Lecture 2: Classical methods and models in MLE with Scikit

> Tomáš Báča

Python Environment

Fanue

Learn

D-----

Preprocessi

Metric

Realwork

Reference:

Classical methods and models in MLE with Scikit Learn BECM33MLE — Machine Learning Engineering

Dr. Tomáš Báča

Multi-Robot Systems group, Faculty of Electrical Engineering Czech Technical University in Prague









Python Virtual Environment

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča Python

Environment

Panda

Scikit Learn

Regression nineline

Metrics

Realworl example

References

- always recommend to work in a Virtual Environment
- it stores copies of the particular dependencies locally
- one solution to a dependency hell

Dependencies (requirements.txt)

```
numpy==2.3.3
matplotlib==3.10.6
scikit-learn==1.7.2
pandas==2.3.2
jupyter
jupyterlab-vim
```

 the version is fixed for some of the packages to make the code within the examples work in the future

Setting up python environment (Ubuntu 24.04)

```
# install the python3's virutal environment
sudo apt get install python3-venv

# create the virtual environment
python3 -m venv python-env

# activate the environment
source ./python-env/bin/activate

# install dependencies manually
pip install numpy ...
```

or install dependencies in bulk

```
# install deps from requirements.txt
python3 -m pip install -r requirements.txt
```

Jupyter Notebook & Jupyter Lab

Lecture 2: Classical methods and models in MLE with Scikit

> Tomáš Báča

Python Environment

Panda:

Scikit Learn

Toy dataset Regression pipeline

Metrics

Realworld example

References

- web-interface with an editor
- allows execution of portions of your code
- remembers the state of variables
- remembers where you left it
- embedded visualization
- well-suited for learning and prototyping
- (bonus) vim-like key bindings
- Jupyter Notebook
 - single document, simple
- Jupyter Lab
 - multiple documents at once
 - similar to an IDE
 - better file management

Running Jupyter Notebook

```
# 1. activate your python env
source ./python-env/bin/activate

# 2. run
jupyter notebook
```

3. check your browser (localhost:8888)

Running Jupyter Lab

```
# 1. activate your python env
source ./python-env/bin/activate

# 2. run
jupyter-lab
```

check your browser (localhost:8888)

Jupyter Notebook & Jupyter Lab

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Python Environment

Panda

Learn

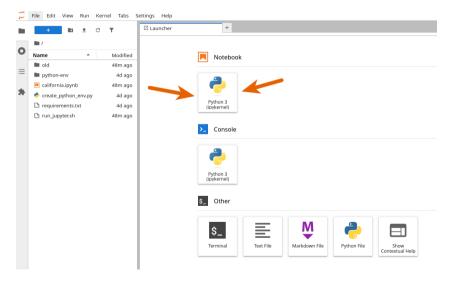
Regressio

Preprocess

. . .

Realwor

Reference



Jupyter Notebook & Jupyter Lab

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Python Environment

Panda

Scikit Learn

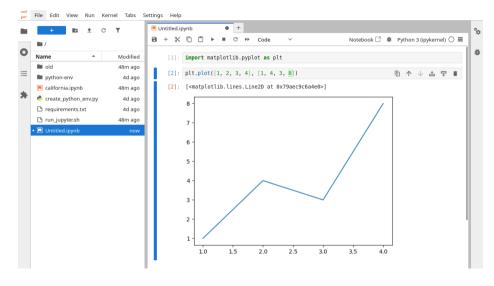
Regressio

Preprocess

Metri

Realworle example

Reference



Pandas — manipulating data in Python

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš

Báča Python Environ-

Pandas

Callele

Learn

Pogression

pipeline Preproces

Metric

Realworl example

References

Python library

- manipulation of tabular data
- "Excel in Python"
- can load common formats (CSV, Excel, SQL database, JSON, ...)
- suitable for:
 - data cleaning and processing
 - merging and joining of datasets
 - visualization (Matplotlib integration)
 - handles temporal data
 - integrates with Numpy and SkLearn

Pandas dataframe

- the main workhorse of Pandas
- 2D, size-mutable data structure

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25
20635	1.5603	25.0	5.045455	1.133333	845.0	2.560606	39.48	-121.09
20636	2.5568	18.0	6.114035	1.315789	356.0	3.122807	39.49	-121.21
20637	1.7000	17.0	5.205543	1.120092	1007.0	2.325635	39.43	-121.22
20638	1.8672	18.0	5.329513	1.171920	741.0	2.123209	39.43	-121.32
20639	2.3886	16.0	5.254717	1.162264	1387.0	2.616981	39.37	-121.24

20640 rows × 8 columns

Scikit Learn

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča Python

Panda

Scikit

Learn Toy datase

Regression pipeline

Matrica

Realworl example

Reference

- free and open source Python library
- implements many classical ML methods suitable for work with small data
- perfect for introduction into practical ML
- elegant interface that leads to clean and simple code
- plenty of examples and community support
- https://scikit-learn.org/stable/user_guide.html

Scikit learn can do

- Classification
 - NN, SVM, random forests, logistic regression, ...
- Regression
 - gradient boosting, random forests, logistic regression, ...
- Clustering
 - k-Means, hierarchical clustering, ...
- Dimensionality reduction
 - PCA, linear & nonlinear manifold learning
- Model selection, hyper parameter optimization
 - grid search, cross validation, metrics, ...
- Data preprocessing
 - data preprocessing, feature extraction, ...

