

JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May

'22

18

Wednesday

Week 21 Day (138-227)

09.00 — DECISION TREES —

• TREE BASED MODELS CAN BE USED FOR REGRESSION AND

10.00 CLASSIFICATION.

• THEY INVOLVED IN DIVIDING THE PREDICTIONS SPACE INTO A

11.00 NUMBER OF REGIONS. THE SET OF SPLITTING RULES CAN BE
SUMMARIZED IN A TREE.

12.00 → HENCE IT IS CALLED DECISION TREE.

* A SINGLE DECISION TREE IS OFTEN NOT BETTER THAN LINEAR REGRESSION

13.00 LOGISTIC REGRESSION OR LDA.

* BAGGING, RANDOM FOREST AND BOOSTING DRAMATICALLY IMPROVE

14.00 PERFORMANCE.

15.00 BASIC TERMINOLOGIES

→ Trees are drawn upside down.

16.00 → The final regions are called leaves. The point where split occurs is called a node.

17.00 → The segment that nodes are (The segment that connects the nodes are called BRANCHES).

18.00

HOW REGRESSION TREES WORKS →

on next page →



19

22

May

Week 21 Day (139-226)

Thursday

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

09.00 TO CREATE A REGRESSION TREE

- Divide the predictor space into J distinct and non-overlapping regions.
- For every observation that falls in a region, predict the mean of the response value in that region.
- Each REGION IS SPLIT TO MINIMISE THE RSS (Residual sum of squares)
- It uses a top-down greedy approach called recursive binary splitting.

14.00 TOP DOWN GREEDY APPROACH -

- TOP DOWN - ALL OBSERVATIONS ARE IN A SINGLE REGION BEFORE THE FIRST SPLIT.
- GREEDY - THE BEST SPLIT OCCURS AT A PARTICULAR STEP TO MAKE THE BEST PREDICTION AT THAT STEP, RATHER THAN LOOKING FORWARD AND MAKING A SPLIT THAT WILL GIVE A BETTER RESULT IN A FUTURE STEP.

18.00 Regression Tree

$$R_1(j,s) = \{x | X_j < s\} \text{ and } R_2(j,s) = \{x | X_j \geq s\}$$

we seek j and s to minimise :

$$\sum_{i: x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2$$

JUNE 2022

SUN	MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31											

May 20
Friday

Regression Tree in Python

```
from sklearn.tree import  
DecisionTreeRegressor  
clf = DecisionTreeRegressor()  
clf.fit(X-train, y-train)
```

CLASSIFICATION TREE

- Predict the most commonly occurring class in a region.
- RSS cannot be used, each split must minimize the classification error rate.

CLASSIFICATION ERROR RATE.

$$E = 1 - \max(\hat{P}_{mK})$$

- Not sensitive enough for tree growing
- So instead we use Gini Index

21

'22

May

Week 21 Day (141-224)

Saturday

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31											

09.00

GINI INDEX - MEASURE OF TOTAL VARIANCE ACROSS ALL CLASSES.

10.00

$$G_i = \sum_{k=1}^K \hat{P}_{mk} (\hat{P}_{mk})^2$$

11.00

12.00

13.00

- Small if proportion is close to 0 or 1.
- Good measure of node purity.

CROSS-ENTROPY - IT CAN ALSO BE USED FOR TREE GROWING.

14.00

$$D = - \sum_{k=1}^K \hat{P}_{mk} \log (\hat{P}_{mk})$$

15

16

17

Sunday 22

18.00

Classification tree in Python

from sklearn.tree import DecisionTreeClassifier

clf = DecisionTreeClassifier()

clf.fit(X_train, Y_train)



JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26

27 28 29 30

May

'22

23

Monday

Week 22 Day (143-222)

09.00 ADVANCE TOPICS ON DECISION TREE.

- BAGGING - BOOTSTRAP AGGREGATION
- BOOTSTRAP CAN COMPUTE STANDARD DEVIATION OF ANY QUANTITY.
FOR DECISION TREES, VARIANCE TENDS TO BE HIGH.
- BAGGING CAN REDUCE THE VARIANCE AND IMPROVE THE PERFORMANCE OF A DECISION TREE.
- REPEATEDLY DRAW SAMPLES FROM THE DATASET, GENERATING B DIFFERENT BOOTSTRAP TRAINING SETS.
- TRAIN ON ALL BOOTSTRAPPED TRAINING SETS TO GET A PREDICTION FOR EACH SET AND AVERAGE THE PREDICTION.

14.00

BAGGING -

15.00

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x)$$

16.00

17.00

18.00

- We can construct high number of trees with high variance and low bias.
- Then, average the prediction to reduce variance and improve performance.

