big picture – frame the problem2)Get the data

- 3)Discover and visualize data to gain insights
- 4) Prepare data for ML algorithms
- 5)Select a model and train it
- 6)Launch, monitor and maintain



#### Get the data





#### Get the data

Parse the data

ho	housing.head()										
	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income			
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252			
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014			
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574			
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431			
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462			



housing - nd read csv(csv nath)
- Parse the data
- Take a quick look
  - Туре
  - Missing Values

#### [9]: housing.info()

<class 'pandas.core.frame.dataframe'=""> RangeIndex: 20640 entries, 0 to 20639 Data columns (total 10 columns):</class>									
#	Column Non-Null Count Dtype								
Θ	longitude	20640 non-null	float64						
1	latitude	20640 non-null	float64						
2	housing_median_age	20640 non-null	float64						
3	total rooms	20640 non-null	float64						
4	total_bedrooms	20433 non-null	float64 <						
5	population	20640 non-null	float64						
6	households	20640 non-null	float64						
7	median_income	20640 non-null	float64						
8	median house value	20640 non-null	float64	4					
9	ocean_proximity	20640 non-null	object 🖪						
dtyp	es: float64(9), obje	ct(1)							
memo	memory usage: 1.6+ MB								



- Parse the data
- Take a quick look
  - Туре
  - Missing Values
  - Basic stats

[10]: housing["ocean\_proximity"].value\_counts()
[10]: ocean\_proximity
 <1H OCEAN 9136
 INLAND 6551</pre>

NEAR OCEAN 2658 NEAR BAY 2290 ISLAND 5 Name: count, dtype: int64



- Parse the data
- Take a quick look
  - Туре
  - Missing Values
  - Basic stats

nousi	ing.describe()	)							
	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value
count	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.000000	20640.000000
mean	-119.569704	35.631861	28.639486	2635.763081	537.870553	1425.476744	499.539680	3.870671	206855.816909
std	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.899822	115395.615874
min	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.499900	14999.000000
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.563400	119600.000000
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.534800	179700.000000
75%	-118.010000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.743250	264725.000000
max	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.000100	500001.000000

- Parse the data
- Take a quick look
  - Туре
  - Missing Values
  - Basic stats



- Parse the data
- Take a quick look
  - Туре
  - Missing Values
  - Basic stats
- Train-test split

from sklearn.model\_selection import train\_test\_split

train\_set, test\_set = train\_test\_split(housing, test\_size=0.2, random\_state=42)





1)Look at the big picture – frame the problem2)Get the data

- 3)Discover and visualize data to gain insights
- 4) Prepare data for ML algorithms
- 5)Select a model and train it
- 6)Launch, monitor and maintain





 Visualize geographical data



Saving figure housing\_prices\_scatterplot



- Visualize geographical data
- Look for correlation

- [42]: corr\_matrix = housing.corr(numeric\_only=True)
- [43]: corr\_matrix["median\_house\_value"].sort\_values(ascending=False)

43]:	median_house_value	1.000000	
	median income	0.687151 <	
	total_rooms	0.135140	
	housing median age	0.114146	
	households	0.064590	
	total bedrooms	0.047781	
	population	-0.026882	
	longitude	-0.047466	
	latitude	-0.142673	
	Name: median house	value, dtype:	float64

- Visualize geographical data
- Look for correlation

Saving figure income\_vs\_house\_value\_scatterplot



- Visualize geographical data
- Look for correlation
- Experiment with attribute combinations

housing["rooms\_per\_household"] = housing["total\_rooms"]/housing["households"]
housing["bedrooms\_per\_room"] = housing["total\_bedrooms"]/housing["total\_rooms"]
housing["population per household"]=housing["population"]/housing["households"]

corr\_matrix = housing.corr(numeric\_only=True)
corr\_matrix["median\_house\_value"].sort\_values(ascending=False)

median_house_value	1.000000
median_income	0.687151
rooms_per_household	0.146255
total_rooms	0.135140
housing median age	0.114146
households	0.064590
total_bedrooms	0.047781
population_per_household	-0.021991
population	-0.026882
longitude	-0.047466
latitude	-0.142673
bedrooms_per_room	-0.259952
Name: median house value,	dtype: float64