Importing a toy dataset — California housing

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Python Environ

Panda

Scikit

Toy dataset

Regression

reproce

Realworle example

References

data matrix:

- median house age
- median income
- avg number of rooms
- avg number of bedrooms
- population in census block
- occupancy
- latitude
- longitude
- target value:
 - \$ house price USD

Importing a dataset

```
from sklearn.datasets import
    fetch_california_housing
from pandas import pd

housing = fetch_california_housing(as_frame=True)
```

What is in the dataset?

```
# plaintext print
print(housing.data)
print(housing.target)

# print in jupyter notebook
housing.data
housing.target
```

California housing dataset

Lecture 2: Classical methods and models in MLF with Scikit Learn

Tomáš

Báča

Toy dataset

Data	matrix

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25
20635	1.5603	25.0	5.045455	1.133333	845.0	2.560606	39.48	-121.09
20636	2.5568	18.0	6.114035	1.315789	356.0	3.122807	39.49	-121.21
20637	1.7000	17.0	5.205543	1.120092	1007.0	2.325635	39.43	-121.22
20638	1.8672	18.0	5.329513	1.171920	741.0	2.123209	39.43	-121.32
20639	2.3886	16.0	5.254717	1.162264	1387.0	2.616981	39.37	-121.24

20640 rows × 8 columns

Target value

0 1 2 3 4	4.526 3.585 3.521 3.413 3.422	
20635 20636 20637 20638 20639	0.781 0.771 0.923 0.847 0.894	

California housing dataset — targets

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Enviror

Panda: Scikit

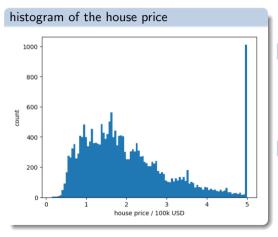
Toy dataset

Regression pipeline

Metrics

Realwor example

Reference



Plotting the histogram

```
plt.hist(housing.target, bins=100)
plt.xlabel("house price / 100k USD")
plt.ylabel("count")
```

- note that the maximum price is 5k USD
- the histogram has a clear spike in that price

California housing dataset — locations

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča Python

Enviro ment

Scikit

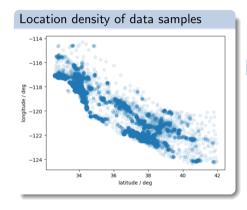
Toy dataset

pipeline Preproces:

Metrics

example

Reference



Plotting the histogram

```
plt.scatter(housing.data["Latitude"],
housing.data["Longitude"], alpha=0.1)
plt.xlabel("latitude / deg")
plt.ylabel("longitude / deg")
```

California housing dataset — loading for learning

Lecture 2: Classical methods and models in

MLE with Scikit Learn

Tomáš Báča

Enviror ment

Panda

Scikit

Toy dataset

pipeline Preprocessii

Metric

Realworl example

Reference

Loading X and y of the dataset

$return_X_y=True$

```
from sklearn.datasets import fetch_california_housing

X, y = fetch_california_housing(return_X_y=True)
```

```
\mathbf{X}
                  8.3252
                                41
                                                6.98412698 ...
                                                                  2.5555556
                  37.88
                              -122.23
                  8.3014
                                21.
                                                6.23813708 ...
                                                                  2.10984183
                  37.86
                              -122.22
                  7.2574
                                52.
                                               8.28813559 ...
                                                                  2.80225989
                  37.85
                              -122.24
               f 1.7
                                17.
                                                5.20554273 ...
                                                                  2.3256351
                  39.43
                              -121.22
                  1.8672
                                18
                                               5.32951289 ...
                                                                  2 12320917
                              -121.32
                  39.43
                  2.3886
                                16
                                                5.25471698 ...
                                                                  2.61698113
                  39.37
                              -121.24
```

```
Y
[4.526 3.585 3.521 ... 0.923 0.847 0.894]
```

The simplest regression model

that does something

using KNN to predict the price by finding

nearest neighbors in

the dataset and averaging their price

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Python Enviror

Pand

Scikit

Toy datas

Regression pipeline

.

Realworle example

References

• the minimal example

```
from sklearn.datasets import fetch_california_housing
from sklearn.neighbors import KNeighborsRegressor

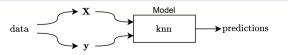
X, y = fetch_california_housing(return_X_y=True)

model = KNeighborsRegressor()

# training
model.fit(X, y)

# making predictions
prediction = model.predict(X)
```

"The ML model"



The simplest regression model

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Python Environ ment

Panda

Learn

Regression pipeline

Metrics

Realworl example

References

- the minimal example that does something
- using KNN to predict the price by finding nearest neighbors in the dataset and averaging their price

Possible flaws?

