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

> Tomáš Báča

Python Enviror ment

Pandas

Learn

Regression pipeline

Preprocessing

Metrics

example

Referenc

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







MULTI-ROBOT SYSTEMS GROUP



Tomáš Báča (CTU in Prague)

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

Santambar 20th 2026

1/4

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

Canada methoda and medida in MEE and Scient Learn METHODIC TO Allows Learning Engineering

On Travell Editor

Main Alexan Learning Engineering

Allows Learning Allows Learning Learnin

025-10-11

Python Virtual Environment

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

> Tomáš Báča

Python Environment

Pandas

Learn

Regression pipeline

Preprocessi

Realworl

example

- 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

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

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
```

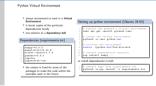
Tomáš Báča (CTU in Prague)

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

Python Environment

Python Enviro

Python Virtual Environment



- there are other options to make sure you dependencies are met and are not colliding with other projects:
 - Conda package manager,
 - Docker,
 - Singularity / Apptainer container system,
 - Nix package manager.

Jupyter Notebook & Jupyter Lab

Lecture 2 Classical methode and models in MIE with Scikit Learn

> Tomáš Ráča

Python Environment

Regression

- 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
3
Λ
   # 2. run
   jupyter notebook
```

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)

Tomáš Báča (CTU in Prague)

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

Jupyter Notebook & Jupyter Lab

Python Environment

Jupyter Notebook & Jupyter Lab

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

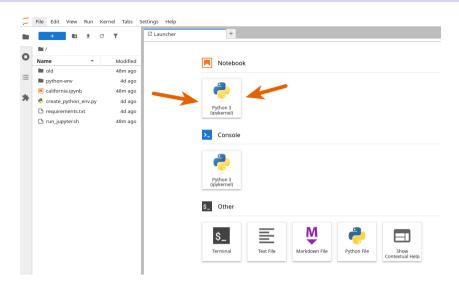
- You can, of course, run your python scripts manually. However, this makes development and tweaking of the algorithms
- In practice, you can design the first part of the architecture in Jupyter and then move to more traditional program architecture.

Jupyter Notebook & Jupyter Lab



ment

Regression Preprocessing



Tomáš Báča (CTU in Prague)

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

Jupyter Notebook & Jupyter Lab

Jupyter Notebook & Jupyter Lab



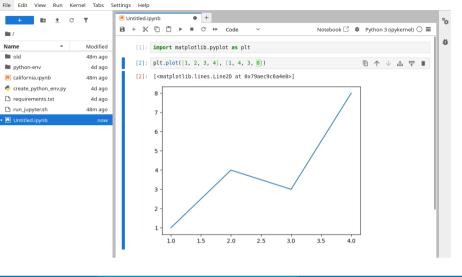
Python Environment

Tomáš Báča (CTU in Prague) Lecture 2: Classical methods and models in MLE with Scikit Learn

Python Environment

2025-10-11

Jupyter Notebook & Jupyter Lab





- use <tab> for code completion
- use shift+<tab> for explanation of objects and their properties
- use ctrl+enter to evaluate current cell

Pandas — manipulating data in Python

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

> Tomáš Báča

Environment

Pandas

Toy datas

Regression pipeline Preprocessi

Metrics

example

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

Fomáš Báša (CTII in Prague)

acture 2: Classical methods and models in MLE with Scikit Lear

September 30th, 2025

CLAS

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

Pandas — manipulating data in Python

Scikit Learn

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

Tomáš Báča

Environ

Panda

Scikit Learn

Regression pipeline Preprocess

WELFICS

example

- 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, ...

Tomás Báfa (CTU in Prague)

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

Scikit

Importing a toy dataset — California housing

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

Tomáš Báča

Python Environ

Panda:

Learn

Toy dataset Regression

Preprocessir

ivietrics

example

data matrix:

- median house age
- median income
- avg number of rooms
- avg number of bedrooms
- population in census block
- occupancy
- latitude
- longitude
- target value:
 - \$ $\frac{\text{house price}}{100000}$ 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
```



• some other datasets: https://scikit-learn.org/stable/datasets/toy_dataset.html

California housing dataset

Lecture 2: Classical methods models in MLE with Scikit Learn

Tomáš Báča

Toy dataset

Regression

Preprocessing

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 r	ows × 8 c	olumns						

Target value	9	
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	

Tomáš Báča (CTU in Prague) Lecture 2: Classical methods and models in MLE with Scikit Learn Scikit Learn Toy dataset California housing dataset



• this dataset is small enough for quick testing and tuning

California housing dataset — targets

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

Tomáš Báča

Environ

i anda.

Learn Toy dataset

Regression pipeline Preprocess

Metric

example

histogram of the house price

Plotting the histogram

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

Observation?

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

Tomáš Báča (CTU in Prague)

and the second second

California housing dataset — targets

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

Scikit Learn

Toy dataset

California housing dataset — targets

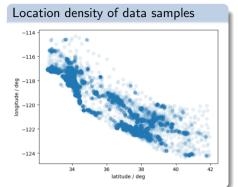
2025-10-11

California housing dataset — locations

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

Báča

Toy dataset



Plotting the histogram

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

Tomáš Báča (CTU in Prague)

Scikit Learn Toy dataset

California housing dataset — locations

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

California housing dataset — loading for learning

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

Báča Python Environment

Pandas

Learn Toy dataset

Regression pipeline Preprocessin

Metrics

Realworle example