1)Look at the big picture – frame the problem2)Get the data

- 3)Discover and visualize data to gain insights
- 4) Prepare data for ML algorithms
- 5)Select a model and train it
- 6)Launch, monitor and maintain

## Most ML algorithms work best with scaled numerical features





# Most ML algorithms work best with scaled numerical features

 Handling missing data

housing.dropna(subset=["total\_bedrooms"]) # option 1
housing.drop("total\_bedrooms", axis=1) # option 2
median = housing["total\_bedrooms"].median() # option 3
housing["total\_bedrooms"].fillna(median, inplace=True)





# Most ML algorithms work best with scaled numerical features

- Handling missing data
- Handling categorical features

<pre>housing_cat = housing[["ocean_proximity"] housing_cat.head(10)</pre>									
	ocean_proximity								
12655	INLAND								
15502	NEAR OCEAN								
2908	INLAND								
14053	NEAR OCEAN								
20496	<1H OCEAN								
1481	NEAR BAY								
18125	<1H OCEAN								
5830	<1H OCEAN								
17989	<1H OCEAN								
4861	<1H OCEAN								



## Most ML algorithms work best with scaled numerical features

#### Ordinal encoding

from sklearn.preprocessing import OrdinalEncoder

1.0 4.0

1.0 4.0

 Handling missing data

Handling categorical features

housing housing	cat = housing cat.head(10)	index= column		
	ocean_proximity		oce	an_proximity
12655	INLAND		12655	1.0
15502	NEAR OCEAN		15502	4.0
2908	INLAND		2908	1.0
14053	NEAR OCEAN		14053	4.0
20496	<1H OCEAN		20496	0.0
1481	NEAR BAY		1481	3.0
18125	<1H OCEAN		18125	0.0
5830	<1H OCEAN		5830	0.0
17989	<1H OCEAN		17989	0.0
4861	<1H OCEAN		4861	0.0

ordinal encoder = OrdinalEncoder() housing cat encoded = ordinal encoder.fit transform(housing cat) pd.DataFrame(housing cat encoded[:10, 0], dex=housing cat.head(10).index, lumns=["ocean proximity"])



## Most ML algorithms work best with scaled numerical features

One-hot encoding

## Handling missing data

## Handling categorical features

housing housing	_cat = housin _cat.head(10)
	ocean_proximity
12655	INLAND
15502	NEAR OCEAN
2908	INLAND
14053	NEAR OCEAN
20496	<1H OCEAN
1481	NEAR BAY
18125	<1H OCEAN
5830	<1H OCEAN
17989	<1H OCEAN
4861	<1H OCEAN

from sklearn.preprocessing import OneHotEncoder

# One-hot encode the categorical data onehot encoder = OneHotEncoder() housing cat lhot = onehot encoder.fit transform(housing cat)

# Pretty formatting for the presentation housing cat 1hot df = pd.DataFrame(housing cat 1hot.toarrav(). columns=onehot encoder.categories [0], index=housing cat.index) housing cat 1hot df.columns = [f"is {col.replace(' ', ' ').lower()}" for col in housing cat 1hot df.columns] housing cat 1hot df.head(10)

	is_<1h_ocean	is_inland	is_island	is_near_bay	is_near_ocean
12655	0.0	1.0	0.0	0.0	0.0
15502	0.0	0.0	0.0	0.0	1.0
2908	0.0	1.0	0.0	0.0	0.0
14053	0.0	0.0	0.0	0.0	1.0
20496	1.0	0.0	0.0	0.0	0.0
1481	0.0	0.0	0.0	1.0	0.0
18125	1.0	0.0	0.0	0.0	0.0
5830	1.0	0.0	0.0	0.0	0.0
17989	1.0	0.0	0.0	0.0	0.0
4861	1.0	0.0	0.0	0.0	0.0