- we are using all the data to train the model at once
- we are doing our predictions on the same data we used for training

KNN regression pipeline

```
from sklearn.datasets import fetch_california_housing
from sklearn.neighbors import KNeighborsRegressor

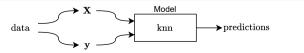
X, y = fetch_california_housing(return_X_y=True)

model = KNeighborsRegressor()

# training
model.fit(X, y)

# making predictions
prediction = model.predict(X)
```

"The ML model"



The simplest regression model — results

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Environment

Panda

Scikit Learn

Toy datase Regression

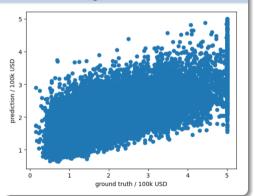
pipeline Preproces

Metrics

Realwor

Reference

Predictions vs. ground truth



Showing the predictions

```
plt.scatter(y, pred)
plt.xlabel("ground truth / 100k USD")
plt.ylabel("prediction / 100k USD")
```

- resemblance of a predictive behavior is there
- ideally, the data would form a line with slope = 1
- see the pattern caused by the capped house price at 500k USD?
- far from ideal, but this is a good start

The simplest regression model — improving?

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Python Enviror ment

Pand

Scikit

Toy data:

Regression pipeline

Matelan

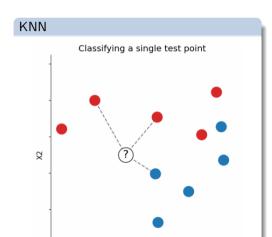
Realworl

References

- our data has varying physical quantities in different dimensions
- one axis is in feet², other in degrees, other in multiples of 100k USD

Solution?

- scaling data such that their distribution is similar
- many ways how to do that, but let's start with a simple practical example



X1

The simplest regression model — creating a pipeline

- Lecture 2: Classical methods and models in MLE with Scikit
- Learn Tomáš Báča

Python Enviror ment

Pand

Scikit Learn

Regression

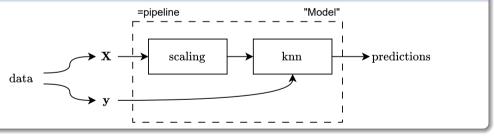
Preprocessi

Realwor

References

- let's add a scaling block in front of our KNN regressor
- with two block, we are already forming a pipeline

The simplest pipeline diagram



Adding preprocessing and pipeline

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Python Enviror ment

Panda:

Scikit

Toy datas

Regression pipeline

i reprocessi

Realworl

References

 the Pipeline has the same interface as the KNN regressor

```
• .fit(X, y)
```

.predict(X)

KNN regression pipeline

```
from sklearn.datasets import fetch_california_housing
   from sklearn.neighbors import KNeighborsRegressor
   from sklearn.preprocessing import StandardScaler
   from sklearn.pipeline import Pipeline
   X, v = fetch_california_housing(return_X_v=True)
   model = Pipeline([
     ("scaler", StandardScaler()),
     ("predictor", KNeighborsRegressor()),
   1)
   # training
   model.fit(X, v)
   # making predictions
16
   prediction = model.predict(X)
```

The simplest regression model — results

Lecture 2: Classical methods and models in MLE with Scikit

Learn

Tomáš Báča

Enviro

Panda

Scikit Learn

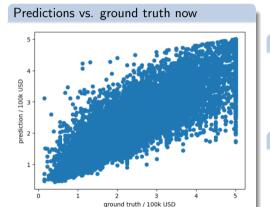
Regression

Preproces

Metrics

Realwor

Reference



Showing the predictions

```
plt.scatter(y, pred)
plt.xlabel("ground truth / 100k USD")
plt.ylabel("prediction / 100k USD")
```

- better performance than before
- as simple to use as before

The simplest regression model — results

Lecture 2: Classical methods and models in MLE with

Learn Tomáš Báča

Python Environment

Panda

Scikit

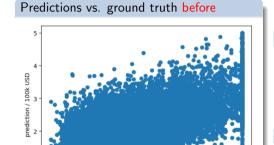
Toy data

Regression pipeline

Metrics

Realwo example

Reference



ground truth / 100k USD

Showing the predictions

```
plt.scatter(y, pred)
plt.xlabel("ground truth / 100k USD")
plt.ylabel("prediction / 100k USD")
```

- better performance than before
- as simple to use as before

Problems with current approach?

Lecture 2: Classical methods

models in MLE with Scikit Learn

Tomáš Báča

Pytho Environment

Panda

Scikit Learn

Toy data:

Regression pipeline

reprocess

Realwor

Reference

1. model parameters

• our model might have parameters that might need changing to improve "our performance"

Let's add n_neighbors=1

```
model = Pipeline([
("scaler", StandardScaler()),
("predictor", KNeighborsRegressor(
n_neighbors=1)),
]
...
```

Problems with current approach?

Lecture 2: Classical methods and

models in MLE with Scikit Learn

Tomáš Báča

Python Enviror ment

Panda

Scikit

Toy datase Regression

pipeline

Metrics

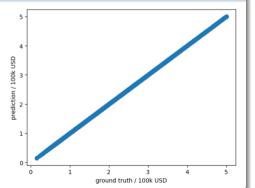
example

Reference:

1. model parameters

• our model might have parameters that might need changing to improve "our performance"

Predictions vs. ground truth for 1 neighbor



Let's add n_neighbors=1

This leads to problem #2

our performance evaluation method in not very good

Problems with current approach?