```
Realworld
```

```
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}
                                                6.98412698 ...
                              -122.23
                 8.3014
                                               6.23813708 ...
                                                                  2.10984183
                              -122.22
                 37.86
                  7.2574
                                               8.28813559 ...
                                                                  2.80225989
                                52.
                  37.85
                              -122.24
                              17.
-121.22
                  1.7
                                               5.20554273 ...
                                                                  2.3256351
                  39.43
                                               5.32951289 ...
                                                                  2.12320917
                  1.8672
                                18.
                              -121.32
                  39.43
                                               5.25471698 ...
                                                                  2.61698113
                              -121.24
                                           11
```



Tomáš Báča (CTU in Prague)

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

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

Scikit Learn

Toy dataset

California housing dataset — loading for learning



- $\bullet \ X$ is a 2D matrix with individual attributes in the columns
- y is a 1D vector
- the scikit API is build to work with corresponding X and y

The simplest regression model

the minimal example

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 Learn

Tomáš Báča

Regression pipeline

```
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()
8
   # training
9
   model.fit(X, y)
   # making predictions
   prediction = model.predict(X)
```

"The ML model" Model data knn → predictions

Tomáš Báča (CTU in Prague)

Scikit Learn Regression pipeline

The simplest regression model

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

- The model can be easily swapped for other model, e.g., the LinearRegression. All the sklearn models have the same interface.
- This approach has may other flaws besides the ones mentioned here. We are going to get to some of them :-).

The simplest regression model

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?

Lecture 2 Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Regression pipeline

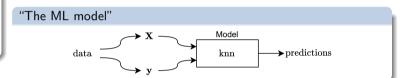
2025-10-11

 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()
8
   # training
9
   model.fit(X, y)
   # making predictions
   prediction = model.predict(X)
```



Tomáš Báča (CTU in Prague) Lecture 2: Classical methods and models in MLE with Scikit Learn Scikit Learn Regression pipeline The simplest regression model

- The model can be easily swapped for other model, e.g., the LinearRegression. All the sklearn models have the same interface.
- This approach has may other flaws besides the ones mentioned here. We are going to get to some of them :-).

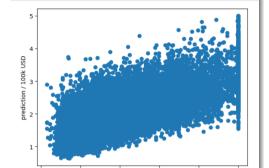
The simplest regression model — results

Predictions vs. ground truth

Lecture 2: Classical methods models in MLE with Scikit Learn

Tomáš Báča

Regression pipeline



ground truth / 100k USD

Showing the predictions

```
plt.scatter(y, pred)
2
  plt.xlabel("ground truth / 100k USD")
  plt.ylabel("prediction / 100k USD")
```

Observations

- 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

Tomáš Báča (CTU in Prague)

Scikit Learn

Regression pipeline

The simplest regression model — results

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



The *simplest* regression model — improving?

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

Tomáš Báča

Python Environ ment

Pandas

Learn

Regression pipeline

Preprocess

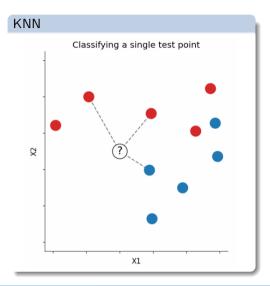
Deelissed

example

- 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





• KNN uses Minkowski metric with power = 2 (= 12 norm). If the attributes have different scales, then the metric will not be influenced evenly by the different attributes.

