

L1: Looking at the big picture

Steps in ML projects

- ✓ 1. Look at the big picture.
- ✓ 2. Get the data. → data viz.
- ✓ 3. Discover and visualize the data to gain insights.
- ✓ 4. Prepare the data for Machine Learning algorithms.
- ✓ 5. Select a model and train it. → EDA → Scaling, LR, Tree pipelines
- ✓ 6. Fine-tune your model. → hyper parameters tuning
- ✓ 7. Present your solution.
- ✓ 8. Launch, monitor and maintain your system.

- ML is usually a small piece in a big project. e.g. wine quality prediction is a small piece in setting up the manufacturing process.
- Typically (10-15%) of time is spent on ML.
- A lot more time is spent on capturing and processing data needed for ML and taking decisions based on output of ML module.
- Needs strong collaboration with domain experts, product managers and eng-teams for successful execution.

- Is this a supervised, unsupervised or a RL problem?
- Is this a classification, regression or some other task?
- What is the nature of the output: single or multiple outputs?
- Does system need continuous learning or periodic updates?
- What would be the learning style: batch or online? import pandas as pd

data.head() → 5 rows of data

↳ how does it look like

↳ training columns → features attribute

1 X_train = train_df.drop('Price', axis=1)
2 y_train = train_df['Price']

↳ final output

↳ info. of columns, how many total entries, empty entries

data.describe()

↳ statistical attributes abt. numerical data

data['quality'].value_counts() → 1 → 10

↳ col name 2 → 20 ...

from sklearn.model_selection import train_test_split X_train, Y_train

X_t, X_e, Y_t, Y_e = train_test_split(X_train, Y_train, test_size=0.2)

1) Data snooping

↳ processes

Y_train, test_size=0.2)

↳ Leads to biased estimation on test sets

↳ Increases the risk of false positives

FP

Leads to better estimation on training sets

Reduces the risk of false positives

2) Which of the following statements are correct?

Random sampling shuffles all of the data points and take it whole as a train set.

Stratified sampling samples test examples such that they are representative of overall distribution

Stratified sampling samples test examples such that they are representative of the training dataset

3) Mention TRUE or FALSE: Stratified sampling gives us test distribution closer to the overall distribution as compared to the random sampling

TRUE

FALSE

sklearn

About Random and Stratified Sampling

↓
k no. of random samples choose

↳ represent overall distribution of the dataset

valida" dataset

L2: Data Visualization

- Performed on training set.
- In case of large training set -
 - Sample examples to form exploration set.
- Enables to understand features and their relationship among themselves and with output label.

In our case, we have a small training data and we use it all for data exploration. There is no need to create a separate exploration set.

data_c = data.copy()

It's a good idea to create a copy of the training set so that we can freely manipulate it without worrying about any manipulation in the original set.

eg. sns.scatterplot(x='fixed acidity', y='density', data=exploration_set) → sns seaborn

eg. exploration_set.plot(kind='scatter', x='fixed acidity', y='density', alpha=0.5, c='quality', cmap=plt.get_cmap('jet'))

• Standard correlation coefficient between features.

- Ranges between -1 to +1
- Correlation = +1: Strong positive correlation between features
- Correlation = -1: Strong negative correlation between features
- Correlation = 0: No linear correlation between features

▪ Visualization with heat map

▪ Only captures linear relationship between features.

▪ For non-linear relationship, use rank correlation

blue, red

corr_matrix['quality']

fig.

plt.figure(figsize=(14,7)) sns.heatmap(corr_matrix, annot=True)

1) A standard correlation coefficient value of +0.98 represents-

Strong positive correlation between features

Weak positive correlation between features

No correlation between features

None of these

3) matplotlib and seaborn are used for which of the following purposes?

data loading

data pre-processing

data visualization

data manipulation

2) Why is it a good idea to create a copy of the training set?

So that we can freely manipulate it without altering original training data set

So that changes in copy of the dataset are reflected in original dataset

so that we can mix it with test dataset

4) No linear correlation between the features is demonstrated by which of the following correlation values?

0

1

-1

2

5) Mention TRUE or FALSE: A heatmap can be used to visualize correlation matrix

TRUE

FALSE

→ EDA (exploratory data analysis)

L3: Data Preparation

We often need to preprocess the data before using it for model building due to variety of reasons:

- Due to errors in data capture, data may contain outliers or missing values.
- Different features may be at different scales.
- The current data distribution is not exactly amenable to learning.

Typical steps in data preprocessing are as follows:

1. Separate features and labels.
2. Handling missing values and outliers.
3. Feature scaling to bring all features on the same scale.
4. Applying certain transformations like log, square root on the features. (optional)

It's a good practice to make a copy of the data and apply preprocessing on that copy. This ensures that in case something goes wrong, we will at least have original copy of the data intact.

wine dataset
wine_features.isna().sum()

sum of all the entries of every column that has missing values.

{ dropna()
drop() } pandas (pd)

new missing values drop.

Imputer → Simple, KNN, Kmeans

```
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy="median")
```

(obj)
imputer.fit_transform <value>
Do fit_transform (wf) → numpy array
arrayify

wine_features_tr = pd.DataFrame(tr_features, columns=wine_features.columns)

String → object M,F → 1
from sklearn.preprocessing import OrdinalEncoder / One Hot Encoder
ordinal_encoder = OrdinalEncoder()

scikit-learn → fit_transform(dataset)

pandas dataframe
feature scaling sklearn. impute

4.5.1 Min-max scaling or Normalization

- We subtract minimum value of a feature from the current value and divide it by the difference between the minimum and the maximum value of that feature.
- Values are shifted and scaled so that they range between 0 and 1.
- Scikit-Learn provides MinMaxScaler transformer for this.
- One can specify hyperparameter feature_range to specify the range of the feature.

↓ better

4.5.2 Standardization

- We subtract mean value of each feature from the current value and divide it by the standard deviation so that the resulting feature has a unit variance.
- While normalization bounds values between 0 and 1, standardization does not bound values to a specific range.
- Standardization is less affected by the outliers compared to the normalization.
- Scikit-Learn provides StandardScaler transformation for feature standardization.
- Note that all these transformers are learnt on the training data and then applied on the training and test data to transform them.
- Never learn these transformers on the full dataset.

range
0 → 1

no fixed range

Pipeline, Feature Union, Column Transformer

```
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
transform_pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy="median")),
    ('std_scaler', StandardScaler())
])
wine_features_tr = transform_pipeline.fit_transform(wine_features)
```

type("") ↓ pd

- The real world data has both categorical as well as numerical features and we need to apply different transformations to them.
- Scikit-Learn introduced ColumnTransformer for this purpose.

```
1 from sklearn.compose import ColumnTransformer
```

- 1) The function .isna() is used for which of the following purposes?
- to detect missing (NaN) values
- to generate heatmap
- to plot standard deviation of available data in the place of NaN values

- 2) Which of the following methods from pandas may be used to drop rows containing missing values?

- dropna()
- drop()
- na_drop()
- nAn()

- 3) Mention TRUE or FALSE : Sklearn provides SimpleImputer class for filling up missing values.

- FALSE
- TRUE

- 4) fit_transform method on OrdinalEncoder object is called for which purpose?

- convert numbers to texts
- convert text to numbers
- convert decimal to binary
- convert real numbers to rational numbers

- 5) Mention TRUE or FALSE: Scikit-Learn provides a Pipeline class to line up transformations in an intended order.

- TRUE
- FALSE

L4: Selection and Training of ML Models

```
from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(wine_features_tr, wine_labels)
```

from sklearn.metrics import mean_squared_error

quality_predictions = lin_reg.predict(wine_features_tr)
mean_squared_error(wine_labels, quality_predictions)

neg - ()

Score = ?

Ex. from sklearn.tree import DecisionTreeRegressor
tree_reg = DecisionTreeRegressor()
tree_reg.fit(wine_features_tr, wine_labels)

Note that the training error is 0, while the test error is 0.58. This is an example of an overfitted model.

→ training ↓, test ↑

We can use cross-validation (CV) for robust evaluation of model performance.

```
1 from sklearn.model_selection import cross_val_score
```

- Cross validation provides a separate MSE for each validation set, which we can use to get a mean estimation of MSE as well as the standard deviation, which helps us to determine how precise is the estimate.
- The additional cost we pay in cross validation is additional training runs, which may be too expensive in certain cases.

Lin Reg main
scores = cross_val_score(lin_reg, wine_features_tr, wine_labels,
scoring="neg_mean_squared_error", cv=10)
lin_reg_mse_scores = -scores
display_scores(lin_reg_mse_scores)

- Forest / Trees
- Random forest model builds multiple decision trees on randomly selected features and then average their predictions.
 - Building a model on top of other model is called ensemble learning, which is often used to improve performance of ML models.

Model diagnosis	Remedy
Underfitting train ↑	→ <u>Forest, Tree, MLP, NN</u> Models with more capacity Less constraints/regularization
Overfitting train ↓ test ↑	✓ More data ✓ Simpler model LR, R, L ✓ More constraints/regularization HPT

- 1) Which of the following can be used as an evaluation measure for regression models?

- pipeline
- confusion matrix
- Mean square error
- None of the above

- 2) Mention TRUE or FALSE : Cross validation provides a separate MSE for each validation set.

- TRUE
- FALSE

- 3) Building a model on top of another model is called-

- random forest
- deep learning
- stepwise learning
- ensemble learning

- 4) What is the remedy for underfitting?

- to have less constraints/regularization
- to have more constraints/regularization
- delete a significant portion of data
- None of the above

L5: Finetuning ML Models → Hyperparameter tuning (HPT)

- Usually there are a number of hyperparameters in the model, which are set manually.
- Tuning these hyperparameters lead to better accuracy of ML models.
- Finding the best combination of hyperparameters is a search problem in the space of hyperparameters, which is huge.

Grid search

- Scikit-Learn provides a class GridSearchCV that helps us in this pursuit.

```
1 from sklearn.model_selection import GridSearchCV
```

We need to specify a list of hyperparameters along with the range of values to try.
It automatically evaluates all possible combinations of hyperparameter values using cross-validation.

For example, there are number of hyperparameters in RandomForest regression, such as:

- 2 → Number of estimators → n_estimators
2 → Maximum number of features → max_features

```
param_grid = [
    {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
    {'bootstrap': [False], 'n_estimators': [3, 10], 'max_features': [2, 3, 4]},
]
```

better

Randomised Search CV → Randomised search in all the parameters

grid_search.best_params_

grid_search.best_estimator_

1) The GridSearchCV class is a part of which library?

- Pillow
- Scikit-Learn
- NumPy
- BeautifulSoup

2) Which of the following are the hyperparameters in the RandomForest regression?

- number of estimators
- maximum number of features
- precision
- recall

3) Mention TRUE or FALSE : Initializing GridSearchCV with refit = True option retains the best estimators on the full training set.