BAGGING IN PYTHON

```
from sklearn.ensemble import BaggingClassifier  
bagging_clf = BaggingClassifier()  
bagging_clf.fit(X_train, y_train.ravel())
```



24

22

May

Tuesday

Week 22 Day 1 (44-22)

MAY 2022

SUN	MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
							1	2	3	4	5	6	7	8
9	10	11	12	13	14	15	16	17	18	19	20	21	22	
23	24	25	26	27	28	29	30	31						

09.00

RANDOM FOREST -

- RANDOM FOREST IMPROVES UPON BAGGING, BY MAKING A SMALL TWEAK THAT DECORATES THE TREES.
- MULTIPLE DECISION TREES ARE BUILT. AT EACH SPLIT, ONLY A RANDOM SAMPLE OF ' m ' PREDICTORS IS CHOSEN FROM ALL ' p ' PREDICTORS.
- THE SPLIT IS ONLY ALLOWED TO USE ONE OF THE ' m ' PREDICTORS.

13.00

Randomforest -

14.00

$$m = \sqrt{p}$$

15.00

16.00

- IN BAGGING, IF THERE IS A STRONG PREDICTOR, IT WILL LIKELY BE THE TOP SPLIT, AND ALL TREES WILL BE SIMILAR.
- VARIANCE WILL NOT BE REDUCED

18.00

- RANDOM FOREST - ONLY A SUBSET OF PREDICTORS IS FORCED FOR EACH SPLIT.
- If $m = p$, then it is like bagging



JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May

22

Wednesday

25

Week 22 Day (145-220)

09.00 RANDOM FOREST IN PYTHON

10.00 from sklearn.ensemble import

11.00 RandomForestClassifier

12.00 random_clf = RandomForestClassifier(100)

13.00 random_clf.fit(X_train, y_train.ravel())

14.00

15.00

16.00

BOOSTING

17.00 - SIMILAR TO BAGGING, BUT TREES ARE GROWN SEQUENTIALLY

EACH TREE USES INFORMATION FROM THE PREVIOUSLY

18.00 GROWN TREES.

- LEARNS SLOWLY

- EACH TREES FIT THE RESIDUAL INSTEAD OF TARGET VARIABLE.

EACH TREE IS SMALL AND WILL SLOWLY IMPROVE THE PREDICTIONS:

- 3 TUNING PARAMETER FOR BOOSTING -

→ NO. OF TREES (B) - BOOSTING CAN OVERFIT IF B IS 100 LARGE. USE CROSS VALIDATION.

26

'22

May

Week 22 Day (146-219)

Thursday

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8						
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30						

09.00

- SHRINKAGE PARAMETER (d) - small positive number that controls the learning rate.

10.00

Typically 0.01 or 0.001

11.00

- NUMBER OF SPLITS IN EACH TREE (d) - controls the complexity of the boosted ensemble. A Single split works well ($d=1$) Also called interaction depth.

12.00

BOOSTING IN PYTHON -

```

1 from sklearn.ensemble import
1 GradientBoostingClassifier
1 boost_clf = GradientBoostingClassifier(
1     learning_rate = 0.01, n_estimators = 100,
1     max_depth = 2)
1 boost_clf.fit(X_train, y_train.ravel())

```



+91 9876543210



abc@xyz.com



www.abc.com

JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May

22

Friday

27

Week 22 Day (147/218)

09.00 - DECISION TREES PYTHON HANDS ON -

10.00
import numpy as np
11.00 import pandas as pd
11.00 import matplotlib.pyplot as plt
12.00 import seaborn as sns
from sklearn.metrics import plot_confusion-
matrix
13.00 %matplotlib inline

14.00

15.00

16.00

numpy - Library for scientific computing in python.

pandas - library that provides data structure for efficient data manipulation and analysis.

matplotlib - library provides variety of visualization

seaborn - Library provides high level interface for creating attractive graphics.

sklearn.metrics.plot_confusion_matrix - function allow to plot confusion matrix.

%matplotlib inline - The magic command.



28

22

May

Saturday

Week 22 Day (148-2) 7

09.00

10.00

11.00

12.00

13.00

14.0

15.1

16.

17.

18.

IMPORTING THE DATASET.

```
datapath = 'url'
data = pd.read_csv(datapath)
data.head()
```

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

```
def violin_plots(x, y, data):
    for i, col in enumerate(y):
        plt.figure(i)
        sns.set(rc={"figure.figsize": (11.7, 8.27)})
        ax = sns.violinplot(x=x, y=col, data=data)
    y = data.columns[:-1]
    x = data.columns[-1]
    violin_plots(x, y, data)
```

Sunday 29

X is a string that represent column name of independent variable in the dataset.

Y is list of strings

data is pandas datframe.



JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May 30 Monday

Week 23 Day 056-215

09.00 function iterates over columns in y and create a separate figure for each column. It sets figure size using sns.set (rc={"figure.figsize": (11.7, 8.27)})) and creates a Seaborn violin plot using sns.violinplot (x=x, y=col, data=data), where col is the current dependant variable in the iteration. Finally returns the violin plots

10.00

11.00
12.00
13.00
14.00
15.00

```
for col in data.columns:  
    print(f" {col}: {data[col].isnull().sum()}" )
```

16.00 code iterates over columns in data dataframe using for a loop for each column, it prints the column name and number of missing values in the column using the print function and f-string. The number of missing values is computed using isnull() method, which returns a boolean mask indicating whether each value in the column is missing or not., and sum () method which sums up the number of missing values in the column.