# Most ML algorithms work best with scaled numerical features

- Handling missing data
- Handling categorical features
- Feature scaling
- Pipeline

```
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
num_pipeline = Pipeline([
               ('imputer', SimpleImputer(strategy="median")),
               ('attribs_adder', CombinedAttributesAdder()),
               ('std_scaler', StandardScaler()),
        ])
housing num tr = num pipeline.fit transform(housing num)
```



# Most ML algorithms work best with scaled numerical features

- Handling missing data
- Handling categorical features
- Feature scaling
- Pipeline

```
from sklearn.compose import ColumnTransformer
num_attribs = list(housing_num)
cat_attribs = ["ocean_proximity"]
full_pipeline = ColumnTransformer([
        ("num", num_pipeline, num_attribs),
        ("cat", OneHotEncoder(), cat_attribs),
])
```

housing\_prepared = full\_pipeline.fit\_transform(housing)

1)Look at the big picture – frame the problem2)Get the data

- 3)Discover and visualize data to gain insights
- 4) Prepare data for ML algorithms
- 5)Select a model and train it
- 6)Launch, monitor and maintain

#### Select a model and train it





 Train a model Predict

```
from sklearn.linear_model import LinearRegression
# Choose the model
lin_reg = LinearRegression()
# Fit to training data
lin_reg.fit(housing_prepared, housing_labels)
# 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: [ 85657.90192014 305492.60737488 152056.46122456 186095.70946094
```

244550.679660891



 Train a model Predict

#### Evaluate performance

from sklearn.metrics import mean\_squared\_error

housing\_predictions = lin\_reg.predict(housing\_prepared)
lin\_mse = mean\_squared\_error(housing\_labels, housing\_predictions)
lin\_rmse = np.sqrt(lin\_mse)
lin\_rmse

68627.87390018745



- Train a model Predict
- Evaluate performance

from sklearn.tree import DecisionTreeRegressor

tree\_reg = DecisionTreeRegressor(random\_state=42)
tree\_reg.fit(housing\_prepared, housing\_labels)

DecisionTreeRegressor

DecisionTreeRegressor(random\_state=42)

housing\_predictions = tree\_reg.predict(housing\_prepared)
tree\_mse = mean\_squared\_error(housing\_labels, housing\_predictions)
tree\_rmse = np.sqrt(tree\_mse)
tree\_rmse

0.0





- Train a model Predict
- Evaluate performance
- Avoid overfit with cross-validation



- Train a model Predict
- Evaluate performance
- Avoid overfit with cross-validation



### Select a model and train it

 Train a model Predict

 Evaluate performance

- Avoid overfit with cross-validation
- Try out different models/parameters

```
# Define model and parameters
models = {
    'DecisionTree': (DecisionTreeRegressor(), {'max_depth': [None, 10, 20], 'min_samples_leaf': [1, 2, 4]}),
    'RandomForest': (RandomForestRegressor(), {'n_estimators': [100, 200], 'max_features': ['log2', 'sqrt']}),
    'GradientBoosting': (GradientBoostingRegressor(), {'n_estimators': [100, 200], 'learning_rate': [0.1, 0.05, 0.01]})
```

```
# Initialize MLflow
mlflow.set_experiment('Regression_Models_Exploration')
```

```
for model_name, (model, params) in models.items():
    for param_set in ParameterGrid(params):
        with mlflow.start_run(run_name=f"{model_name}"):
            # Set the parameters and train the model
            model.set_params(**param_set)
```

```
# Log parameters, metrics, and model
mlflow.log_params(param_set)
mlflow.log_metric("mean_rmse", np.mean(rmse_scores))
mlflow.log_metric("std_rmse", np.std(rmse_scores))
```

# Log the model
mlflow.sklearn.log\_model(model, f"model\_{model\_name}")

print(f"Logged {model\_name} with params {param\_set} with mean RMSE: {np.mean(rmse\_scores)}")