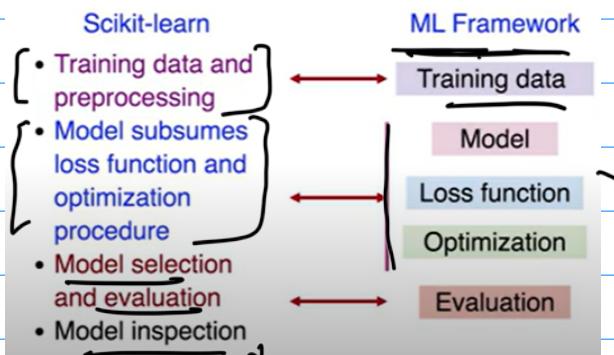
- TRUE
- FALSE

4) For a large hyperparameter space, it is desirable to try which of the following options-

- RandomizedSearchCV
- GridSearchCV
- MyCV
- StratifiedCV

L6: Intro to Scikit Learn

sklearn APIs are organized on the lines of our ML framework.



standard principle
sklearn APIs are well designed with the following principles:

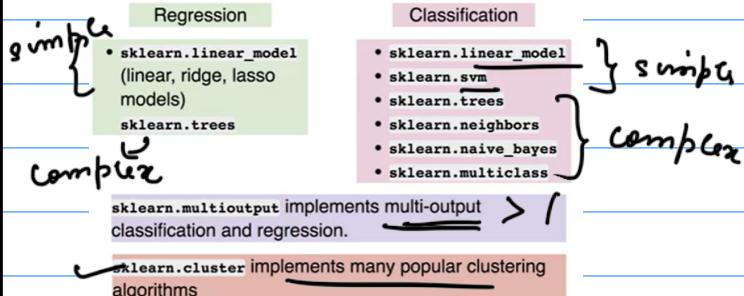
- Consistency:** All APIs share a simple and consistent interface.
- Inspection:** The learnable parameters as well as hyperparameters of all estimator's are accessible directly via public instance variables.
- Nonproliferation of classes:** Datasets are represented as Numpy arrays or Scipy sparse matrix instead of custom designed classes.
- Composition:** Existing building blocks are reduced as much as possible. pipeline, f1, CT
- Sensible defaults** values are used for parameters that enables quick baseline building.

Types of sklearn objects

Transformers	Estimators	Predictors
<ul style="list-style-type: none"> transforms dataset transform() for transforming dataset. fit() learns parameters. fit_transform() fits parameters and transform() the dataset. 	<ul style="list-style-type: none"> Estimates model parameters based on training data and hyper parameters. fit() method 	<ul style="list-style-type: none"> Makes prediction on dataset predict() method that takes dataset as an input and returns predictions. score() method to measure quality of predictions.

Module	Functionality
sklearn.datasets	Loading datasets - custom as well as popular reference dataset.
sklearn.preprocessing	Scaling, centering, normalization and binarization methods
sklearn.impute	Filling missing values
sklearn.feature_selection	Implements feature selection algorithms
sklearn.feature_extraction	Implements feature extraction from raw data.

Implements supervised and unsupervised models



Hint

```
1 import sklearn.linear_model import LogisticRegression
```

- Estimators
- Predictors
- Transformers
- Random forest

2) Mark the correct statements:

- predict() method takes a feature matrix as an input and returns predictions
- score() method measures quality of predictions for a given test set
- fit_transform() is used on the test data X
- fit_transform() method fits the parameters and uses them to transform() the dataset

3) The function of sklearn.linear model is to implement:

- classical algorithms like logistic regression, linear regression etc.
- the k-nearest neighbors algorithm
- multi-class classification models
- multilabel classification and regression

? LR

4) The function to implement the Naive-bayes classification is-

- sklearn.naive_Bayes
- sklearn.naive_bayes
- sklearn.naiveBayes
- sklearn.naive bayes

learn

5) Suppose we use the following code for creating a random dataset of 10 samples with 2 features each. What is the 'y' or label vector?

```
{ from sklearn.datasets import make_blobs  
X,y = make_blobs(n_samples=10, centers=3, n_features=2, random_state=42)  
○ y = array([0,1,1,1,0,0,1,1,0,0])  
○ y = array([1,1,0,1,1,1,1,1,1])  
○ y = array([1,0,0,0,1,0,1,1,0,0])  
 y = array([2,2,1,2,0,0,0,1,1,0])
```

6) Which of the following is a supervised learning problem?

- Recommending relevant product ads to a website visitor based on previous purchases or product search made on the website
- You have been given heights of different students of the class. You have to group them into three categories i.e. tall, medium and short.

A bank looking to decide the amount of personal loan to its customers based on transactions in their account

○ A robot has been placed in one corner of a grid-world with a 3x3 grid. The grid-world has an obstacle in the centre and goal in the other corner (diametrically opposite to the agent). There is a reward of +3, -1 and -10 for reaching the goal, one step in any direction and for hitting the obstacle.

7) Which of the following is a regression task?

- Recommending relevant product ads to a website visitor based on previous purchases or product search made on the website.
 - You have been given heights of different students of the class. You have to group them into three categories i.e. tall, medium and short.
- A bank looking to decide the amount of personal loan to its customers based on transactions in their account. It has already identified potential customers but now wants to ascertain loan amounts for each customer.
- Predicting rainfall this year based on data available for the last 30 years.

8) Which one of the following statement is true?

- I. A regression model that uses L1 regularization technique is called Ridge Regression
- II. A regression model that uses L2 regularization technique is called Lasso Regression
- Statement I is true but II is false
- Statement II is true but I is false
- Both statements are true
- Both statements are false

9) We need to divide the sample data into train and test data because

- We do not need to train on the entire data
- The training happens faster on smaller set
- For ensuring that model has low bias and generalizes better on test/unseen data

For ensuring the model has low variance and generalizes better on test/unseen data

low variance, high bias → DL

10) Which function is used for fully describing the dataset?

- .shape → (1000, 1)
- .head ✓ 5 rows
- .fulldescription ↗
- .DESCR / describe() ↗ statistical all data

Watch Practical Lectures of Data Loading and Demo of sklearn dataset API

1) In the iris dataset, what is the output of the following code?

```
from sklearn.datasets import load_iris  
data = load_iris()  
data.feature_names
```

['sepal_length (cm)', 'sepal_width (cm)', 'petal_length (cm)', 'petal_width (cm)']

['sepal_width (cm)', 'sepal_length (cm)', 'petal_width (cm)', 'petal_thickness (cm)']

['sepal_width (cm)', 'petal_length (cm)', 'petal_width (cm)', 'petal_thickness (cm)']

['petal_width (cm)', 'petal_length (cm)', 'petal_thickness (cm)', 'sepal_thickness (cm)']

2) With which of the following codes, we can examine the shape of the feature matrix?

- data=load_iris()
data.shape
- data=load_iris()
data.data.shape
- data=load_iris()
data.data-data.shape
- data=load_iris()
data.shape

1 point

3) The first four examples of the feature matrix can be accessed via which of the following codes?

- data.data[5]
- data.data[4]
- data.data[4:1]
- data.data[1:4]

1 point

4) Which of the following statements is true about openml.org?

- It is a public repository for machine learning data and experiments
- It is a private completely restricted repository for machine learning data and experiments
- openml.org allows everybody to upload open datasets.
- OpenML does not automatically analyze the data.

2 points

5) Mention TRUE or FALSE:
make_blobs enables us to generate random data for clustering.

- TRUE
- FALSE

1 point

Ridge → L2 reg.
lasso → L1 reg.

AQ 1.7

MLP Week 2 Notes

L1: Data Preprocessing

The real world training data is usually not clean and has many issues such as **missing values** for certain features, features on different scales, non-numeric attributes etc.

Often there is a need to **pre-process** the data to make it amenable for training the model.

Sklearn provides a rich set of transformers for this job.

The **same pre-processing** should be applied to both training and test set.

Sklearn provides **pipeline** for making it easier to chain multiple transforms together and apply them **uniformly** across train, eval and test sets.

Once you get the training data, the first job is to **explore the data** and list down **preprocessing** needed.

Typical problems include

Missing values in features

Numerical features are **not on the same scale**.

Categorical attributes need to be represented with sensible numerical representation.

Too many features, reduce them.

Extract features from non-numeric data.

- Data cleaning (**sklearn.preprocessing**) such as **standardization**, **missing value imputation**, etc.
- Feature extraction (**sklearn.feature_extraction**)
- Feature reduction (**sklearn.decomposition.pca**)
- Feature expansion (**sklearn.kernel_approximation**)

Each transformer has the following methods:

• **fit()** method learns model parameters from a training set.

• **transform()** method applies the learnt transformation to the new data.

• **fit_transform()** performs function of both **fit()** and **transform()** methods and is more convenient and efficient to use.

sklearn.feature_extraction has useful APIs to extract features from data:

DictVectorizer

FeatureHasher

DictVectorizer

Converts lists of mappings of feature name and feature value, into a matrix.

Original data

```
data = [{"age": 4, "height":96.0},  
        {"age": 1, "height":73.9},  
        {"age": 3, "height":88.9},  
        {"age": 2, "height":81.6}]
```

Transformed feature matrix \mathbf{X}'

$$\mathbf{X}'_{4 \times 2} = \begin{bmatrix} 4 & 96.0 \\ 1 & 73.9 \\ 3 & 88.9 \\ 2 & 81.6 \end{bmatrix}$$

`dv = DictVectorizer(sparse=False)
dv.fit_transform(data)`

FeatureHasher

- High-speed, low-memory vectorizer that uses **feature hashing** technique.
- Instead of building a hash table of the features, as the vectorizers do, it **applies a hash function** to the features to **determine their column index** in sample matrices directly.
- This results in **increased speed and reduced memory usage**, at the **expense of inspectability**; the hasher does not remember what the input features looked like and has **no inverse_transform** method.
- Output of this transformer is **scipy.sparse** matrix.

• **sklearn.feature_extraction.image.*** has useful APIs to extract features from image data. Find out more about them in sklearn user guide at the following link: [Feature Extraction from Images](#).

• **sklearn.feature_extraction.text.*** has useful APIs to extract features from text data. Find out more about them in sklearn user guide at the following link: [Feature Extraction from Text](#).

Missing values occur due to **errors in data capture** such as **sensor malfunctioning, measurement errors** etc.

Many ML algorithms do not work with missing data and need all features to be present.

Discarding records containing missing values would result in loss of valuable training samples.

sklearn.impute API provides functionality to fill missing values in a dataset.

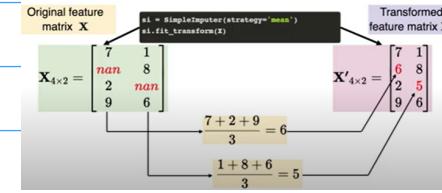
SimpleImputer

KNNImputer

MissingIndicator provides indicators for missing values.

SimpleImputer

- Fills missing values with one of the following strategies: **'mean'**, **'median'**, **'most_frequent'** and **'constant'**.



KNNImputer

- Uses **k-nearest neighbours** approach to fill missing values in a dataset.

▪ The missing value of an attribute in a specific example is filled with the **mean** value of the same attribute of **n_neighbors** closest neighbors.