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

ment

Pandas

Learn Toy datas

Regression

Metric

Realwor example

References

2. our evaluation method is wrong

- we are checking performance of the model using the same data we trained on
- this makes selecting the right parameters impossible (for other data than the one we are training on)

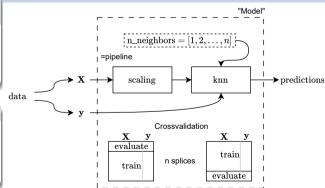
Grid search / hyperparameter optimization

- selecting parameter of block that can not be trained by the block itself
- iterates over all combinations of parameters

Crossvalidation

- splicing data to several pieces
- training on majority of the data while testing on the rest
- then swapping the splice

Using grid search with crossvalidation



GridSearch with Crossvalidation

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Environment

Panda

Scikit Learn

Regression

Prepro

Realwo

example

References

Observations

- GridSearchCV has the same interface
 - .fit(X, y)
 - .predict(X)
- whole parts of the pipeline can be disabled/enabled

Checking the results

in Jupyter notebook

```
pd.DataFrame(model.cv_results_)
```

Introducing GridSearchCV

```
. . .
   from sklearn.model selection import
        GridSearchCV
   pipeline = Pipeline([
      ("scaler", StandardScaler()),
      ("predictor", KNeighborsRegressor()),
8
   1)
g
   model = GridSearchCV(
     pipeline,
     param_grid={'predictor_n_neighbors':
                   Γ1. 2. 3. 4. 5. 6]
14
     cv = 3
16
   model.fit(X, y)
18
20
      . . .
```

GridSearch with Crossvalidation

Lecture 2: Classical methods and models in MLE with Scikit

Learn

Tomáš Báča Python

Panda

0.11.11

Learn

Regression

Preproce

Metric

Realwor example

References

Grid search results

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_predictorn_neighbors	params	split0_test_score	split1_test_score	split2_test_score	mean_test_score	std_test_score	rank_test_score
0	0.009950	0.000354	0.218222	0.016866	1	{'predictor_n_neighbors': 1}	0.324068	0.334830	0.323371	0.327423	0.005245	6
1	0.009837	0.000220	0.246559	0.021540	2	{'predictor_n_neighbors': 2}	0.468788	0.503457	0.424388	0.465544	0.032361	5
2	0.009659	0.000020	0.269420	0.019367	3	{'predictor_n_neighbors': 3}	0.518547	0.543340	0.473595	0.511827	0.028867	4
3	0.009642	0.000061	0.275874	0.021485	4	{'predictor_n_neighbors': 4}	0.540323	0.564974	0.499827	0.535041	0.026857	3
4	0.010023	0.000212	0.289105	0.020307	5	{'predictor_n_neighbors': 5}	0.551149	0.579313	0.511781	0.547414	0.027696	2
5	0.009678	0.000049	0.305267	0.028973	6	{'predictor_n_neighbors': 6}	0.558435	0.586185	0.521134	0.555251	0.026652	1

- the best KNN model is the one with 6 neighbors (makes sense)
- the model won with a mean_test_score of 0.555
 - what score is it
 - how is it calculated?
 - how do we change it?
 - we will get back to this

Data preprocessing — more about sklearn's transformers

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Enviror

Panda

Scikit

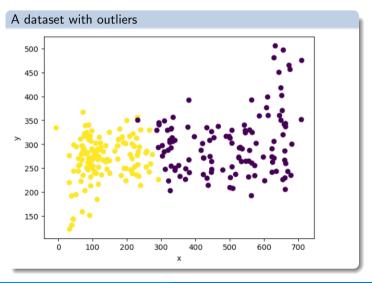
Toy datas

Preprocessing

Realworld example

References

- two classes
- should be nicely separable
- outliers near the edges
- data in arbitrary range on both axes



Data preprocessing — more about sklearn's transformers

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

ment Pandas

Scikit

Toy datase Regression

Preprocessing

Realworle

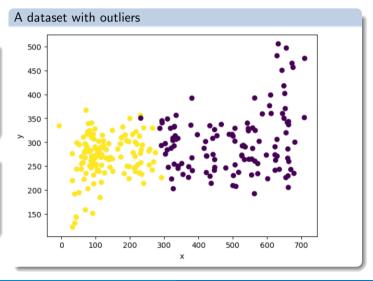
References

Observations

- two classes
- should be nicely separable
- outliers near the edges
- data in arbitrary range on both axes

Let's try standard scaler

- what does it even do?
- how will it perform?
- will it be influenced by the outliers?



Standard scaler

$$z = \frac{x - \mu}{\sigma},$$

where

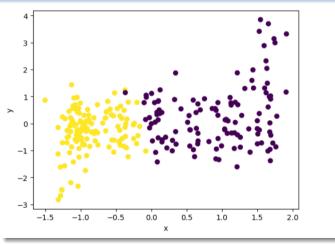
$$\mu = \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (x_i - \mu)^2}$$

Observations

- outliers still visible
- axes' range is near zero
- the shape of the data is similar

Dataset after transformation



Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Environ ment

Panda

.