 Train a model Predict

#### Evaluate performance

- Avoid overfit with cross-validation
- Try out different models/parameters

Qm	netrics.rmse < 1 and params.mo	del = "tree" (	i) Time	created 🗸	State: Active 🗸	Datasets 🗸	≞↑ Sort: me	ean_rmse 👻	Columns 🗸	🗄 Group by 🖌	:	C + New run
Table	Chart Evaluation Experim	nental Traces Experime	ntal									
						Metrics		Parameters				
	Run Name	Created	Duration	Source	Models	mean_rmse <u>=</u> ↑	std_rmse	max_depth	max_features	min_samples_leaf	n_estimators	
	RandomForest	2 minutes ago	1.6min	ipykernel	🐝 sklearn	49024.8782	1963.92759	-	log2	-	200	
	RandomForest	2 minutes ago	48.8s	ipykernel	🕵 sklearn	49487.4028	2059.16654	-	log2	-	100	
	DecisionTree	<ul> <li>4 minutes ago</li> </ul>	2.8s	ipykernel	🕵 sklearn	61076.8455	2099.94073	10	÷	4	÷	
	DecisionTree	3 minutes ago	3.0s	ipykernel	🕵 sklearn	61187.6778	2076.22989	10	-	4	-	
	DecisionTree	3 minutes ago	3.1s	ipykernel	🐝 sklearn	61866.0497	2297.44510	10		2	-	
	DecisionTree	⊘ 4 minutes ago	3.1s	ipykernel	🐝 sklearn	62191.0119	1829.66052	10		2	-	
	DecisionTree	3 minutes ago	2.9s	ipykernel	🐝 sklearn	62752.0142	1718.43596	10		1	-	
	DecisionTree		2.9s	ipykernel	🐝 sklearn	63078.1181	1961.64812	10	-	1	-	



## Select a model and train it

- Train a model Predict
- Evaluate performance
- Avoid overfit with cross-validation
- Try out different models/parameters
- Save with joblib

```
# Define model and parameters
models = {
    'DecisionTree': (DecisionTreeRegressor(), {'max depth': [None, 10, 20], 'min samples leaf': [1, 2, 4]}),
    'RandomForest': (RandomForestRegressor(), {'n estimators': [100, 200], 'max features': ['log2', 'sgrt']}),
    'GradientBoosting': (GradientBoostingRegressor(), { 'n estimators': [100, 200], 'learning rate': [0.1, 0.05, 0.01]})
best rmse = float('inf')
best model = None
# Initialize MLflow
mlflow.set experiment('Regression Models Exploration')
for model name. (model, params) in models.items():
    for param set in ParameterGrid(params):
        with mlflow.start run(run name=f"{model name}"):
            # Set the parameters and train the model
            model.set params(**param set)
            # Perform cross-validation
            scores = cross val score(model, housing prepared, housing labels,
                                     scoring="neg mean squared error", cv=10)
            rmse scores = np.sqrt(-scores)
            mean rmse = np.mean(rmse scores)
            # Log parameters, metrics, and model
            mlflow.log params(param set)
            mlflow.log metric("mean rmse", mean rmse)
            mlflow.log metric("std rmse", np.std(rmse scores))
            mlflow.sklearn.log model(model, f"model {model name}")
            # Update the best model if the current model is better
            if mean rmse < best rmse:</pre>
                best rmse = mean rmse
                best model = model
            print(f"Logged {model name} with params {param set} with mean RMSE: {mean rmse}")
```

2024-07-09

# Save the best model
joblib.dump(best\_model, 'best\_model.joblib')

Tech Trend Talk: Practical ML 101

1)Look at the big picture – frame the problem2)Get the data

- 3)Discover and visualize data to gain insights
- 4) Prepare data for ML algorithms
- 5)Select a model and train it
- 6)Launch, monitor and maintain

Let's deploy with FastAPI



#### Let's deploy with FastAPI

from fastapi import FastAPI
import joblib
import pandas as pd
from pydantic import BaseModel

# Load the data processing pipeline
full\_pipeline = joblib.load('full\_pipeline.joblib')

# Load the ML model
model = joblib.load('best\_model.joblib')

# Launch API
app = FastAPI()

#### class HousingData(BaseModel):

longitude: float = 119.0
latitude: float = 35.0
housing\_median\_age: float = 30.0
total\_rooms: float = 2635.0
total\_bedrooms: float = 500.0
population: float = 1425.0
households: float = 500.0
median\_income: float = 3.87
ocean proximity: str = "INLAND"