The simplest regression model — creating a pipeline

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

Environment

Panda

Learn Toy data

pipeline Proprocessis

Preprocessi

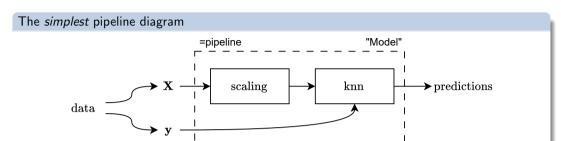
Metrics

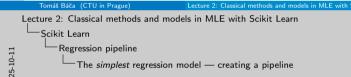
Toy dataset Regression

Realworld example

let's add a scaling block in front of our KNN regressor

• with two block, we are already forming a pipeline





September 30th, 2025

The simplest regression model—creating a pipelines

* We start a coding time in from of the 50th regions

* What the a coding time in from of the 50th regions

* What the simple pipelines degree

The simple pipelines degree

The simple pipelines degree

* T

Adding preprocessing and pipeline

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

Tomáš Báča

Python Environ-

Pandas

Learn

Regression pipeline

Preprocessi

Metrics

Realworld

D 6

 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, y = fetch_california_housing(return_X_y=True)
7
8
   model = Pipeline([
9
      ("scaler", StandardScaler()),
      ("predictor", KNeighborsRegressor()),
10
   1)
12
   # training
   model.fit(X, y)
14
15
16
   # making predictions
   prediction = model.predict(X)
```

Tomáš Báča (CTU in Prague)

Lacture 2: Classical methods and models in MLE with Scikit Lear

September 30th, 2025

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

gression pipeline

Adding preprocessing and pipeline

The simplest regression model — results

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

Tomáš Báča

Environment

.

Toy data

Regression pipeline Preprocess

Metrics

example

2025-10-11

Predictions vs. ground truth now

Showing the predictions

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

Observations

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

Tomáš Báča (CTU in Prague)

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

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

Scikit Learn

Regression pipeline

The simplest regression model — results



- The results already looks better with just the StandardScaler.
- But what does the standard scaler do?
 - more on that later, but in the meantime:
 - the standard scaler subtracts the mean of the data and then divides the data by the standard deviation.

The simplest regression model — results

Predictions vs. ground truth before

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

Tomáš Báča

Environment

randas

Learn

Regression pipeline Preprocessi

Metrics

example

2025-10-11

Showing the predictions

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

Observations

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

Tomáš Báča (CTU in Prague)

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

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

Scikit Learn

Regression pipeline

The simplest regression model — results

ground truth / 100k USD

The simplest regression model — results

Professions w. ground such before

| Description | Descript

- The results already looks better with just the StandardScaler.
- But what does the standard scaler do?
 - more on that later, but in the meantime:
 - the standard scaler subtracts the mean of the data and then divides the data by the standard deviation.

Problems with current approach?

Lecture 2: Classical methode models in MLE with Scikit Learn

> Tomáš Báča

Regression pipeline

```
1. model parameters
```

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

Let's add n_neighbors=1

```
3
   model = Pipeline([
     ("scaler", StandardScaler()),
4
     ("predictor", KNeighborsRegressor(
          n_neighbors=1)),
   1)
6
```

Tomáš Báča (CTU in Prague) Lecture 2: Classical methods and models in MLE with Scikit Learn Lecture 2: Classical methods and models in MLE with Scikit Learn Scikit Learn Let's add n_neighbors=1 Regression pipeline Problems with current approach?

- When setting the KNN to use only 1 neighbor, the methods suddenly produces perfect predictions.
- That is caused by the KNN being asked to predict using the same samples it remembered during the training.
- So now we know we have a problem:
 - There might be hyper parameters that influence the performance we don't know which values to pick.
 - We can not use the same dataset for training as well as for evaluation, that gives us false sense of success.

Problems with current approach?

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

Tomáš Báča

Python Environment

Pandas

Learn

Regression

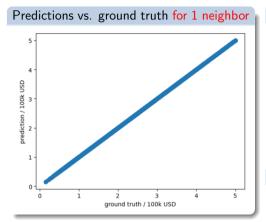
Matrics

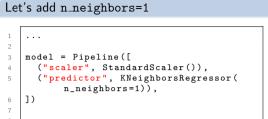
example

2025-10-11

1. model parameters

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

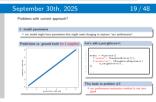




This leads to problem #2

• our performance evaluation method in not very good

Lecture 2: Classical methods and models in MLE with Scikit Learn
Scikit Learn
Regression pipeline
Problems with current approach?



- When setting the KNN to use only 1 neighbor, the methods suddenly produces perfect predictions.
- That is caused by the KNN being asked to predict using the same samples it remembered during the training.
- So now we know we have a problem:
 - There might be hyper parameters that influence the performance we don't know which values to pick.
 - We can not use the same dataset for training as well as for evaluation, that gives us false sense of success.

Problems with current approach?

Lecture 2: Classical methode models in MLE with Scikit Learn

> Tomáš Báča

Regression pipeline

2. our evaluation method is wrong

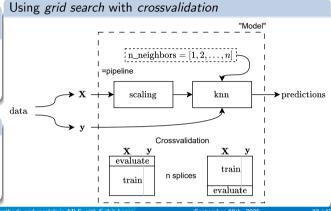
- we are checking performance of the model using the same data we trained on
- this makes selecting the right parameters impossible (for other data than the one we are training on)

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

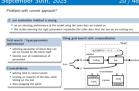


Tomáš Báča (CTU in Prague)

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

Scikit Learn Regression pipeline

Problems with current approach?



GridSearch with Crossvalidation

Regression pipeline

GridSearch with Crossvalidation

```
Lecture 2:
Classical
                                                                     Introducing GridSearchCV
methode
 and
models in
MLE with
 Scikit
            Observations
 Learn
                                                                         from sklearn.model_selection import
                                                                     3
 Tomáš
                                                                               GridSearchCV

    GridSearchCV has the same interface

 Ráča
                                                                     4

    .fit(X, y)