Learn

Regressio

Preprocessing

Metrics

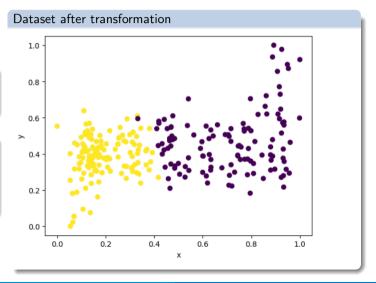
example

References



$$z = \frac{x - \min}{\max - \min},$$

- outliers still visible
- axes' range is between 0 and1
- distribution is preserved



Data preprocessing — Quantile transformer

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Enviror ment

Panda

Scikit

Learn Toy data

Regression pipeline

Preprocessing

Realworl

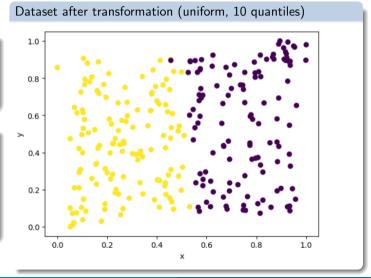
Defenses

References

Quantile transformer

- reshapes the data to have uniform (or normal) distribution
- robust scaler

- uniform output (default)
- outliers are not prominent anymore
- non-linear transformation
 - will break linear correlations between features



Data preprocessing — Quantile transformer

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Python Environ ment

Panda

Scikit

Toy datas

Preprocessing

Realwork

References

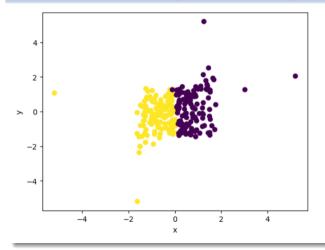
Quantile transformer

- reshapes the data to have uniform (or normal) distribution
- robust scaler

Observations

- uniform output (default)
- outliers are not prominent anymore
- non-linear transformation
 - will break linear correlations between features

Dataset after transformation (normal, 10 quantiles)



Linear separability

Lecture 2: Classical methods and models in MLE with

> Scikit Learn Tomáš Báča

Python Enviror ment

Pand

Scikit

Toy datas

Regression pipeline

Preprocessing Metrics

Realwo

example

References

Problem?

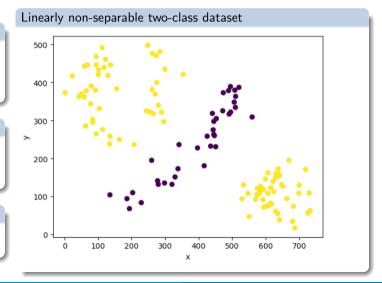
- classes are not linearly separable
- a.k.a., a single separating hyperplane can not be found

Solution?

- use more complex classifier
 - nonlinear
 - e.g., boosting
 - decision tree

Easier solution?

- use more complex classifier
- lift feature to new dimension



Linear separability

Lecture 2: Classical methods and models in MLF with

Scikit Learn Tomáš Báča

Python Environment

Pand

Scikit Learn

Toy datas Regression

Preprocessing

Realwo

example

References

Problem?

- classes are not linearly separable
- a.k.a., a single separating hyperplane can not be found

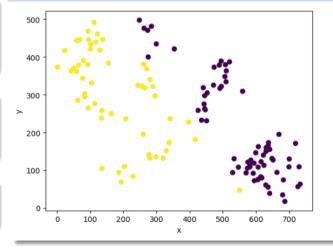
Solution?

- use more complex classifier
 - nonlinear
 - e.g., boosting
 - decision tree

Easier solution?

- use more complex classifier
- lift feature to new dimension





Polynomial features

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Pythor

Panda

Scikit

Toy data

Regression pipeline
Preprocessing

Metrics

Realwor example

References

Lifting features using polynomial interactions

• example for order = 2

$$\begin{aligned} \mathbf{x} &= \left[x_1, x_2\right]^\intercal \\ \mathbf{x}_{\mathsf{new}} &= \Phi(\mathbf{x}) = \left[1, x_1, x_2, x_1^2, x_1 x_2, x_2^2\right]^\intercal \end{aligned}$$

Observations

- leads to simple classifier, which can be robust
- e.g., all favourite SVM
- makes computations more expensive during production
- can be simplified with the kernel trick

Using polynomial features

```
. . .
from sklearn.preprocessing import
     PolynomialFeatures
from sklearn.linear model import
     LogisticRegression
pipeline = Pipeline([
  ("scale", PolynomialFeatures()),
  ("classifier", LogisticRegression(
       max_iter=1000, class_weight="balanced
       ")),
1)
      pipeline.fit(X, y).predict(X)
  . . .
```

Polynomial features

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Environment

Panda

0.00

Learn Toy datas

Regression

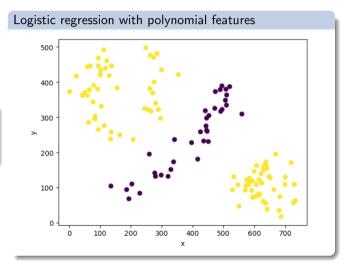
Preprocessing

Metrics

Realwor

Reference

- simple to use
- just like any other Sklearn transformer



Polynomial features

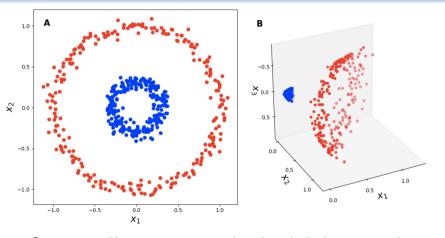
Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Preprocessing



Illustration of lifting from 2D to 3D



Source: https://gregorygundersen.com/blog/2019/12/10/kernel-trick/

Encoding text and categorical attributes

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Panda

Learn
Toy data

Regression pipeline Preprocessing

Metrics

Realwor example

References

How to encode text?