- The nearest neighbours are decided based on **Euclidean distance**.

Let's fill the missing value in first sample/row. $\mathbf{X}_{4 \times 3} = \begin{bmatrix} 1. & 2. & \text{nan} \\ 3. & 4. & 3. \\ \text{nan} & 6. & 5. \\ 8. & 8. & 7. \end{bmatrix}$

Distance with [1. 2. nan.]

$\sqrt{(1-3)^2 + (2-4)^2} \approx 2.82$
 $\sqrt{(2-6)^2} = 4$
 $\sqrt{(1-8)^2 + (2-8)^2} \approx 9.21$

2 nearest neighbours

Values of the feature from 2 nearest neighbours

$\frac{3+5}{2} = 4$
of neighbours

Original feature matrix \mathbf{X}

$$\mathbf{X}_{4 \times 4} = \begin{bmatrix} 1. & 2. & \text{nan} \\ 3. & 4. & 3. \\ \text{nan} & 6. & 5. \\ 8. & 8. & 7. \end{bmatrix}$$

Transformed feature matrix \mathbf{X}'

$$\mathbf{X}'_{4 \times 4} = \begin{bmatrix} 1. & 2. & 4. \\ 3. & 4. & 3. \\ 5.5 & 6. & 5. \\ 8. & 8. & 7. \end{bmatrix}$$

`knni = KNNImputer(n_neighbors=2, weights="uniform")
knni.fit_transform(X)`

1) What will be the output of the following code:

```
data = [{"#courses_opted": 4, "GPA": 6.7}, {"#courses_opted": 4, "GPA": 8.5},  
        {"#courses_opted": 2, "GPA": 8.1}, {"#courses_opted": 1, "GPA": 7.2},  
        {"#courses_opted": 1, "GPA": 9.1}]
```

```
from sklearn.feature_extraction import DictVectorizer  
dv = DictVectorizer(sparse=False)  
data_transformed = dv.fit_transform(data)  
data_transformed[2]
```

4) import numpy as np
from sklearn.impute import SimpleImputer
imp = SimpleImputer(missing_values=np.nan, strategy='median')
 $X = [[4, 1], [np.nan, 5], [8, 0]]$
imp.fit_transform(X)

What will be the missing value in X after the execution of the above code:

- 4
 12
 6
 5

2) Which API provides functionality to fill missing values in a dataset?

- sklearn.impute
 sklearn.FillMissingValue
 sklearn.extraction
 sklearn.DictVectorizer
 SimpleImputer(strategy='median')
 Imputer(strategy='median')
 SimpleImputer(strategy='mean')
 Imputer(strategy='mean')

L2: Numeric Transformers

Numerical features with different scales leads to slower convergence of iterative optimization procedures.

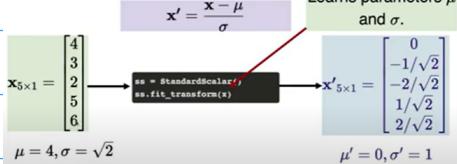
It is a good practice to scale numerical features so that all of them are on the same scale.

Let's learn how to scale numerical features with sklearn API.

Three feature scaling APIs are available in sklearn

StandardScaler

Transforms the original features vector x into new feature vector x' using following formula



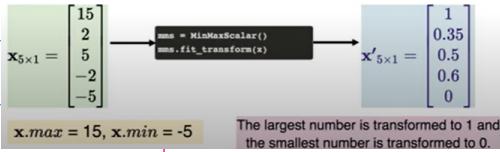
Note that the transformed feature vector x' has mean (μ) = 0 and standard deviation (σ) = 1.

MinMaxScaler

It transforms the original feature vector x into new feature vector x' so that all values fall within range [0, 1]

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

where x_{\max} and x_{\min} are largest and smallest values of that feature respectively, of the original feature vector x .

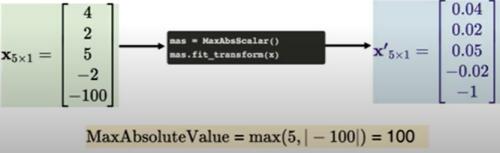


MaxAbsScaler

It transforms the original features vector x into new feature vector x' so that all values fall within range [-1, 1]

$$x' = \frac{x}{\text{MaxAbsoluteValue}}$$

where MaxAbsoluteValue = $\max(x_{\max}, |x_{\min}|)$

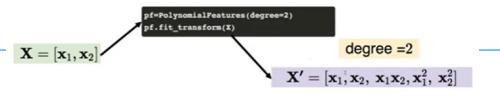


FunctionTransformer

Constructs transformed features by applying a user defined function.

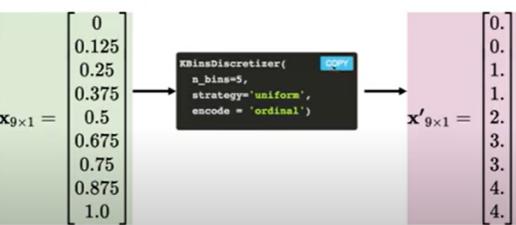


Generates a new feature matrix consisting of all polynomial combinations of the features with degree less than or equal to the specified degree.



KBinsDiscretizer

- Divides a continuous variable into bins.
- One hot encoding or ordinal encoding is further applied to the bin labels.



1) Consider the following feature vector:

$$X = \begin{bmatrix} 0 \\ 3 \\ 6 \end{bmatrix}$$

What will be the transformed feature matrix X' after the StandardScaler() transformation?

- $\begin{bmatrix} 0 \\ 3/\sqrt{6} \\ 6 \end{bmatrix}$
- $\begin{bmatrix} 0 \\ -3/\sqrt{6} \\ 3/\sqrt{6} \end{bmatrix}$
- $\begin{bmatrix} 3/\sqrt{6} \\ 0 \\ 3/\sqrt{6} \end{bmatrix}$
- $\begin{bmatrix} 0 \\ 3/\sqrt{6} \\ -3/\sqrt{6} \end{bmatrix}$

2) The new feature matrix x' using the transformer MinMaxScaler() is obtained as

- $x' = \frac{x}{x_{\max} - x_{\min}}$
- $x' = \frac{x_{\max} - x}{x_{\max} - x_{\min}}$
- $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$
- $x' = \frac{x}{x_{\max} + x_{\min}}$

4) What will be the output of the following code:

```
from sklearn.preprocessing import MaxAbsScaler
Y = [[2], [8], [-4]]
mas= MaxAbsScaler()
print(mas.fit_transform(Y))
```

[[0.2] [2.0] [-1]]

L3: Categorical Transformers

OneHotEncoder

- Encodes categorical feature or label as a one-hot numeric array.
- Creates one binary column for each of K unique values.
- Exactly one column has 1 in it and rest have 0.



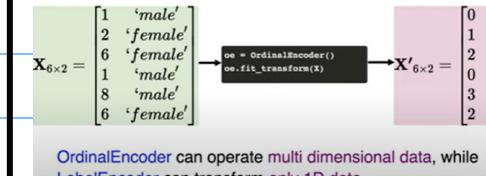
LabelEncoder

Encodes target labels with value between 0 and $K - 1$, where K is number of distinct values.



OrdinalEncoder

Encodes categorical features with value between 0 and $K - 1$, where K is number of distinct values.

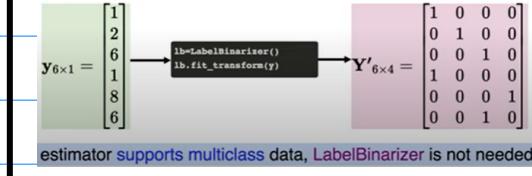


LabelBinarizer

Several regression and binary classification can be extended to multi-class setup in one-vs-all fashion.

This involves training a single regressor or classifier per class.

For this, we need to convert multi-class labels to binary labels, and LabelBinarizer performs this task.



MultLabelBinarizer

Encodes categorical features with value between 0 and $K - 1$, where K is number of classes.

In this example $K = 4$, since there are only 4 genres of movies.

```
movie_genres =
[{'action', 'comedy'},
{'comedy'},
{'action', 'thriller'},
{'science-fiction', 'action', 'thriller'}]
```

$$X'_{4 \times 4} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

`m1b = MultLabelBinarizer()`

`m1b.fit_transform(movie_genres)`

`add_dummy_feature`

Augments dataset with a column vector, each value in the column vector is 1.

$$X'_{4 \times 2} = \begin{bmatrix} 7 & 1 \\ 1 & 8 \\ 2 & 0 \\ 9 & 6 \end{bmatrix} \xrightarrow{\text{add_dummy_feature}(X)} X'_{4 \times 3} = \begin{bmatrix} 1 & 7 & 1 \\ 1 & 1 & 8 \\ 1 & 2 & 0 \\ 1 & 9 & 6 \end{bmatrix}$$

1) Consider the following feature matrix:

$$X = \begin{bmatrix} 4 \\ 3 \\ 2 \\ 3 \\ 1 \end{bmatrix}$$

What will be the dimension of transformed feature matrix X' after the OneHotEncoder() transformation?

- 4 × 4
- 4 × 5
- 5 × 5
- 5 × 4

2) Consider the following categorical feature matrix X and its transformed feature matrix X' after the MultLabelBinarizer() transformation:

$$X = \begin{bmatrix} \{\text{Heart-disease}\} & \{\text{Diabetes}\} \\ \{\text{High-BP}\} & \{\text{Diabetes}\} \\ \{\text{Heart-disease}\} & \{\text{High-BP}\} \\ \{\text{Diabetes}\} & \{\text{High-BP}\} \end{bmatrix}$$

$$X' = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Which column in X' refers to the label 'Diabetes'?

3

3) What does the transformer LabelBinarizer() do?

- Converts multi-class labels to binary labels.
- Converts multi-label feature to the binary label.
- Encodes categorical features with value between 0 and $K - 1$, where K is number of classes.
- Encodes categorical features with values -1 and 1.

4) What will be the output of the following code:

```
import numpy as np
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import add_dummy_feature

imp = SimpleImputer(missing_values=np.nan, strategy='mean')
X = [[3, 8], [1, 4], [8, np.nan]]
print(add_dummy_feature(imp.fit_transform(X)))
```

- [[3, 8], [1, 4], [8, 6]]
- [[1, 3, 8], [1, 1, 4], [1, 8, np.nan]]
- [[1, 3, 8], [1, 1, 4], [1, 8, 6]]
- [[1, 3, 8], [1, 1, 4], [1, 8, 4, 1]]

5) What will be the output of the following code:

```
from sklearn.preprocessing import LabelBinarizer
lb = LabelBinarizer()
print(lb.fit_transform(['yes', 'yes', 'no', 'yes']))
```

- [[1], [0], [0], [1]]
- [[1], [1], [0], [1]]
- [[0], [1], [1], [0]]
- [[0], [0], [1], [1]]

Univariate feature selection

Univariate feature selection selects features based on univariate statistical tests.