                                                                         pipeline = Pipeline([
                                                                     5
                   • .predict(X)
                                                                            ("scaler", StandardScaler()),
                                                                     6

    whole parts of the pipeline can be

                                                                     7
                                                                            ("predictor", KNeighborsRegressor()),
                                                                         1)
                                                                     8
                 disabled/enabled
                                                                     9
                                                                    10
                                                                         model = GridSearchCV(
                                                                           pipeline,
            Checking the results
                                                                           param_grid={'predictor__n_neighbors':
Regression
                                                                                           [1, 2, 3, 4, 5, 6]
pipeline

    in Jupyter notebook

                                                                                          },
                                                                    14
                                                                            cv=3
                pd.DataFrame(model.cv_results_)
                                                                         )
                                                                    16
                                                                         model.fit(X, y)
                                                                    18
                                                                    19
                                                                    20
          Tomáš Báča (CTU in Prague)
                                           Lecture 2: Classical methods and models in MLE with Scikit Learn
       Lecture 2: Classical methods and models in MLE with Scikit Learn
           Scikit Learn
```

 whole parts of the pipeline can be skipped by replacing them using the passthrough step in the param_grid: 'scaler': [StandardScaler(), 'passthrough']

Checking the results

GridSearch with Crossvalidation

mean fit time std fit time mean score time std score time param predictor in neighbors

0.016866

0.021540

0.019367

0.021485

0.020307

0.028973

0.218222

0.246559

0.269420

0.275874

0.289105

Grid search results

0.000354

0.000020

0.000061

0.000212

0.009950

0.009837

0.009659

0.009642

0.010023

Observations

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

Tomáš Báča

Environment

Pandas

Learn

Regression pipeline

Metrics

example example

Preprocessir

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



1 {'predictor_n_neighbors': 1}

2 {'predictor_n_neighbors': 2}

3 {'predictor_n_neighbors': 3}

4 {'predictor_n_neighbors': 4}

5 ('predictor n neighbors': 5)

6 {'predictor_n_neighbors': 6}

0.324068

0.468788

0.518547

0.540323

0.551149

0.558435

params split0 test score split1 test score split2 test score mean test score std test score rank test score

0.323371

0.473595

0.400927

0.511781

0.521134

0.327423

0.465544

0.511827

0.535041

0.547414

0.005245

0.032361

0.028867

0.026857

0.027696

0.026652

0.334830

0.543340

0.564974

0.579313

Data preprocessing — more about sklearn's transformers



Tomáš Báča

Regression

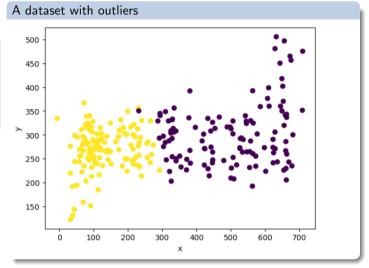
Preprocessing

2025-10-11

two classes

Observations

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



Tomáš Báča (CTU in Prague)

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

Scikit Learn Preprocessing

Data preprocessing — more about sklearn's transformers



Data preprocessing — more about sklearn's transformers



Tomáš Báča

ment

Scikit

Toy datase

Preprocessing

Realwor

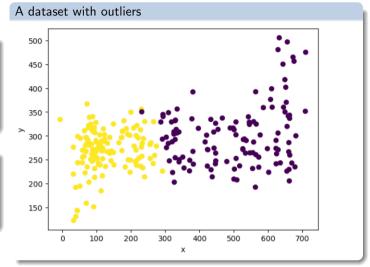
example

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?



Tomáš Báča (CTU in Prague)

A CLASS OF THE PROPERTY OF THE

September 30th, 2025

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

Scikit Learn
Preprocessing

Data preprocessing — more about sklearn's transformers



Lecture 2: Classical methods models in MLE with Scikit Learn

> Tomáš Báča

Regression Preprocessing

2025-10-11

Observations

Standard scaler

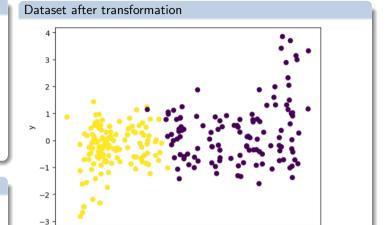
where

- outliers still visible
- axes' range is near zero

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

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

• the shape of the data is similar



Tomáš Báča (CTU in Prague)

-1.0

-1.5

-0.5

0.0

0.5

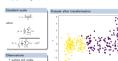
1.5

1.0

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

Scikit Learn Preprocessing

Data preprocessing — Standard Scaler



2.0

Data preprocessing — MinMax scaler



Tomáš Báča

Python Environment

Pandas

Scikit

Regression pipeline

Preprocessing

example

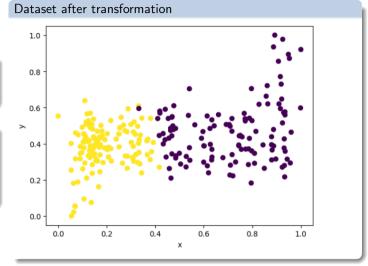
2025-10-11

MinMax scaler

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

Observations

- outliers still visible
- axes' range is between 0 and 1
- distribution is preserved



Tomáš Báča (CTU in Prague)

-- / -

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

-Scikit Learn

Preprocessing

Data preprocessing — MinMax scaler



Data preprocessing — Quantile transformer

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

Tomáš Báča

Environment

Pandas

Learn

Regression pipeline

Preprocessing

example

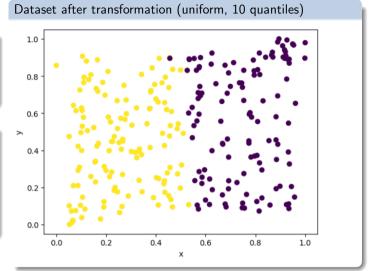
2025-10-11

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



Tomáš Báča (CTU in Prague)

A COLUMN TO A STATE OF CASE

September 30th, 2025

26./4

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

Preprocessing

Data preprocessing — Quantile transformer



Data preprocessing — Quantile transformer

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

Tomáš Báča

Environment

Panda

Learn

Regression pipeline

Preprocessing

Deelissed

example

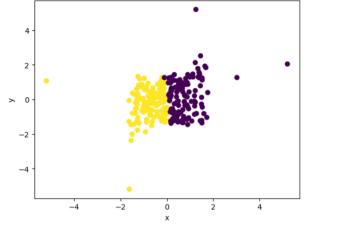
Quantile transformer

- 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





Tomáš Báča (CTU in Prague)

A COLUMN TO THE MET WE SHOULD

eptember 30th, 2025

Scikit Learn

Preprocessing

Data preprocessing — Quantile transformer

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



Linear separability

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Regression Preprocessing

Problem?

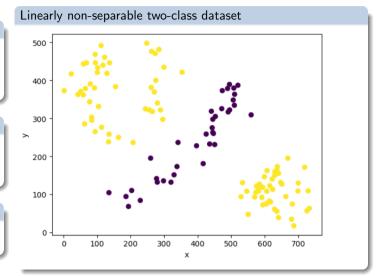
- classes are not linearly separable
- a.k.a., a single separating hyperplane can not be found

Solution?

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

Easier solution?

- use more complex classifier
- lift feature to new dimension



Tomáš Báča (CTU in Prague)

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

Preprocessing

Linear separability



Linear separability

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

Tomáš Báča

Regression Preprocessing

Problem?

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

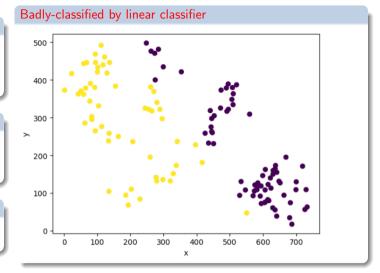
Solution?

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

Easier solution?

- use more complex classifier
- lift feature to new dimension

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



Tomáš Báča (CTU in Prague)

Scikit Learn

Preprocessing

Linear separability



Polynomial features

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

Tomáš Báča

Environment

Call it

Toy data

Regression pipeline Preprocessing

Metrics

Realworl example

Lifting features using polynomial interactions

example for order = 2

$$\mathbf{x} = \begin{bmatrix} x_1, x_2 \end{bmatrix}^{\mathsf{T}}$$

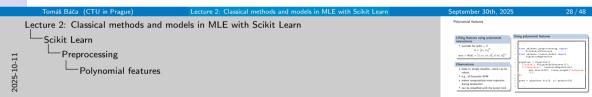
$$\mathbf{x}_{\mathsf{new}} = \Phi(\mathbf{x}) = \left[1, x_1, x_2, x_1^2, x_1 x_2, x_2^2\right]^\mathsf{T}$$

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