- added uunique ID in 1D space might cause problems
- better way?
 - give each class a unique unit vector in its own dimension
 - this this creates as many dimensions (+1) as the # of classes

Even better way?

- learn unique vectors in lower dimensional space
- called embedding
- heart (mouth) of transformers

Using OneHot encoder

```
from sklearn.preprocessing import OneHotEncoder

X = [['red'], ['green'], ['blue'], ['red']]

encoder = OneHotEncoder(handle_unknown="ignore")
encoder.fit(X)

new_X = encoder.transform(X).todense()
```

Inverse:

```
encoder.inverse_transform([[0, 1, 0]])
```

```
[['green']]
```

Objects in Scikit

Lecture 2: Classical methods and models in MLE with Scikit

Learn

Tomáš Báča

Pandas

Scikit Learn

Regression pipeline

Preprocessing Metrics

Realwor example

References

Estimator

- implements .fit()
- only consumes data
- e.g., some statistics collector

Transformer

- implements .fit(), .transform()
- scalers
- dimensionality reducers (PCA)
- imputers (handle missing values)

Predictor

- implements .fit(), .predict()
- regressors, classifiers

Pipeline

is a Predictor

GridSearchCV

• is a Predictor

Binary classification metrics — Precision, Accuracy and Recall

Lecture 2: Classical methods and models in

MLE with Scikit Learn

Tomáš Báča

Python Enviror ment

Panda

Scikit

Toy dat

Regressi

Preprocess

Metrics

Realwo example

References

Precision

$$\mathsf{precision} = \frac{\mathsf{TP}}{\mathsf{TP} + \mathsf{FP}}$$

 given I mark the sample as relevant, high often am I right

Recall

$$\mathsf{recall} = \frac{\mathsf{TP}}{\mathsf{TP} + \mathsf{FN}}$$

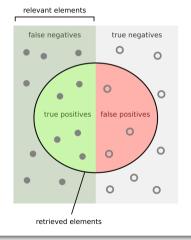
• did I marked all the relevant instances?

Accuracy

$$\mathsf{accuracy} = \frac{\mathsf{TP} {+} \mathsf{TN}}{\mathsf{TP} {+} \mathsf{TN} {+} \mathsf{FP} {+} \mathsf{FN}}$$

how precise am I in marking correctly both instances

${\sf Diagram\ of\ possible\ results}$



Metrics example

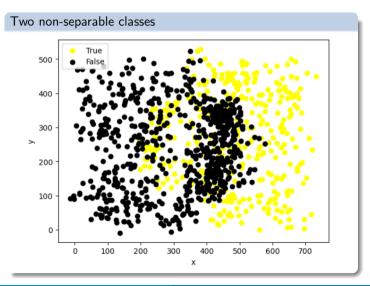
Lecture 2: Classical methods and models in MLF with Scikit Learn

Tomáš Báča

Metrics

Is there a single best classifier?

- no, it depends on our metrics
- do we care more about
 - precision?
 - accuracy?
 - recall?
- or some combination of these?



Logistic regression with default metrics

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Panda

Scikit

Toy data

Preprocessing

Metrics

Realwor example

References

Observations

- we are trying Logistic Regression
- we are varying the class weight
- we are searching for the best-performing predictor

Let's run this

- let's run it as it is
- the default score for logistic regression is Accuracy

Trying logistic regression

```
. . .
   from sklearn.linear_model import LogisticRegression
   from sklearn.model_selection import GridSearchCV
   pipeline = Pipeline([
      ("classifier", LogisticRegression(max_iter=1000.
          verbose=0)).
Q
   model = GridSearchCV(
10
      estimator=pipeline, cv=4, param grid={
      'classifier__class_weight': [{0: 5, 1: v} for v in
           range(1, 10)]},
14
   pred = model.fit(X, v).predict(X)
16
    . . .
```

Metrics example — Accuracy

Lecture 2: Classical methods and models in MLE with

Learn Tomáš Báča

Enviror ment

Panda

Scikit

Toy data

Regressio pipeline

Metrics

Realworl example

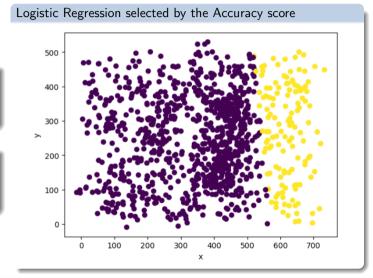
References

Observations

- the decision is *cut* at a reasonable compromise
- many FP as well as FN

Metrics on the dataset

- accuracy = 0.77
- precision = 0.97
- recall = 0.35



Logistic regression with Precision metrics

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Environ ment

Pand

Scikit Learn

Toy dataset Regression pipeline

Metrics

Realworl example

References

Observations

- we are trying Logistic Regression
- we are varying the class weight
- we are searching for the best-performing predictor

Scoring

- new scoring argument
 - we can add arbitrary scoring functions
 - all will be evaluated during the grid search
- new refit argument
 - the pipeline will be re-trained using parameters that won using the pre-defined score