There are three APIs for univariate feature selection:

SelectKBest

Removes all but the k highest scoring features

SelectPercentile

Removes all but a user-specified highest scoring percentage of features

GenericUnivariateSelect

Performs univariate feature selection with a configurable strategy, which can be found via hyper-parameter search.

sklearn provides one more class of univariate feature selection methods that work on common univariate statistical tests for each feature:

`SelectFpr` selects features based on a false positive rate test.

`SelectFdr` selects features based on an estimated false discovery rate.

`SelectFwe` selects features based on family-wise error rate.

Univariate scoring function

• Each API need a scoring function to score each feature.

• Three classes of scoring functions are proposed:

Mutual information (MI) Chi-square F-statistics

• MI and F-statistics can be used in both classification and regression problems.

`mutual_info_regression`

`mutual_info_classif`

`f_regression`

`f_classif`

• Chi-square can be used only in classification prob

`chi2`

Chi-square

Mutual information (MI)

- Measures dependency between two variables.
- It returns a non-negative value.
 - MI = 0 for independent variables.
 - Higher MI indicates higher dependency.

- Measures dependence between two variables.
- Computes chi-square stats between non-negative feature (boolean or frequencies) and class label.
- Higher chi-square values indicates that the features and labels are likely to be correlated.

SelectKBest

```
skb = SelectKBest(chi2, k=20)
X_new = skb.fit_transform(X, y)
```

Selects 20 best features based on chi-square scoring function.

SelectPercentile

```
sp = SelectPercentile(chi2, percentile=20)
X_new = sp.fit_transform(X, y)
```

Selects top 20 percentile best features based on chi-square scoring function. .

GenericUnivariateSelect

```
transformer = GenericUnivariateSelect(chi2, mode='k_best', param=20)
X_new = transformer.fit_transform(X, y)
```

- Selects set of features based on a feature selection mode and a scoring function.
- The mode could be 'percentile' (default), 'k_best', 'fpr', 'fdr', 'fwe'.
- The param argument takes value corresponding to the mode.

L4: Filter Based Feature Selection

VarianceThreshold

Removes all features with variance below a certain threshold, as specified by the user, from input feature matrix

By default removes a feature which has same value, i.e. zero variance.

1) Which of the following is/are the APIs for the feature selection based on the filter based methods?

- RFE
- VarianceThreshold
- SelectFromModel
- SelectPercentile

2) What task does the API VarianceThreshold do?

- Removes all features with variance below a certain threshold from input feature matrix.
- Removes all features with variance above a certain threshold from input feature matrix.
- Only removes all features with zero variance.
- Normalize the features to reduce the variance.

3) What task does the API SelectKBest do?

- Removes all the k highest scoring features.
- Removes all but the k highest scoring features.
- Removes all but a user-specified highest scoring percentage of features.
- Removes all user-specified highest scoring percentage of features.

4) What will the following code implement?


```
transformer = GenericUnivariateSelect(chi2, mode='k_best', param=20)
X_new = transformer.fit_transform(X, y)
```

- Selects top 20 percentile best features based on chi-square scoring function.
- Selects top 20 percent best features based on chi-square scoring function.
- Selects 20 best features based on chi-square scoring function.
- Selects 20 best features based on F-statistics scoring function.

5) Which of the following option(s) is/are correct about the scoring function in univariate feature selection methods?

- Chi-square can be used only in classification problems.
- Chi-square can be used in both regression and classification problems.
- F-statistics can only be used in classification problems.
- F-statistics can be used in both classification and regression problems.

6) What will be the output of the following code?


```
from sklearn.feature_selection import VarianceThreshold
import numpy as np

X = np.array([[1, 1], [1, 3, 4], [1, 2, 4]])

vt = VarianceThreshold()
vt.fit_transform(X)
```

- [[1 1 1]]
- [[1 3 4]]
- [[1 2 4]]
- [[1 3 4]]
- [[3 4]]
- [[2 4]]
- [[1 1]]
- [[3 4]]
- [[2 4]]

The direction parameter controls whether forward or backward SFS is used.
 In general, forward and backward selection do not yield equivalent results.
 Select the direction that is efficient for the required number of selected features:
 When we want to select 7 out of 10 features,
 Forward selection would need to perform 7 iterations.
 Backward selection would only need to perform 3.
 Backward selection seems to be a reasonable choice here.
 SFS does not require the underlying model to expose a `coef_` or `feature_importances_` attributes unlike in RFE and SelectFromModel.
 SFS may be slower than RFE and SelectFromModel as it needs to evaluate more models compared to the other two approaches.

For example in backward selection, the iteration going from m features to $m - 1$ features using k -fold cross-validation requires fitting $m \times k$ models, while
 RFE would require only a single fit, and
 SelectFromModel performs a single fit and requires no iterations.

1) Which of the following options are correct about the recursive feature elimination(RFE)?

- It is a filter based feature selection method.
- It is a wrapper based feature selection method.
- RFE works by searching for a subset of features by starting with all features in the training dataset and successively removing features until the desired number remains.
- Selects desired number of important features (as specified with max_features parameter) above certain threshold of feature importance as obtained from the trained estimator

2) The argument of the parameter threshold in the class SelectFromModel can be

- integer
- 'mean'
- 'median'
- Can not be a string like 'mean'.

3) Which approach is used in the class SequentialFeatureSelection()?

- Only forward selection approach
- Only backward selection approach
- Both forward and Backward
- None of the above

4) If we want to select three features out of ten features using SequentialFeatureSelection(), which approach will be more reasonable?

- Forward selection approach
- Backward selection approach

L5: Wrapper Based Feature Selection

Recursive Feature Elimination (RFE)

- Uses an estimator to recursively remove features.
 - Initially fits an estimator on all features.
- Obtains feature importance from the estimator and removes the least important feature.
- Repeats the process by removing features one by one, until desired number of features are obtained.

- Use `RFECV` if we do not want to specify the desired number of features in `RFE`.
- It performs `RFE` in a cross-validation loop to find the optimal number of features.

SelectFromModel

Selects desired number of important features (as specified with `max_features` parameter) above certain threshold of feature importance as obtained from the trained estimator.

- The feature importance is obtained via `coef_`, `feature_importances_` or an importance_getter callable from the trained estimator
- The feature importance threshold can be specified either numerically or through string argument based on built-in heuristics such as 'mean', 'median' and float multiples of these like '0.1*mean'.

Let's look at a concrete example of SelectFromModel

```
clf = LinearSVC(C=0.01, penalty="l1", dual=False)
clf = clf.fit(X, y)
clf.coef_
model = SelectFromModel(clf, prefit=True)
X_new = model.transform(X)
```

- Here we use a linear support vector classifier to get coefficients of features for `SelectFromModel` transformer.
- It ends up selecting features with non-zero weights or coefficients.

Sequential feature selection

Performs feature selection by selecting or deselecting features one by one in a greedy manner.

Uses one of the two approaches

Forward selection	Backward selection
Starting with a zero feature, it finds one feature that obtains the best cross validation score for an estimator when trained on that feature.	Starting with all features and removes least important features one by one following the idea of forward selection.
Repeats the process by adding a new feature to the set of selected features.	
Stops when reach the desired number of features.	

L6: Heterogenous Feature Transformations

Generally training data contains diverse features such as numeric and categorical.

Different feature types are processed with different transformers.

Need a way to combine different feature transformers seamlessly.

Composite Transformer

`sklearn.compose` has useful classes and methods to apply transformation on subset of features and combine them:

ColumnTransformer

- It applies a set of transformers to columns of an array or `pandas.DataFrame`, concatenates the transformed outputs from different transformers into a single matrix.
- It is useful for transforming heterogenous data by applying different transformers to separate subsets of features.
- It combines different feature selection mechanisms and transformation into a single transformer object.

- `ColumnTransformer()`
- Each tuple has format
 - (`estimatorName`, `estimator(...)`, `columnIndices`)

```
column_trans = ColumnTransformer(
    [('ageScaler', CountVectorizer(), [0]),
     ('genderEncoder', OneHotEncoder(dtype='int'), [1])],
    remainder='drop', verbose_feature_names_out=False)
```

Consider following feature matrix, which represent weight and gender of a class of students.

$$\mathbf{X}_{6 \times 2} = \begin{bmatrix} 20.0 & \text{'male'} \\ 11.2 & \text{'female'} \\ 15.6 & \text{'female'} \\ 13.0 & \text{'male'} \\ 18.6 & \text{'male'} \\ 16.4 & \text{'female'} \end{bmatrix}$$

Here, first column is numeric, however, second column is categorical, therefore different transformers have to be applied on them.

In this example, lets apply `MaxAbsScaler` on the numeric column and `OneHotEncoder` on categorical column.

```
column_trans = ColumnTransformer([
    ('scaler', MaxAbsScaler(), [0]),
    ('encoder', OneHotEncoder(dtype='int'), [1])
], remainder='drop', verbose_feature_names_out=True)
column_trans.fit_transform(X)
```

$$\mathbf{X}'_{6 \times 2} = \begin{bmatrix} 20.0 & \text{'male'} \\ 11.2 & \text{'female'} \\ 15.6 & \text{'female'} \\ 13.0 & \text{'male'} \\ 18.6 & \text{'male'} \\ 16.4 & \text{'female'} \end{bmatrix} \quad \mathbf{X}'_{6 \times 3} = \begin{bmatrix} 1. & 0. & 1. \\ 0.56 & 1. & 0. \\ 0.78 & 1. & 0. \\ 0.65 & 0. & 1. \\ 0.93 & 0. & 1. \\ 0.82 & 1. & 0. \end{bmatrix}$$

Transforming Target for Regression

TransformedTargetRegressor

- Transforms the target variable y before fitting a regression model.
- The predicted values are mapped back to the original space via an inverse transform.
- `TransformedTargetRegressor` takes `regressor` and `transformer` to be applied to the target variable as arguments.

```
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.compose import TransformedTargetRegressor
tt = TransformedTargetRegressor(regressor=LinearRegression(),
                                func=np.log, inverse_func=np.exp)
X = np.arange(4).reshape(-1, 1)
y = np.exp(2 * X).ravel()
tt.fit(X, y)
```

1) Which of the following modules in sklearn deals with applying transformation on subset of diverse features (such as numerical and categorical both) and combine them?

- sklearn.impute
 sklearn.feature_selection
 sklearn.model_selection
 sklearn.compose

2) What will be the output of the following code:

```
import numpy as np
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import MinMaxScaler, OrdinalEncoder
X = np.array([[1, yes], [2, no], [3, no]])
ct = ColumnTransformer([('scaler', MinMaxScaler(), [0]), ('pass', 'passthrough', [1]), ('encoder', OrdinalEncoder(), [2])])
print(ct.fit_transform(X))
[[0.1 1 1] [0.5 2 0] [0.0 3 0]
 [[0.0 1 1] [0.25 2 0] [0.1 3 0]
 [[0.0 1 1] [0.5 2 1] [1.0 3 0]
 [[0.0 1 1] [0.5 2 0] [1.0 3 0]]]
```

1) Which module in sklearn is used for dimensionality reduction techniques?

- sklearn.impute
 sklearn.feature_selection
 sklearn.compose
 sklearn.decomposition

2) Which of the following statements is/ are correct about the principle component analysis (PCA)?