```
from sklearn.preprocessing import
        PolynomialFeatures
   from sklearn.linear model import
        LogisticRegression
5
6
   pipeline = Pipeline([
      ("scale", PolynomialFeatures()),
8
      ("classifier", LogisticRegression(
          max iter=1000, class weight="balanced
          ")),
   1)
q
   pred = pipeline.fit(X, y).predict(X)
11
```

Using polynomial features



• The Kernel trick with the SVM (or other classifier that uses a dot product between the attributes and its weights) is about doing just the dot product using the lifted dimensions but training the classified with the original dimension.

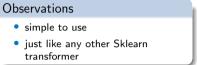
Polynomial features

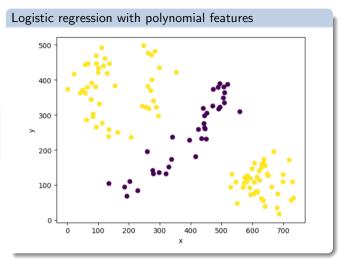
Lecture 2: Classical methods and models in MLE with Scikit Learn

> Tomáš Báča

Regression

Preprocessing









Polynomial features



Tomáš Báča

ment

Fallua

Toy dat

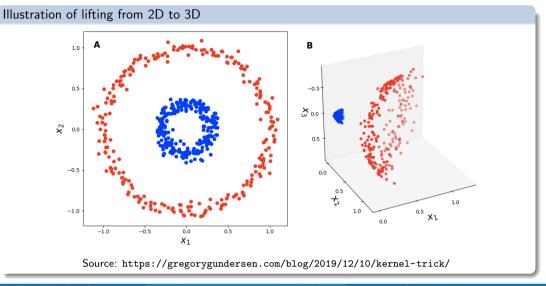
Regression pipeline

Preprocessing

IVICTICS

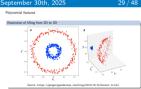
example

2025-10-11



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

eprocessing Polynomial features



Encoding text and categorical attributes

Lecture 2: Classical methode and models in MIE with Scikit

Tomáš Báča

Preprocessing

How to encode text?

- added uunique ID in 1D space might cause problems
- better wav?
 - 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