Logistic regression precision

```
from sklearn.linear model import LogisticRegression
   from sklearn.model_selection import GridSearchCV
   from sklearn.metrics import precision_score,
3
        accuracy_score, recall_score
4
   pipeline = Pipeline([
     ("classifier", LogisticRegression(max_iter=1000,
6
          verbose=0)).
   model = GridSearchCV(
Q
     estimator=pipeline, cv=4, param_grid={
     'classifier__class_weight': [{0: 5, 1: v} for v in
           range(1, 10)]}.
     scoring={'precision': make_scorer(precision_score)
          . 'recall': make_scorer(recall_score). '
          accuracy': make scorer(accuracy score)}.
     refit='precision'.
14
15
   pred = model.fit(X, v).predict(X)
16
```

Metrics example — Precision

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča Python

ment Panda

Panda

Scikit Learn

Toy datase Regression

Metrics

Realwo

Reference

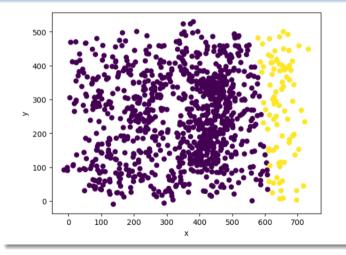
Observations

- the decision is cut at the very right
- minimizes false positives
- lot of missed true sample

Metrics on the dataset

- accuracy = 0.73
- precision = 1.0
- recall = 0.23

Logistic Regression selected by the Precision score



Logistic regression with Recall metrics

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Environ ment

Pand

Learn

Toy dataset Regression pipeline

Metrics

Realworl example

References

Observations

- we are trying Logistic Regression
- we are varying the class weight
- we are searching for the best-performing predictor

Scoring

- new scoring argument
 - we can add arbitrary scoring functions
 - all will be evaluated during the grid search
- new refit argument
 - the pipeline will be re-trained using parameters that won using the pre-defined score

Logistic regression recall

```
from sklearn.linear model import LogisticRegression
   from sklearn.model_selection import GridSearchCV
   from sklearn.metrics import precision_score,
3
        accuracy_score, recall_score
Λ
   pipeline = Pipeline([
     ("classifier", LogisticRegression(max_iter=1000,
6
          verbose=0)).
   model = GridSearchCV(
Q
     estimator=pipeline, cv=4, param_grid={
     'classifier__class_weight': [{0: 5, 1: v} for v in
           range(1, 10)]}.
     scoring={'precision': make_scorer(precision_score)
          . 'recall': make_scorer(recall_score). '
          accuracy': make scorer(accuracy score)}.
     refit='recall'.
14
15
   pred = model.fit(X, v).predict(X)
16
```

Metrics example — Recall

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

ment

Panda

Scikit Learn

Regression pipeline

Metrics

Realwor example

Reference

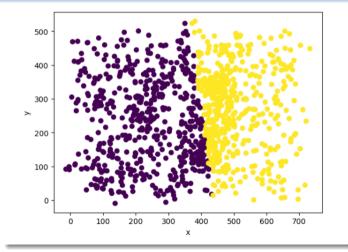
Observations

- the decision is cut at the very right
- minimizes false positives
- lot of missed true sample

Metrics on the dataset

- accuracy = 0.67
- precision = 0.52
- recall = 0.72

Logistic Regression selected by the Precision score



Storing a trained sklearn model

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Enviror ment

Panda

rand

Callde

Toy dat

Regressio

Preprocessi

Metrics

Realworl example

Reference

- as simple as dumping the object into a file
- alternatively, you can use pickle

Storing a model

```
from joblib import dump

dump(model, "my_model.joblib")

...
```

Loading a model

```
1
2
3
4
5 model = load("my_model.joblib")
6
7
```

Realworld example — Realtime classification of ionizing particles

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Python Environ ment

Panda

Toy datas

pipeline Preproces

Metric

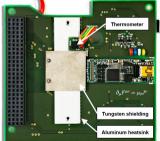
Realworld example

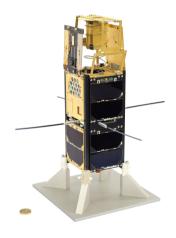
References

VZLUSAT-1 — The first Czech CubeSat [1], [2]

- Custom-designed embedded Timepix board and software [2]
- Onboard image processing, filtering, automated acquisition
- The longest-operating Czech(-oslovakian) sat. (2017–2023)







Timepix - Ionizing radiation dosimetry and imaging

Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

ment Panda

ranue

Learn

Regression

Metric

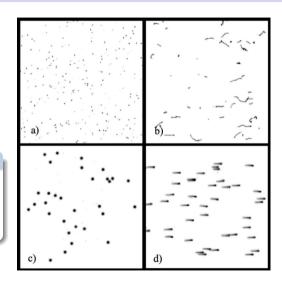
Realworld example

References

- real-time recording of the interaction of an incoming particle and the sensor
- machine learning was applied to distinguish particle types [3]
- no dark-current noise (the images are spotless besides the actual data)