- It is a linear dimensionality reduction technique.
 It uses singular value decomposition (SVD) to project the feature matrix or data to a lower dimensional space.
 It uses stochastic methods to project the feature matrix or data to a lower dimensional space.
 The subsequent PCs are orthogonal to the first PC.

L8: Chaining Transformers

The preprocessing transformations are applied one after another on the input feature matrix.

```
si = SimpleImputer()
X_imputed = si.fit_transform(X)
ss = StandardScaler()
X_scaled = ss.fit_transform(X_imputed)
```

It is important to apply exactly same transformation on training, evaluation and test set in the same order.

Failing to do so would lead to incorrect predictions from model due to distribution shift and hence incorrect performance evaluation.

The `sklearn.pipeline` module provides utilities to build a composite estimator, as a chain of transformers and estimators.

There are two classes: (i) `Pipeline` and (ii) `FeatureUnion`.

Class	Usage
Pipeline	Constructs a chain of multiple transformers to execute a fixed sequence of steps in data preprocessing and modelling.
FeatureUnion	Combines output from several transformer objects by creating a new transformer from them.

sklearn.pipeline.Pipeline

Sequentially apply a list of `transformers` and `estimators`.

Intermediate steps of the pipeline must be 'transformers' that is, they must implement `fit` and `transform` methods.

The final estimator only needs to implement `fit`.

The purpose of the pipeline is to assemble several steps that can be cross-validated together while setting different parameters.

Pipeline()

- It takes a list of `('estimatorName', estimator(...))` tuples.
- The pipeline object exposes interface of the last step.

```
estimators = [
    ('simpleimputer', SimpleImputer()),
    ('pca', PCA()),
    ('regressor', LinearRegression())
]
pipe = Pipeline(steps=estimators)
```

Accessing individual steps in Pipeline

```
estimators = [
    ('simpleimputer', SimpleImputer()),
    ('pca', PCA()),
    ('regressor', LinearRegression())
]
pipe = Pipeline(steps=estimators)
```

Total # steps: 3

- SimpleImputer
- PCA
- LinearRegression

The second estimator can be accessed in following 4 ways:

- `pipe.named_steps.pca`
- `pipe.steps[1]`
- `pipe[1]`
- `pipe['pca']`

Accessing parameters of each step in Pipeline

Parameters of the estimators in the pipeline can be accessed using the `<estimator>__<parameterName>` syntax, note there are two underscores between `<estimator>` and `<parameterName>`

```
estimators = [
    ('simpleimputer', SimpleImputer()),
    ('pca', PCA()),
    ('regressor', LinearRegression())
]
pipe = Pipeline(steps=estimators)
pipe.set_params(pca__n_components = 2)
```

In above example `n_components` of `PCA()` step is set after the pipeline is created.

L7: Dimensionality Reduction by PCA

Another way to reduce the number of feature is through unsupervised dimensionality reduction techniques.

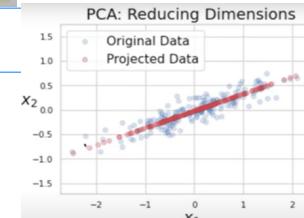
`sklearn.decomposition` module has a number of APIs for this task.

We will focus on how to perform feature reduction with principle component analysis (PCA) in sklearn.

PCA 101

- PCA, is a linear dimensionality reduction technique.
- It uses singular value decomposition (SVD) to project the feature matrix or data to a lower dimensional space.
- The first principle component (PC) is in the direction of maximum variance in the data.
 - It captures bulk of the variance in the data.
- The subsequent PCs are orthogonal to the first PC and gradually capture lesser and lesser variance in the data.
- We can select first k PCs such that we are able to capture the desired variance in the data.

`sklearn.decomposition.PCA` API is used for performing PCA based dimensionality reduction.



Performing grid search with pipeline

By using naming convention of nested parameters, grid search can implemented.

```
param_grid = dict(imputers='passthrough',
                  simpleImputer=SimpleImputer(),
                  ENNImputer=ENNImputer(),
                  clf=[SVC(), LogisticRegression()],
                  clf_C=[0.1, 10, 100])
grid_search = GridSearchCV(pipe, param_grid=param_grid)
```

- **c** is an inverse of regularization, lower its value stronger the regularization is.
- In the example above **clf_c** provides a set of values for grid search.

- Combines multiple steps of end to end ML into single object such as missing value imputation, feature scaling and encoding, model training and cross validation.
- Enables joint grid search over parameters of all the estimators in the pipeline.
- Makes configuring and tuning end to end ML quick and easy.
- Offers convenience, as a developer has to call `fit()` and `predict()` methods only on a `Pipeline` object (assuming last step in the pipeline is an estimator).
- Reduces code duplication: With a `Pipeline` object, one doesn't have to repeat code for preprocessing and transforming the test data.

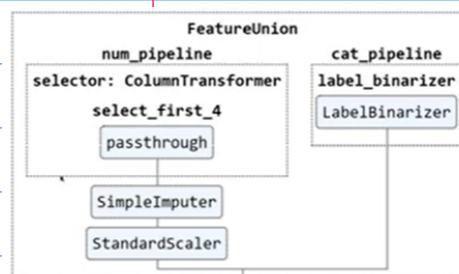
sklearn.pipeline.FeatureUnion

- Concatenates results of multiple transformer objects.
- Applies a list of transformer objects in parallel, and their outputs are concatenated side-by-side into a larger matrix.
- `FeatureUnion` and `Pipeline` can be used to create complex transformers.

Combining Transformers and Pipelines

- `FeatureUnion()` accepts a list of tuples.
- Each tuple is of the format:
 - ('estimatorName', estimator(...))

```
num_pipeline = Pipeline([('selector', ColumnTransformer([('select_first_4',
                                                       'passthrough',
                                                       slice(0,4))]),
                        ('imputer', SimpleImputer(strategy="median")),
                        ('std_scaler', StandardScaler())),
                       ])
cat_pipeline = ColumnTransformer([('label_binarizer', LabelBinarizer(), [4]),
                                  ])
full_pipeline = FeatureUnion(transformer_list=
    [("num_pipeline", num_pipeline),
     ("cat_pipeline", cat_pipeline),])
```



- 1) What task does the module `sklearn.pipeline` do?
- It provides utilities to build a composite estimator, as a chain of transformers and estimators.
 - It provides the utilities to build the estimators only on test sets.
 - It provides the utilities to build the estimators only on training sets.
 - It provides utilities to reduce the dimension of training dataset.

- 2) What will be the output of the following code:

```
from sklearn.decomposition import PCA
from sklearn.impute import SimpleImputer
from sklearn.linear_model import LinearRegression
from sklearn.pipeline import Pipeline
estimators = [
    (simpleImputer, SimpleImputer()),
    (pca, PCA()),
    (regressor, LinearRegression())
]
pipe = Pipeline(steps=estimators)
print(len(pipe.steps))
```

- 3) What are the two ways to create a pipeline object?

- `Pipeline()` and `Create_Pipeline()`
- `Pipeline()` and `Pipe()`
- `Pipeline()` and `Join_Pipeline()`
- `Pipeline()` and `make_pipeline()`

- 4) Which of the following statements is/are correct about the class `sklearn.pipeline.FeatureUnion`?

- It provides utilities to build a composite estimator, as a chain of transformers and estimators.
- It concatenates results of multiple transformer objects.
- The list of transformer objects can be executed in parallel.
- It provides the utilities to build the estimators only on test sets.

L9, 10, 11: Demo Lectures

L1: Linear Regression

How to build baseline regression model?

DummyRegressor helps in creating a baseline for regression.

```
1 from sklearn.dummy import DummyRegressor
2
3 dummy_regr = DummyRegressor(strategy="mean")
4 dummy_regr.fit(X_train, y_train)
5 dummy_regr.predict(X_test)
6 dummy_regr.score(X_test, y_test)
```

- It makes a prediction as specified by the strategy.
- Strategy is based on some statistical property of the training set or user specified value.

Strategy mean median quantile constant

LR How is Linear Regression model trained?

Step 1: Instantiate object of a suitable linear regression estimator from one of the following two options

Normal equation

```
1 from sklearn.linear_model import LinearRegression
2 linear_regressor = LinearRegression()
```

Iterative optimization

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor()
```

Step 2: Call fit method on linear regression object with training feature matrix and label vector as arguments.

```
1 # Model training with feature matrix X_train and
2 # label vector or matrix y_train
3 linear_regressor.fit(X_train, y_train)
```

SGDRegressor Estimator

- Implements stochastic gradient descent
- Use for large training set up ($> 10k$ samples)
- Provides greater control on optimization process through provision for hyperparameter settings.

loss = 'squared_error' penalty = 'l1'
 loss = 'huber' penalty = 'l2'
 loss = 'epsilon_insensitive' penalty = 'elasticnet'

SGDRegressor
 learning_rate = 'constant'
 learning_rate = 'optimal'
 learning_rate = 'invscaling'
 learning_rate = 'adaptive'

It's a good idea to use a random seed of your choice while instantiating SGDRegressor object. It helps us get reproducible results.

Set random_state to seed of your choice.

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor(random_state=42)
```

Note: In the rest of the presentation, we won't set the random seed for sake of brevity. However while coding, always set the random seed in the constructor.

SGD is sensitive to feature scaling, so it is highly recommended to scale input feature matrix.

```
1 from sklearn.linear_model import SGDRegressor
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import StandardScaler
4
5 sgd = Pipeline([
6     ('feature_scaling', StandardScaler()),
7     ('sgd_regressor', SGDRegressor())
8 ])
9 sgd.fit(X_train, y_train)
```

Note Feature scaling is not needed for word frequencies and indicator features as they have intrinsic scale.

- Features extracted using PCA should be scaled by some constant c such that the average L2 norm of the training data equals one.

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor(shuffle=True)
```

learning_rate = 'constant' learning_rate = 'invscaling'

learning_rate = 'adaptive'

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor(random_state=42)
```

What is the default setting?

learning_rate = 'invscaling' eta0 = 1e-2 power_t = 0.25

Learning rate reduces after every iteration:

eta = eta0 / pow(t, power_t)

Note: You can make changes to these parameters to speed up or slow down the training process.

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor(learning_rate='adaptive',
3                                 eta0=1e-2)
```

The learning rate is kept to initial value as long as the training loss decreases.

When the stopping criterion is reached, the learning rate is divided by 5, and the training loop continues.

The algorithm stops when the learning rate goes below 10^{-6} .

quick baseline

MLP Notes Week 3

LR (most simple model)
 SGD (more complex)

Set max_iter to desired #epochs. The default value is 1000.
 my pc p. → 1 full iteran over the entire training set

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor(max_iter=1000)
```

Remember one epoch is one full pass over the training data.