```
1
   from sklearn.preprocessing import OneHotEncoder
2
3
   X = [['red'], ['green'], ['blue'], ['red']]
4
   encoder = OneHotEncoder(handle_unknown="ignore")
5
   encoder.fit(X)
6
7
   new_X = encoder.transform(X).todense()
```

```
matrix([[0., 0., 1.],
            [0., 1., 0.],
2
            [1., 0., 0.],
            [0., 0., 1.]])
```

Inverse:

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

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

Tomáš Báča (CTU in Prague)

Lecture 2: Classical methods and models in MLE with

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

Encoding text and categorical attributes

Objects in Scikit

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

> Tomáš Báča

Python Environment

Pandas

Learn

Regression pipeline

Preprocessing

Metrics

Estimator

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

Transformer

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

Pipeline

• is a Predictor

GridSearchCV

is a Predictor

Predictor

• implements .fit(), .predict()

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

regressors, classifiers

Tomáš Báča (CTU in Prague)

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

Objects in Scibit

Estimator

* anjuneses.1510

* only consensed size

* a.g., wave estimities collector

| Pigaline

—Scikit Learn

Preprocessing

Objects in Scikit

2025-10-11

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

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

Recall

Precision

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

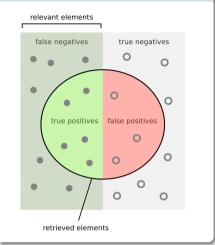
did I marked all the relevant instances?

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

Diagram of possible results



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

Binary classification metrics — Precision, Accuracy and Recall



Metrics example

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

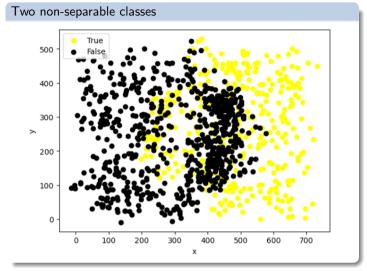
Regression

Metrics

2025-10-11

Is there a single best classifier?

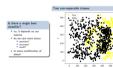
- no, it depends on our metrics
- do we care more about
 - precision?
 - accuracy?
 - recall?
- or some combination of these?



Tomáš Báča (CTU in Prague)

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

Metrics example



Logistic regression with default metrics

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

Tomáš Báča

Environment

randas

Toy data

Regression

Metrics

example

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

```
2
3
   from sklearn.linear_model import LogisticRegression
   from sklearn.model_selection import GridSearchCV
4
5
6
   pipeline = Pipeline([
      ("classifier", LogisticRegression(max_iter=1000,
          verbose=0)),
8
   1)
9
   model = GridSearchCV(
      estimator=pipeline, cv=4, param_grid={
      'classifier__class_weight': [{0: 5, 1: v} for v in
           range(1, 10)]},
14
   pred = model.fit(X, y).predict(X)
16
```

Tomáš Báča (CTU in Prague)

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

September 30th, 2025

34 / 49

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

Logistic regression with default metrics

Metrics example — Accuracy

• the decision is cut at a

Metrics on the dataset

• accuracy = 0.77

precision = 0.97

recall = 0.35

reasonable compromise many FP as well as FN

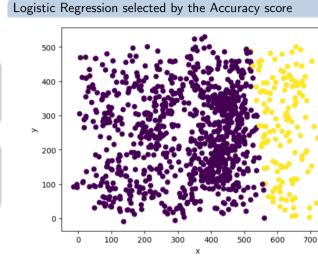
Lecture 2: Classical methods and models in MLE with Scikit Learn

Tomáš Báča

Regression

Metrics

2025-10-11



Tomáš Báča (CTU in Prague)

Observations

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

Metrics example — Accuracy



Logistic regression with Precision metrics

Lecture 2: Classical methods and models in MLE with Scikit

Tomáš Báča

Python Environment

Pandas

Scikit

Regression pipeline

Metrics

Realworl example

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