31

22

May

Tuesday

Week 23 Day (151-214)

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31											

```
from sklearn.preprocessing import
```

```
LabelEncoder
```

```
le = LabelEncoder()
```

```
data['Classification'] = le.fit_transform
```

```
(data['classification'])
```

```
data.head()
```

→ 1

→ 2

→ 3

LabelEncoder class from scikit learn Library to encode a

16.00 categorical variable in pandas dataframe

1) converts categorical variables into numeric labels

17.00 2) le is instance of LabelEncoder class

3) Classification is a string that represent column

18.00 name of categorical variable

→ Method first fits the encoder to the unique values

in the column and then transform the column by

replacing each unique value with corresponding

integer label.

→ The resulting encoded values are then stored back

in Classification column of data using
statement data (3) statement

JULY 2022

SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8
9	10	11	12	13	14	15	16
17	18	19	20	21	22	23	24

25 26 27 28 29 30 31

June

22

02

Thursday

Week 23 Day (153-212)

09.00

```
from sklearn.model_selection  
import train_test_split
```

10.00

```
y = data['Classification'].values.reshape  
(-1,1)
```

11.00

```
X = data.drop(['Classification'], axis=1)
```

12.00

```
X_train, X_test, y_train, y_test =  
train_test_split(X,y, test_size = 0.1,  
random_state=42)
```

13.00

14.00

15.00

1) Split data into training and testing

16.00 2) Classification is string that represent the column name of the target variable in data datframe.

17.00 y contains values of target variable , reshaped into a 2D array using reshape method.

18.00 X is pandas datframe that contains feature variable of the dataset , with target variable column dropped using the drop method and axis parameter set to 1. (column-wise)

test_size is float that represent proportion of the dataset to include in testing-set.

random_state is integer that sets the seed for the random number generator used to split the dataset.

Code is useful for preparing dataset for Machine Learning

03

22

June

Week 23 Day (156-211)

Friday

JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

09.00

- Baseline Model -

10.00

from sklearn.tree import DecisionTreeClassifier

11.00

```
clf = DecisionTreeClassifier()
```

12.00

```
clf.fit(X-train, y-train)
```

13.00

```
plot_confusion_matrix(clf, X-test,  
y-test, cmap=plt.cm.Blues)
```

14.00

```
plt.grid(False)
```

```
plt.show()
```

15.00

16.00

DecisionTreeClassifier is class from sklearn.tree module that implements a decision tree classifier.

clf is an instance of Decision Tree Classifier class.

X-train and y-train are numpy arrays that contain the feature and target variable of training set.

X-test and y-test — testing set.

cmap — colormap is used to color cells of confusion matrix.
plt.grid(False) turns off the grid lines in the plot..



JULY 2020

SUN	MON	TUE	WED	THU	FRI	SAT
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

June 22

Saturday

04

Week 21 Sat (106-216)

```
from sklearn.tree import plot_tree  
plot_tree(clf, max_depth=5)  
  
// plot-tree function is scikit-learn  
library to visualize decision tree-  
classifiers.  
// clf is an instance that has been trained  
on dataset.  
// max-depth integer that specifies the  
maximum depth of tree to be plotted
```

Code uses the plot-tree function to plot decision trees classifiers using maximum depth of 5. Result shows each node representing decision based on feature variables and each leaf node representing predicted target variable. Sunday 05 The plot also displays impurity score and number of samples at each node.

Code is useful for visualizing the decision making process of a decision tree classifier, and can help to identify patterns and insights in the data that classifier is using to make predictions.

06

Jun

Monday

WEEK 2022

Mon	Tue	Wed	Thu	Fri	Sat	Sun
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28

10:00

BAGGING -

from sklearn.ensemble import BaggingClassifier

bagging_clf = BaggingClassifier()

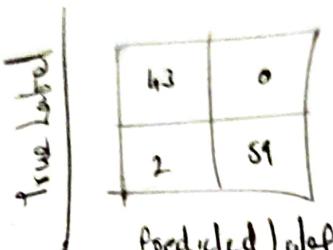
bagging_clf.fit(X_train, y_train, ravel())

plot_confusion_matrix(bagging_clf,

X_train, y_train, cmap=plt.cm.Blues)

plt.grid(False)

plt.show()

Explanation in ipynb fileRandom Forest Classifier

from sklearn.ensemble import RandomForestClassifier

random_clf = RandomForestClassifier(100)

random_clf.fit(X_train, y_train, ravel())

plot_confusion_matrix(random_clf,

X_test, y_test, cmap=plt.cm.Blues)

plt.grid(False)

plt.show()

JULY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
					1	2	3	4	5	6	7	8	9
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30	31							

June

22

07

Tuesday

Week 24 Day (158-207)

09.00

Boosting -

10.00

from sklearn.ensemble import
GradientBoostingClassifier

11.00

boost_clf = GradientBoostingClassifier()

12.00

boost_clf.fit(X_train, Y_train.ravel())

13.00

plot_confusion_matrix(boost_clf, X_test,

y_test, cmap=plt.cm.BuGn)

plt.grid(False)

plt.show()

14.00

15.00

16.00

Explanation in Source file -

17.00

18.00