Examples in the Figure

- photons (gamma)
- electrons (beta)
- Helium nuclei (alpha)
- protons



Random forests for particle track classification

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča Python

ment

Panda

Scikit

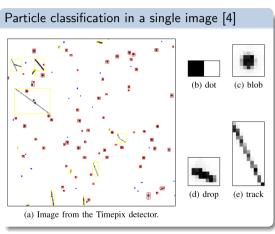
Toy data

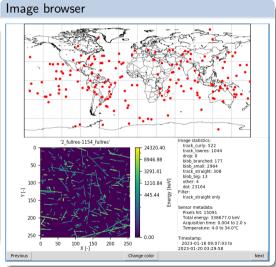
Regressio pipeline

Metri

Realworld example

Reference





Low-Earth orbit particle classification

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

Pythor Enviro ment

Panda

Scilit

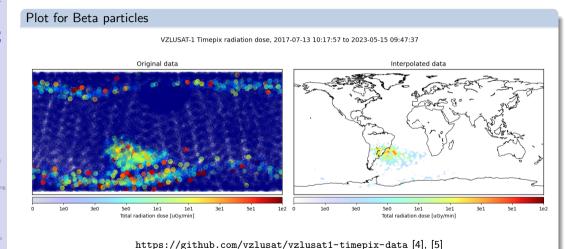
Learn Toy dat

Regression pipeline

Metri

Realworld example

Reference



Realtime classification onboard UAVs

Lecture 2: Classical methods and models in MLE with Scikit

Learn Tomáš Báča

ment

Panda

Learn Toy dat

Preproc

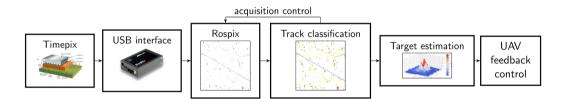
Metrics

Realworld example

References

Proof of concept for UAVs

- Real-time particle track classifier for mobile robots and drones.
- We have finished well-evaluated TACR grant to develop this tech.



[4] T. Baca, M. Jilek, P. Manek, P. Stibinger, V. Linhart, J. Jakubek, et al., "Timepix Radiation Detector for Autonomous Radiation Localization and Mapping by Micro Unmanned Vehicles," in 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems, IEEE, 2019, pp. 1–8

Realtime classification onboard UAVs

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Pythor Environment

Pand

Scikit Learn

Toy data Regressia

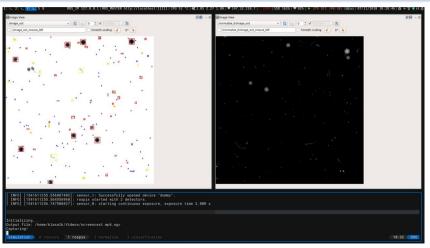
pipeline Preproce

N.A. and an

Realworld example

Reference

Realtime classification of ionizing particles in ROS



Video: https://youtu.be/8rZrbmDHG4E

Lecture 2: Classical methods and models in MLE with Scikit Learn

Realtime classification onboard UAVs

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Pythor Environment

Pand

Scikit

Toy data

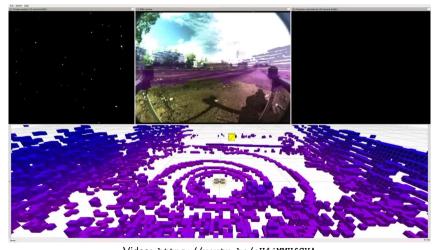
Regression pipeline

Metric

Realworld example

References

Endgame: realtime search for ionizing radiation sources



Conclusion

Classical methods and models in MLE with Scikit Learn

Lecture 2:

Tomáš Báča

Python Enviror ment

Pand:

Scikit

Learn

Regressie

Proprocess

Metric

Realworld example

Reference:

Thanks for your attention

References I

- Lecture 2: Classical methods and models in MLE with Scikit
- Tomáš Báča Python
- Panda
-
- Toy data
- Regression pipeline
- Preprocessir
- ivietrics
- Realworl example
- References

- [1] M. Urban, O. Nentvich, V. Stehlikova, T. Baca, V. Daniel, and R. Hudec, "VZLUSAT-1: Nanosatellite with miniature lobster eye X-ray telescope and qualification of the radiation shielding composite for space application," Acta Astronautica, vol. 140, pp. 96–104, 2017.
 - T. Baca, M. Platkevic, J. Jakubek, et al., "Miniaturized X-ray telescope for VZLUSAT-1 nanosatellite with Timepix detector," Journal of Instrumentation, vol. 11, no. 10, p. C10007, 2016.
- [3] T. Baca, M. Jilek, I. Vertat, et al., "Timepix in LEO Orbit onboard the VZLUSAT-1 Nanosatellite: 1-year of Space Radiation Dosimetry Measurements," Journal of Instrumentation, vol. 13, no. 11, p. C11010, 2018.
- [4] T. Baca, M. Jilek, P. Manek, et al., "Timepix Radiation Detector for Autonomous Radiation Localization and Mapping by Micro Unmanned Vehicles," in 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems, IEEE, 2019, pp. 1–8.
- [5] T. Baca, P. Stibinger, D. Doubravova, et al., "Gamma Radiation Source Localization for Micro Aerial Vehicles with a Miniature Single-Detector Compton Event Camera," in 2021 International Conference on Unmanned Aircraft Systems (ICUAS), IEEE, 2021, pp. 1–9.