```
1
   from sklearn.linear_model import LogisticRegression
2
   from sklearn.model_selection import GridSearchCV
   from sklearn.metrics import precision score.
        accuracy_score, recall_score
4
5
   pipeline = Pipeline([
     ("classifier", LogisticRegression(max_iter=1000,
6
          verbose=0)).
   1)
8
q
   model = GridSearchCV(
     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
   pred = model.fit(X, y).predict(X)
16
```

Tomáš Báča (CTU in Prague)

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

September 30th, 2025

36 / 48

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

Logistic regression with Precision metrics



025-10-11

Metrics example — Precision

Lecture 2: Classical methods models in MLE with Scikit Learn

> Tomáš Báča

Regression

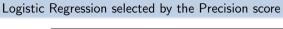
Metrics

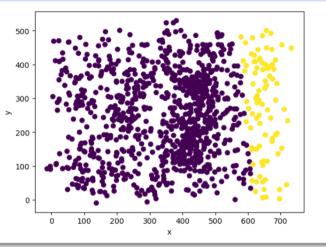
Observations

- 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





Tomáš Báča (CTU in Prague)

2025-10-11

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

Metrics example — Precision



Logistic regression with Recall metrics

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

Tomáš Báča

Environment

Pandas

Scikit

Regression pipeline

Metrics

Realwork

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

```
1
   from sklearn.linear_model import LogisticRegression
2
   from sklearn.model_selection import GridSearchCV
   from sklearn.metrics import precision score.
        accuracy_score, recall_score
4
5
   pipeline = Pipeline([
     ("classifier", LogisticRegression(max_iter=1000,
6
          verbose=0)).
   1)
8
q
   model = GridSearchCV(
     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
   pred = model.fit(X, y).predict(X)
16
```

Tomáš Báča (CTU in Prague)

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

September 30th, 2025

38 / 4

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

Logistic regression with Recall metrics



Metrics example — Recall

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

> Tomáš Báča

Regression

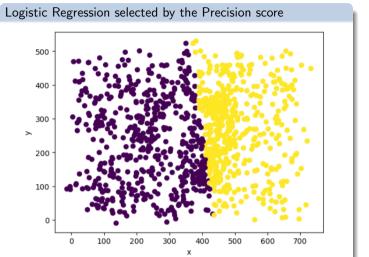
Metrics

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



Tomáš Báča (CTU in Prague)

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

Metrics

Metrics example — Recall



Storing a trained sklearn model

Lecture 2: Classical methods models in MLE with Scikit Learn

Tomáš Báča

Regression

Preprocessing Metrics

```
Storing a model
```

```
2
   from joblib import dump
   dump(model, "my_model.joblib")
6
```

Loading a model