Practical tip

SGD converges after observing approximately 10^6 training samples. Thus, a reasonable first guess for the number of iterations for n sampled training set is

$$n = \text{no. of data pts}$$

$$\text{tot} = \text{tolerance value change}$$

The SGDRegressor stops

- when the training loss does not improve (loss > best_loss - tol) for n_iter_no_change consecutive epochs
- else after a maximum number of iteration max_iter.

when you want to get avg. training loss

```
1 from sklearn.linear_model import SGDRegressor
2 linear_regressor = SGDRegressor(average=True)
```

```
1 sgd_reg = SGDRegressor(max_iter=1, tol=-np.inf, warm_start=True,
2                         penalty=None, learning_rate="constant", eta0=0.0005)
3
4 for epoch in range(1000):
5     sgd_reg.fit(X_train, y_train) # continues where it left off
6     y_val_predict = sgd_reg.predict(X_val)
7     val_error = mean_squared_error(y_val, y_val_predict)
```

gmppt: $\hat{y} = w_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m = w^T x$

The weights w_1, w_2, \dots, w_m are stored in coef_ class variable.

linear_regressor.coef_ → obj. coef_

The intercept w_0 is stored in intercept_ class variable.

linear_regressor.intercept_ → value

Note These code snippets works for both LinearRegression and SGDRegressor, and for that matter to all regression estimators that we will study in this module. Why?

All of them are estimators.

Step 1: Arrange data for prediction in a feature matrix of shape (#samples, #features) or in sparse matrix format.

Step 2: Call predict method on linear regression object with feature matrix as an argument.

```
1 # Predict labels for feature matrix X_te
2 linear_regressor.predict(X_test)
```

= y-predict

Same code works for all regression estimators.

1) Which of the following are valid strategy parameter in DummyRegressor?

- mean
 median
 quantile
 None of the Above

2) Which of the following objects can be used to solve linear regression using iterative optimization?

- LR=LinearRegression()
 LR=SGDRegressor()
 LR=DummyRegressor()
 None of the Above

3) For performing Linear Regression using sklearn, which of the following module will be appropriate?

- sklearn.linear_regression
 sklearn.linear_model
 sklearn.regression_model
 sklearn.linear_fit

4) Which of the following sklearn modules is used to do necessary transformation(s) on the training data before performing polynomial regression?

- sklearn.linear_model
 sklearn.metrics
 sklearn.model_selection
 sklearn.preprocessing

SGD(eta = 42, learning_rate = 'invscaling')

learning_rate = 'invscaling'

eta = eta0 / pow(t, power_t)

no. of iterations

L2: Model Evaluation

STEP 1: Split data into train and test

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
```

STEP 2: Fit linear regression estimator on training set.

STEP 3: Calculate training error (a.k.a. empirical error)

STEP 4: Calculate test error (a.k.a. generalization error)

Compare training and test errors

Using score method on linear regression object:

```
1 # Evaluation on the eval set with
2 # 1. feature matrix
3 # 2. label vector or matrix (single/multi-output)
4 linear_regressor.score(X_test, y_test)
```

The score returns R^2 or coefficient of determination

$R^2 = 1 - \frac{u}{v}$

Sum of squared error (actual and predicted label)

total sum of square

Sum of squared error (actual and mean predicted) $v = (y - \hat{y}_{\text{mean}})^T (y - \hat{y}_{\text{mean}})$

The score returns R^2 or coefficient of determination

$R^2 = 1 - \frac{u}{v}$

When?

- The best possible score is 1.0. $u, \text{sum of squared error} = 0$

- A constant model that always predicts the expected value of y , would get a score of 0.0.

- The score can be negative (because the model can be arbitrarily worse).

sklearn provides a bunch of regression metrics to evaluate performance of the trained estimator on the evaluation set.

mean_absolute_error

```
1 from sklearn.metrics import mean_absolute_error
2 eval_score = mean_absolute_error(y_test, y_predicted)
```

mean_squared_error

```
1 from sklearn.metrics import mean_squared_error
2 eval_score = mean_squared_error(y_test, y_predicted)
```

r2_score

```
1 from sklearn.metrics import r2_score
2 eval_score = r2_score(y_test, y_predicted)
```

mean_squared_log_error

```
1 from sklearn.metrics import mean_squared_log_error
2 eval_score = mean_squared_log_error(y_test, y_predicted)
```

- Useful for targets with exponential growths like population, sales growth etc.

- Penalizes under-estimation heavier than the over-estimation.

mean_absolute_percentage_error

```
1 from sklearn.metrics import mean_absolute_percentage_error
2 eval_score = mean_absolute_percentage_error(y_test, y_predicted)
```

- Sensitive to relative error.

median_absolute_error

```
1 from sklearn.metrics import median_absolute_error
2 eval_score = median_absolute_error(y_test, y_predicted)
```

- Score is a metric for which higher value is better.

- Error is a metric for which lower value is better.

Convert error metric to score metric by adding neg_ suffix.

Function → Scoring

metrics.mean_absolute_error	neg_mean_absolute_error
metrics.mean_squared_error	neg_mean_squared_error
metrics.mean_squared_error	neg_root_mean_squared_error
metrics.mean_squared_log_error	neg_mean_squared_log_error
metrics.median_absolute_error	neg_median_absolute_error

In case we get comparable performance on train and test with this split, is this performance guaranteed on other splits too?

- Is test set sufficiently large?

- In case it is small, the test error obtained may be unstable and would not reflect the true test error on large test set.

- What is the chance that the easiest examples were kept aside as test by chance?

- If this happens would lead to optimistic estimation of the true test error.

We use cross validation for robust performance evaluation.

CV → cross validation

→ default (KFold)

cv = 5

clue = Value

Leave One out

R² version

$$R^2 = 1 - \frac{u}{v}$$

$$u = (Xw - y)^T (Xw - y)$$

$$(Xw - y)^T$$

$$v = (y - \hat{y}_{\text{mean}})^T (y - \hat{y}_{\text{mean}})$$

$$(y - \hat{y}_{\text{mean}})^T$$

$$(y - \hat{y}_{\text{mean}})$$

The score returns R^2 or coefficient of determination

$R^2 = 1 - \frac{u}{v}$

When?

- The best possible score is 1.0. $u, \text{sum of squared error} = 0$

- A constant model that always predicts the expected value of y , would get a score of 0.0.

- The score can be negative (because the model can be arbitrarily worse).

sklearn provides a bunch of regression metrics to evaluate performance of the trained estimator on the evaluation set.

mean_absolute_error

```
1 from sklearn.metrics import mean_absolute_error
2 eval_score = mean_absolute_error(y_test, y_predicted)
```

mean_squared_error

```
1 from sklearn.metrics import mean_squared_error
2 eval_score = mean_squared_error(y_test, y_predicted)
```

r2_score

```
1 from sklearn.metrics import r2_score
2 eval_score = r2_score(y_test, y_predicted)
```

mean_squared_log_error

```
1 from sklearn.metrics import mean_squared_log_error
2 eval_score = mean_squared_log_error(y_test, y_predicted)
```

- Useful for targets with exponential growths like population, sales growth etc.

- Penalizes under-estimation heavier than the over-estimation.

mean_absolute_percentage_error

```
1 from sklearn.metrics import mean_absolute_percentage_error
2 eval_score = mean_absolute_percentage_error(y_test, y_predicted)
```

- Sensitive to relative error.

median_absolute_error

```
1 from sklearn.metrics import median_absolute_error
2 eval_score = median_absolute_error(y_test, y_predicted)
```

- Score is a metric for which higher value is better.

- Error is a metric for which lower value is better.

exponential data pts

Score, better

Error, better

```
1 from sklearn.model_selection import cross_val_score
2 from sklearn.linear_model import LinearRegression
3
4 lin_reg = LinearRegression()
5 score = cross_val_score(lin_reg, X, y, cv=5)
```

- Uses KFold cross validation iterator, that divides training data into 5 folds.

- In each run, it uses 4 folds for training and 1 for evaluation.

```
1 from sklearn.model_selection import cross_val_score
2 from sklearn.model_selection import LeaveOneOut
3 from sklearn.linear_model import LinearRegression
4
5 lin_reg = LinearRegression()
6 loocv = LeaveOneOut()
7 score = cross_val_score(lin_reg, X, y, cv=loocv)
```

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.model_selection import cross_val_score
3 from sklearn.model_selection import ShuffleSplit
4
5 lin_reg = LinearRegression()
6 shuffle_split = ShuffleSplit(n_splits=5, test_size=0.2, random_state=42)
7 score = cross_val_score(lin_reg, X, y, cv=shuffle_split)
```

It is also called random permutation based cross validation strategy.

- Generates user defined number of train/test splits.

- It is robust to class distribution.

In each iteration, it shuffles order of data samples and then splits it into train and test.

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.model_selection import cross_val_score
3 from sklearn.model_selection import ShuffleSplit
4
5 lin_reg = LinearRegression()
6 shuffle_split = ShuffleSplit(n_splits=5, test_size=0.2, random_state=42)
7 score = cross_val_score(lin_reg, X, y, cv=shuffle_split, scoring='neg_mean_absolute_error')
```

scoring parameter can be set to one of the scoring schemes implemented in sklearn as follows

max_error r2
neg_mean_absolute_error neg_mean_squared_error
neg_mean_squared_log_error neg_median_absolute_error

- For trained estimator, set return_estimator = True

- For scores on training set, set return_train_score = True

```
1 from sklearn.model_selection import cross_val_score
2 from sklearn.model_selection import ShuffleSplit
3
4 cv = ShuffleSplit(n_splits=40, test_size=0.3, random_state=0)
5 cv_results = cross_val_score(
6     regressor, data, target,
7     cv=cv, scoring='neg_mean_absolute_error',
8     return_train_score=True,
9     return_estimator=True)
```

- STEP 1: Instantiate an object of learning_curve class with estimator, training data, size, cross validation strategy and scoring scheme as arguments.

```
1 from sklearn.model_selection import learning_curve
2
3 results = learning_curve(
4     lin_reg, X_train, y_train, train_sizes=train_sizes, cv=cv,
5     scoring='neg_mean_absolute_error')
6 train_size, train_scores, test_scores = results[:3]
7 # Convert the scores into errors
8 train_errors, test_errors = -train_scores, -test_scores
```

- STEP 2: Plot training and test scores as function of the size of training sets. And make assessment about model fitment: under/overfitting or right fit.

- STEP 1: Fit linear models with different number of features.

- STEP 2: For each model, obtain training and test errors.

- STEP 3: Plot #features vs error graph - one each for training and test errors.

- STEP 4: Examine the graphs to detect under/overfitting.

1) Consider the following

- Generalization error
- Empirical error
- Training error
- Test error

Which among the following is correct match?

- A-1
- B-2
- C-2
- D-3

- 4) Which of the following evaluation metric is best suited for targets with exponential growth?

- Mean squared error

- Mean absolute error

- Mean absolute percentage error

- Mean squared log error

2) What is the correct numerical range for the R^2 score (coefficient of determination)?

- [-1, 1]

- (-∞, +∞)

- (-1, 1]

- [0, 1]

3) Which of the following python modules implements useful routines for model evaluation?

- sklearn.model_selection

- sklearn.measurements

- sklearn.metrics

- numpy

3 → demo lectures



MLP Week 4 Notes

L1: Polynomial Regression

How is polynomial regression model trained?

- Step 1 Apply polynomial transformation on the feature matrix.
- Step 2 Learn linear regression model (via normal equation or SGD) on the transformed feature matrix.

Implementation tips: Make use of pipeline construct for polynomial transformation followed by linear regression estimator.

Set up polynomial regression model with normal equation

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import PolynomialFeatures
4
5 # Two steps:
6 # 1. Polynomial features of desired degree (here degree=2)
7 # 2. Linear regression
8 poly_model = Pipeline([
9     ('polynomial_transform', PolynomialFeatures(degree=2)),
10    ('linear_regression', LinearRegression())
11])
12 # Train with feature matrix X_train and label vector y_train
poly_model.fit(X_train, y_train)
```

Set up polynomial regression model with SGD

```
1 from sklearn.linear_model import SGDRegressor
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import PolynomialFeatures
4
5 poly_model = Pipeline([
6     ('polynomial_transform', PolynomialFeatures(degree=2)),
7     ('sgd_regression', SGDRegressor())
8])
poly_model.fit(X_train, y_train)
```

Notice that there is a single line code change in two code snippets.