@app.post("/predict/")

def predict(housing\_data: HousingData):
 # Convert the input data to a DataFrame
 data\_dict = housing\_data.dict()
 df = pd.DataFrame([data\_dict])

# Use the loaded pipeline to transform the data
prepared\_data = full\_pipeline.transform(df)

# Predict using the loaded model
prediction = model.predict(prepared\_data)
return {"prediction": prediction.tolist()}



#### Let's deploy with FastAPI

from fastapi import FastAPI
import joblib
import pandas as pd
from pydantic import BaseModel

# Load the data processing pipeline
full\_pipeline = joblib.load('full\_pipeline.joblib')

# Load the ML model
model = joblib.load('best\_model.joblib')

# Launch API
app = FastAPI()

#### class HousingData(BaseModel):

longitude: float = 119.0
latitude: float = 35.0
housing\_median\_age: float = 30.0
total\_rooms: float = 2635.0
total\_bedrooms: float = 500.0
population: float = 1425.0
households: float = 500.0
median\_income: float = 3.87
ocean proximity: str = "INLAND"

#### @app.post("/predict/")

def predict(housing\_data: HousingData):
 # Convert the input data to a DataFrame
 data\_dict = housing\_data.dict()
 df = pd.DataFrame([data\_dict])

# Use the loaded pipeline to transform the data
prepared\_data = full\_pipeline.transform(df)

# Predict using the loaded model
prediction = model.predict(prepared\_data)
return {"prediction": prediction.tolist()}

#### Curl

curl -X 'POST' \
 'http://127.0.0.1:8000/predict/' \
 -H 'accept: application/json' \
 -H 'Content-Type: application/json' \
 -d '{
 "longitude": 119,
 "latitude": 35,
 "housing\_median\_age": 30,
 "total\_rooms": 2635,
 "total\_bedrooms": 500,
 "population": 1425,
 "households": 500,
 "median\_income": 3.87,
 "ocean\_proximity": "INLAND"
}'

#### **Request URL**



#### Let's deploy with FastAPI

from fastapi import FastAPI
import joblib
import pandas as pd
from pydantic import BaseModel

# Load the data processing pipeline
full\_pipeline = joblib.load('full\_pipeline.joblib')

# Load the ML model
model = joblib.load('best\_model.joblib')

# Launch API
app = FastAPI()

```
class HousingData(BaseModel):
```

longitude: float = 119.0
latitude: float = 35.0
housing\_median\_age: float = 30.0
total\_rooms: float = 2635.0
total\_bedrooms: float = 500.0
population: float = 1425.0
households: float = 500.0
median\_income: float = 3.87
ocean proximity: str = "INLAND"

#### @app.post("/predict/")

def predict(housing\_data: HousingData):
 # Convert the input data to a DataFrame
 data\_dict = housing\_data.dict()
 df = pd.DataFrame([data dict])

# Use the loaded pipeline to transform the data
prepared\_data = full\_pipeline.transform(df)

# Predict using the loaded model
prediction = model.predict(prepared\_data)
return {"prediction": prediction.tolist()}

#### District Pricing <sup>042</sup> <sup>045 33</sup>

/openapi.json

Deploying the pipeline in chapter 2 of Hands-on ML





### **Current challenges of ML**





## **Current challenges of ML**

Explainability


- Explainability
- Ethical and Social Implications





- Explainability
- Ethical and Social Implications
- Resource Efficiency





- Explainability
- Ethical and Social Implications
- Resource Efficiency
- Security and Privacy





- Explainability
- Ethical and Social Implications
- Resource Efficiency
- Security and Privacy
- Lack of "world model": machines are stupid





- Explainability
- Ethical and Social Implications
- Resource Efficiency
- Security and Privacy
- Lack of "world model": machines are stupid
- Multimodal learning
  - Type of data: graphs, images, text, ...
  - Joint representation and decision making