```
1
   from joblib import load
3
   model = load("my_model.joblib")
6
```

Tomáš Báča (CTU in Prague)

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

• as simple as dumping the object into a file

alternatively, you can use pickle

Metrics

Storing a trained sklearn model

Realworld example — Realtime classification of ionizing particles

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

> Tomáš Báča

Environment

Pandas

Learn

Regression pipeline

Preprocess

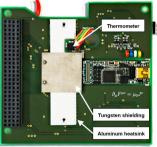
. .

example

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)







Tomáš Báča (CTU in Prague)

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

September 30th, 2025

41 / 48

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

campic

Realworld example — Realtime classification of ionizing particles





Timepix - Ionizing radiation dosimetry and imaging

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

Tomáš Báča

Environment

Pandas

Learn

Regression pipeline

Matrice

Realworl example

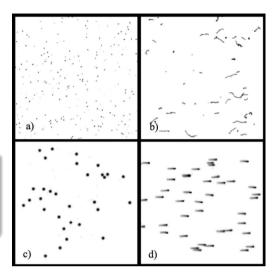
- incoming particle and the sensor
 machine learning was applied to dist
- machine learning was applied to distinguish particle types [3]

• real-time recording of the interaction of an

 no dark-current noise (the images are spotless besides the actual data)

Examples in the Figure

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



Tomás Báča (CTU in Prague)

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

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

Realworld example

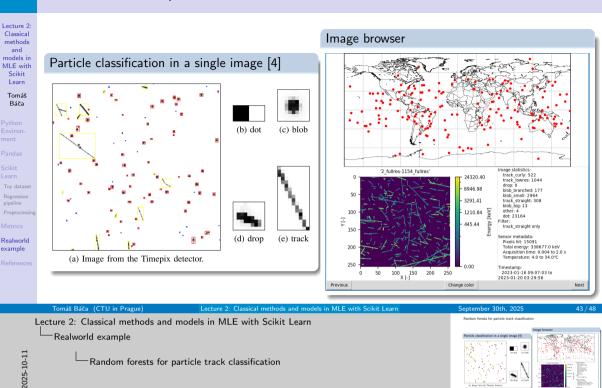
Timepix - Ionizing radiation dosimetry and imaging

Timepix - Ionizing radiation dosimetry and imaging

Timepix - Ionizing radiation dosimetry and imaging

• The sensor is a digital counterpart of the Cloud chamber (https://en.wikipedia.org/wiki/Cloud_chamber)

Random forests for particle track classification



 We have built a classified dataset of the data from the 6-years of deployment in the Low-Earth Orbit: https://github.com/vzlusat/vzlusat1-timepix-data

Low-Earth orbit particle classification

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

> Tomáš Báča

Python Environ ment

Panda:

Learn Toy dat

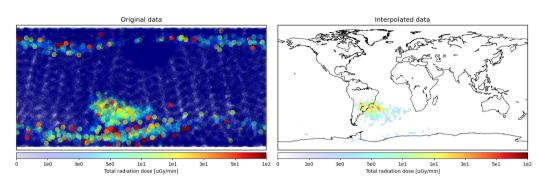
Regression pipeline Preprocess

IVICTIC

example

Plot for Beta particles

VZLUSAT-1 Timepix radiation dose, 2017-07-13 10:17:57 to 2023-05-15 09:47:37



https://github.com/vzlusat/vzlusat1-timepix-data [4], [5]

Tomáš Báča (CTU in Prague)

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

ptember 30th, 2025

.

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

eaiworid exampi

Low-Earth orbit particle classification

Pine for this agent of the control o

Realtime classification onboard UAVs

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

Tomáš

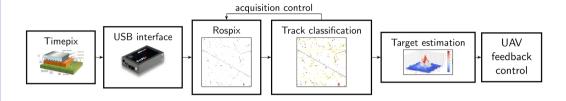
example

2025-10-11

Báča

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

Tomáš Báča (CTU in Prague)

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

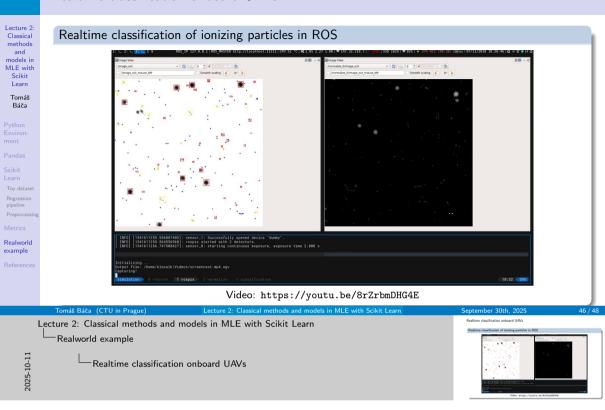
Proof of concept for UAVs

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

Realtime classification onboard UAVs

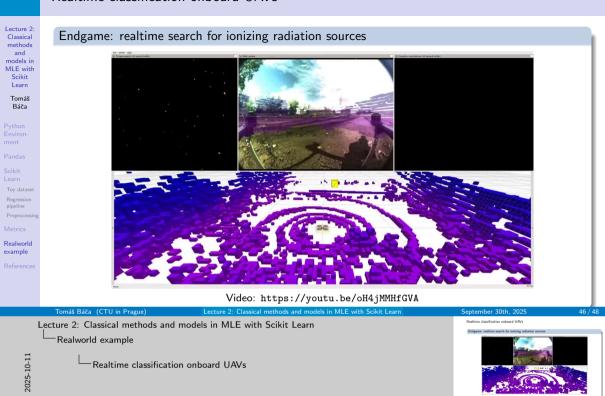
- The particle classifier is a necessary tool for real-time control of the acquisition parameters of the detector (acquisition time, bias, etc.).
- The particle spectrum measurement is also needed for educated decision about what source of radiation we are observing, this leads to better localization with some advanced techniques, such as the Compton camera localization baca2021gamm.

Realtime classification onboard UAVs

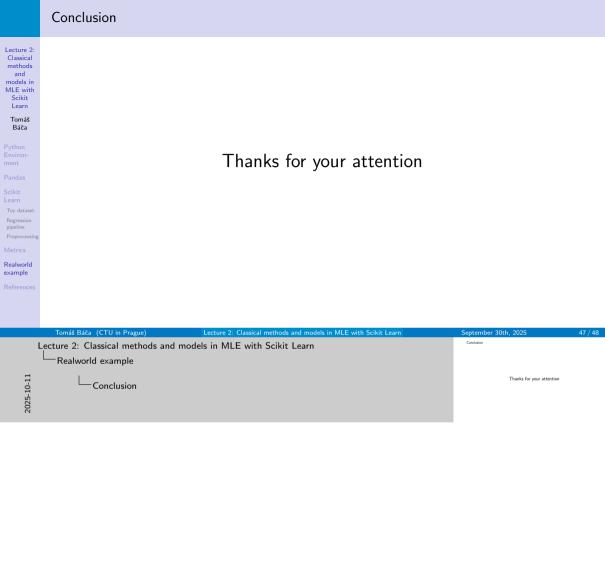


• The classification needs to run for hundreds of particles per second on a limited hardware of a drone.

Realtime classification onboard UAVs



• The classification needs to run for hundreds of particles per second on a limited hardware of a drone.



References I

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

Tomáš Báča

Regression

References

- 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.
- 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,
- 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.
- 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.

Tomáš Báča (CTU in Prague)

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

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

References

- p. Collet, 2016.
 [8] T. Borr, W. Shi, I. Venne, et al., "Empirica Mili Grids enhand the NEXIGE", it Manuschilde is part of Span Reducine Chainney Means. Street, and the Contract of the Contract of