How to use only interaction features for polynomial regression?

```
1 from sklearn.preprocessing import PolynomialFeatures
2 poly_transform = PolynomialFeatures(degree=2, interaction_only=True)
```

$[x_1, x_2]$ is transformed to $[1, x_1, x_2, x_1 x_2]$.

Note that $[x_1^2, x_2^2]$ are excluded.

$$\begin{aligned} \text{degree} &= 2 \\ (x_1, x_2) \xrightarrow{\text{PF}} &(1, x_1, x_2, x_1 x_2, x_1^2, x_2^2) \Rightarrow 6 \\ \text{interaction only} &\rightarrow \text{True} \\ (x_1, x_2) \xrightarrow{\text{PF}} &(1, x_1, x_2, x_1 x_2) \end{aligned}$$

- 1) Which of the following API will generate a new feature matrix consisting of all polynomial combinations of the features with degree less than or equal to the specified degree?
- PolynomialFeatures
 - TransformFeatures
 - Polynomial_Features
 - Transform_Features

- 2) What will following code implement?

```
from sklearn.linear_model import LinearRegression
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import PolynomialFeatures
poly_model = Pipeline([('poly_transformation', PolynomialFeatures(degree=3)),
                      ('linear_regression', LinearRegression())])
poly_model.fit(X, y)
```

where X and y are the feature matrix and label vector, respectively.

- It will train the linear regression model on the feature matrix X .
- It will only apply the polynomial transformation of degree 3.
- It will apply the polynomial transformation of degree 3 on feature matrix X and then will train the linear regression model.
- It will apply the polynomial transformation of degree 2 on feature matrix X and then will train the linear regression model.

- 3) What will be the transformed feature of $[x_1, x_2]$ after execution of following code?

```
from sklearn.preprocessing import PolynomialFeatures
poly_transform = PolynomialFeatures(degree=2, interaction_only=True)
 $\alpha$  [1, x1, x2, x12, x22]  $\alpha$  x12, x22
 $\alpha$  [1, x12, x22]  $\alpha$ 
[1, x1, x2, x1x2]  $\alpha$ 
 $\alpha$  [1, x1, x2, x12, x22, x1x2]  $\alpha$ 
```

polynomial features
Regression (linear, SGD, Ridge, darts, Trees, !)

How to perform ridge regularization with specific regularization rate?

[Option #1]

Step 1: Instantiate object of Ridge estimator

Step 2: Set parameter alpha to the required regularization rate.

```
1 from sklearn.linear_model import Ridge
2 ridge = Ridge(alpha=1e-3)  $\Rightarrow 10^{-3}$ 
```

fit, score, predict work exactly like other linear regression estimators

[Option #2]

Step 1: Instantiate object of SGDRegressor estimator

Step 2: Set parameter alpha to the required regularization rate and penalty = l2.

```
1 from sklearn.linear_model import SGDRegressor
2 sgd = SGDRegressor(alpha=1e-3, penalty='l2')
```

penalty = 'l2' \Rightarrow Ridge

How to search the best regularization parameter for ridge?

[Option #1]

Search for the best regularization rate with built-in cross validation in RidgeCV estimator.

[Option #2]

Use cross validation with Ridge or SGDRegressor to search for best regularization.

- Grid search
- Randomized search

Set up a pipeline of polynomial transformation followed by the ridge regressor.

```
1 from sklearn.linear_model import Ridge
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import PolynomialFeatures
4
5 poly_model = Pipeline([
6     ('polynomial_transform', PolynomialFeatures(degree=2)),
7     ('ridge', Ridge(alpha=1e-3))])
8 poly_model.fit(X_train, y_train)
```

Instead of Ridge, we can use SGDRegressor, as shown on previous slide, to get equivalent formulation.

How to perform lasso regularization with specific regularization rate?

[Option #1]

Step 1: Instantiate object of Lasso estimator

Step 2: Set parameter alpha to the required regularization rate.

```
1 from sklearn.linear_model import Lasso
2 lasso = Lasso(alpha=1e-3)
```

fit, score, predict work exactly like other linear regression estimators

[Option #2]

Step 1: Instantiate object of SGDRegressor estimator

Step 2: Set parameter alpha to the required regularization rate and penalty = l1.

```
1 from sklearn.linear_model import SGDRegressor
2 sgd = SGDRegressor(alpha=1e-3, penalty='l1')
```

Search for the best regularization rate with built-in cross validation in LassoCV estimator.

[Option #1]

Use cross validation with Lasso or SGDRegressor to search for best regularization.

- Grid search
- Randomized search

Set up a pipeline of polynomial transformation followed by the lasso regressor.

```
1 from sklearn.linear_model import Lasso
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import PolynomialFeatures
4
5 poly_model = Pipeline([
6     ('polynomial_transform', PolynomialFeatures(degree=2)),
7     ('lasso', Lasso(alpha=1e-3))])
8 poly_model.fit(X_train, y_train)
```

Instead of Lasso, we can use SGDRegressor to get equivalent formulation.

L2: Regularization

Lasso
Ridge

How to perform both lasso and ridge regularization in polynomial regression?

Set up a pipeline of polynomial transformation followed by the SGDRegressor with `penalty = 'elasticnet'`

```
1 from sklearn.linear_model import Lasso
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import PolynomialFeatures
4
5 poly_model = Pipeline([
6     ('polynomial_transform', PolynomialFeatures(degree=2)),
7     ('elasticnet', SGDRegressor(penalty='elasticnet',
8         l1_ratio=0.3)))
9 poly_model.fit(X_train, y_train)
```

Remember elasticnet is a convex combination of L1 (Lasso) and L2 (Ridge) regularization.

- In this example, we have set `l1_ratio` to 0.3, which means `l2_ratio` = $1 - l1_ratio$ = 0.7. L2 takes higher weightage in this formulation.

1) Which of the following is correct ways to perform ridge regularization with regularization rate of 0.2?

`from sklearn.linear_model import Ridge`
`ridge = Ridge(alpha=0.2)`

`from sklearn.linear_model import Ridge_regressor`
`ridge = Ridge_regressor(alpha=0.2)`

`sgd = SGDRegressor(alpha=0.2, penalty='l1')`

`from sklearn.linear_model import SGDRegressor`
`sgd = SGDRegressor(alpha=0.2, penalty='l2')`

2) Which penalty is used to perform both lasso and ridge regularization in polynomial regression using `SGDRegressor()`?

- l1
 elasticnet
 None of the above
- 0.7
 3.3
 33.3
 0.5

4) Consider the following code:
`sgdr = SGDRegressor(penalty='elasticnet', l1_ratio=0.3)`
 What is the `l2_ratio` for sgdr?

L3: Hyper Parameter Tuning

- Hyper-parameters are parameters that are not directly learnt within estimators
- In sklearn, they are passed as arguments to the constructor of the estimator classes.

For example,

- degree in `PolynomialFeatures`
- learning rate in `SGDRegressor`

Eg. HPT

n-estimator
 $\text{penalty} = [l_1, l_2]$
 elastic

Select hyperparameters that results in the best cross validation scores.

Hyper parameter search consists of

- an estimator (regressor or classifier);
- a parameter space;
- a method for searching or sampling candidates;
- a cross-validation scheme; and
- a score function.

Two generic HPT approaches implemented in sklearn are:

- `GridSearchCV` exhaustively considers all parameter combinations for specified values.

```
1 param_grid = [
2     {'C': [1, 10, 100, 1000], 'kernel': ['linear']}
```

$4 \times 1 = 4$ combinations

`RandomizedSearchCV` samples a given number of candidate values from a parameter space with a specified distribution.

```
1 param_dist = {
2     "average": [True, False],
3     "l1_ratio": stats.uniform(0, 1),
4     "alpha": loguniform(1e-4, 100)}
```

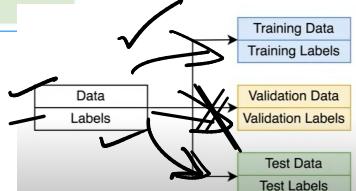
exhaustive

better

Grid search
 Specifies exact values of parameters in grid

Randomized search
 Specifies distributions of parameter values and values are sampled from those distributions. Computational budget can be chosen independent of number of parameters and their possible values.
 The budget is chosen in terms of the number of sampled candidates or the number of training iterations. Specified in `n_iter` argument

Steps in ML

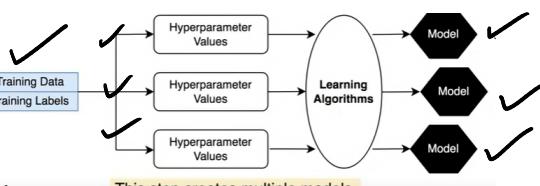


Ridge + Lasso

elastic net

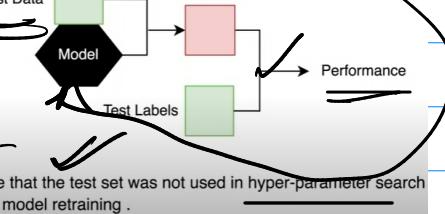
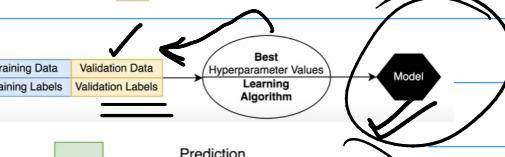
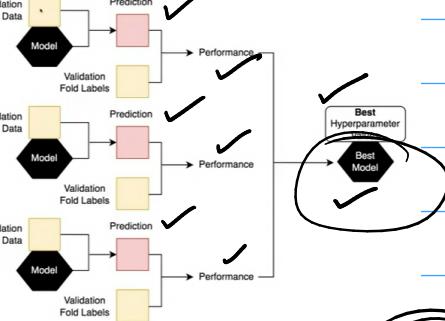
$l_1\text{-ratio} = 0.3$

$\Rightarrow \alpha = 0.2$
 $10^{-3} \Rightarrow \gamma = 10^{-3}$
 $l_2 \Rightarrow 0.7$



This step creates multiple models.

Tips • This step can be run in parallel by setting `n_jobs = -1`.



- Some models can fit data for a range of values of some parameter almost as efficiently as fitting the estimator for a single value of the parameter.
- This feature can be leveraged to perform more efficient cross-validation used for model selection of this parameter.

`linear_model.LassoCV`
 `linear_model.RidgeCV`
 `linear_model.ElasticNetCV`

dim Reg CV
 SGD Reg CV

Cross-val-sq

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.linear_model import SGDRegressor
5
6 param_grid = [
7     {'poly_degree': [2, 3, 4, 5, 6, 7, 8, 9]}
8 ]
9
10 pipeline = Pipeline(steps=[('poly', PolynomialFeatures()),
11                           ('sgd', SGDRegressor())])
12
13 grid_search = GridSearchCV(pipeline, param_grid, cv=5,
14                             scoring='neg_mean_squared_error',
15                             return_train_score=True)
16
17 grid_search.fit(X_train.reshape(-1, 1), y_train)
```

1) Which of the following statements is/are correct about the hyper-parameter?

- Hyper-parameter is a parameter whose value is used to control the learning process.
 Hyper-parameters are parameters that are directly learnt within estimators.
 Hyper-parameters are parameters that are not directly learnt within estimators.
 In linear regression, weight vector is an example of hyper-parameter.

SGD

2) What are the two Hyper-parameter search approaches?

- GridSearchCV
 GridCV
 RandomizedSearchCV
 RandomizedCV

3) Which of the following API is used to find the names and current values for all parameters for a given estimator?

- estimator.get_param()
 estimator.find_param()
 estimator.get_parameter()
 estimator.get_params()

param

- 4) Which of the following options is/are correct about GridSearchCV?
- It exhaustively generates candidates from a grid of parameter values specified with the param_grid parameter.
 - evaluates all possible combination of parameter values when "fitting" on a dataset and retains the best combination.
 - samples a given number of candidates from a parameter space.

- 5) In GridSearchCV, a budget can be chosen independent of the number of parameters and possible values.
- fit_score=0
 - error_score=0
 - error_value=1
 - fit_value=0

5) Some parameter settings may result in a failure to fit one or more folds of the data. Which keyword is used to avoid the same failure?

- fit_score=0
- error_score=0
- error_value=1
- fit_value=0

error_score=0 → theoretical basis

concrete
distribution

4x1

MLP Week 5 Notes

L1: Classification Functions in Sklearn

- In this week we will study sklearn functionality for implementing classification algorithms.

- We will cover sklearn APIs for

- Specific classification algorithms for least square, perceptron and logistic regression classifier.
- ✓ with regularization
- ✓ multiclass, multilabel and multi-output setting
- ✓ Various classification metrics.

- Cross validation and hyper parameter search for classification works exactly like how it works in regression setting.

- However there are a couple of CV strategies that are specific to classification

There are broadly two types of APIs based on their functionality:

- Generic
- SGD classifier

- Specific
- Logistic regression
- Perceptron
- Ridge classifier (for LSC)
- K-nearest neighbours (KNNs)
- Support vector machines (SVMs)
- Naive Bayes

Uses gradient descent for opt
Need to specify loss

Model training

`fit(X, y, coef_init, intercept_init, ...)`

Prediction

`predict(X)` predicts class label for samples

`decision_function(X)` predicts confidence score for samples.

Evaluation

`score(X, y[, sample_weight])`

There are a few common miscellaneous methods as follows:

`get_params([deep])` gets parameter for this estimator.

`set_params(**params)` sets the parameters of this estimator.

`densify()`

→ sparse dataset

↓ dense

Ridge classifier

- RidgeClassifier is a classifier variant of the Ridge regressor.

Binary classification:

- classifier first converts binary targets to {-1, 1} and then treats the problem as a regression task, optimizing the objective of regressor:

- minimize a penalized residual sum of squares

$$\min ||\mathbf{X}\mathbf{w} - \mathbf{y}||_2^2 + \alpha ||\mathbf{w}||_2^2$$

sklearn provides different solvers for this optimization

sklearn uses α to denote regularization rate

predicted class corresponds to the sign of the regressor's prediction

Multiclass classification:

- treated as multi-output regression

predicted class corresponds to the output with the highest value

Step 1: Instantiate a classification estimator without passing any arguments to it. This creates a ridge classifier object.

```
1 from sklearn.linear_model import RidgeClassifier
2 ridge_classifier = RidgeClassifier()
```

Step 2: Call fit method on ridge classifier object with training feature matrix and label vector as arguments.

Note: The model is fitted using X_train and y_train.

```
1 # Model training with feature matrix X_train as input
2 # label vector or matrix y_train
3 ridge_classifier.fit(X_train, y_train)
```

Set alpha to float value. The default value is 0.1.

- alpha should be positive
- Larger alpha values specify stronger regularization.

How to solve optimization problem in RidgeClassifier?

Using one of the following solvers

✓ svd → uses a Singular Value Decomposition of the feature matrix to compute the Ridge coefficients.

✓ cholesky → uses `scipy.linalg.solve` function to obtain the closed-form solution

✓ sparse_cg → uses the conjugate gradient solver of `scipy.sparse.linalg.cg`.

✓ lsqr → uses the dedicated regularized least-squares routine `scipy.sparse.linalg.lsqr` and is fastest.

✓ sag, saga → uses a Stochastic Average Gradient descent iterative procedure. 'saga' is unbiased and more flexible version of 'sag'

✓ newton-cg → uses L-BFGS-B algorithm implemented in `scipy.optimize.minimize`.

✓ lbfgs → can be used only when coefficients are forced to be positive.

→ default

Theo.
theory

optimization procedure

optimization

- For large scale data, use 'sparse_cg' solver.
- When both n_samples and n_features are large, use 'sag' or 'saga' solvers.
 - Note that fast convergence is only guaranteed on features with approximately the same scale.

ridge_classifier = RidgeClassifier(solver='auto')

auto chooses the solver automatically based on the type of data

If data is already centered, set fit_intercept as false, so that no intercept will be used in calculations.

Default: ridge_classifier = RidgeClassifier(fit_intercept=True)

Use predict method to predict class labels for samples

Step 1: Arrange data for prediction in a feature matrix of shape (#samples, #features) or in sparse matrix format.

Step 2: Call predict method on classifier object with feature matrix as an argument.

Predict labels for feature matrix X_test

y_pred = ridge_classifier.predict(X_test)

Perceptron classification → simplest (LR)

- It is a simple classification algorithm suitable for large-scale learning.
- Shares the same underlying implementation with SGDClassifier

Perceptron() → SGDClassifier(loss="perceptron", eta0=1, learning_rate="constant", penalty=None) → similar

Perceptron uses SGD for training.

Step 1: Instantiate a Perceptron estimator without passing any arguments to it to create a classifier object.

from sklearn.linear_model import Perceptron
perceptron_classifier = Perceptron()

Step 2: Call fit method on perceptron estimator object with training feature matrix and label vector as arguments.

Model training with feature matrix X_train and
label vector or matrix y_train
3 perceptron_classifier.fit(X_train, y_train)

Perceptron classifier can be trained in an iterative manner with partial_fit method

False (Default) → perce. run

penalty (default = 'l2') → best

alpha (default = 0.0001) → remains

fit_intercept (default = True) →

n_iter_no_change (default = 5) →

eta0 (default = 1) →

validation_fraction (default = 0.1) →

List of HPTs

Binary classification: → DVA

- This implementation can fit
 - binary classification
 - one-vs-rest (OVR)
 - multinomial logistic regression

Step 1: Instantiate a classifier estimator without passing any arguments to it. This creates a logistic regression object.

from sklearn.linear_model import LogisticRegression
logit_classifier = LogisticRegression()

Step 2: Call fit method on logistic regression classifier object with training feature matrix and label vector as arguments

Model training with feature matrix X_train and
label vector or matrix y_train
3 logit_classifier.fit(X_train, y_train)

For small datasets, 'liblinear' is a good choice, whereas 'sag' and 'saga' are faster for large ones.

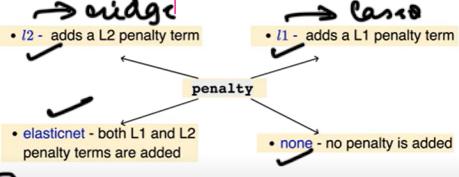
For unscaled datasets, 'liblinear', 'lbfgs' and 'newton-cg' are robust.

For multiclass problems, only 'newton-cg', 'sag', 'saga' and 'lbfgs' handle multinomial loss.

'liblinear' is limited to one-versus-rest schemes

By default, logistic regression uses 'lbfgs' solver.

multi class support X



Regularization is applied by default because it improves numerical stability.

- Not all the solvers support all the penalties.
- Select appropriate solver for the desired penalty.
- L2 penalty is supported by all solvers
- L1 penalty is supported only by a few solvers.

Solver	Penalty
'newton-cg'	['l2', 'none']
'lbfgs'	['l2', 'none']
'liblinear'	✓ ['l1', 'l2']
'sag'	['l2', 'none']
'saga'	✓ ['elasticnet', 'l1', 'l2', 'none']

Note
PYQ

} train

How to control amount of regularization in logistic regression?

- sklearn implementation uses parameter C , which is inverse of regularization rate to control regularization.

Recall

$$\arg \min_{w,C} \text{regularization penalty} + C \text{ cross entropy loss}$$

- C is specified in the constructor and must be positive
- Smaller value leads to stronger regularization.
- Larger value leads to weaker regularization.

LogisticRegression classifier has a class_weight parameter in its constructor.

$$C = \alpha^{-1}$$

$$\alpha \uparrow, \text{reg.} \uparrow$$

$$\alpha \downarrow, \alpha \uparrow, \text{reg.} \uparrow$$

$$0.1 \quad 0.8$$

$$\begin{array}{|c|c|} \hline & \rightarrow \\ \rightarrow & \rightarrow \\ \hline \end{array}$$

concave
convex

What purpose does it serve?

Handles class imbalance with differential class weights.

Mistakes in a class are penalized by the class weight.

- Higher value here would mean higher emphasis on the class.

This parameter is available in classifier estimators in sklearn.

SGDClassifier → general

SGD is a simple yet very efficient approach to fitting linear classifiers under convex loss functions

- This API uses SGD as an optimization technique and can be applied to build a variety of linear classifiers by adjusting the loss parameter.

It supports multi-class classification by combining multiple binary classifiers in a "one versus all" (OVA) scheme.

- Easily scales up to large scale problems with more than 10^5 training examples and 10^3 features. It also works with sparse machine learning problems

Text classification and natural language processing

loss parameter

- 'hinge' - (soft-margin) linear Support Vector Machine
- 'logit' - logistic regression classifier
- 'modified_huber' - smoothed hinge loss brings tolerance to outliers as well as probability estimates
- 'squared_hinge' - like hinge but is quadratically penalized
- 'perceptron' - linear loss used by the perceptron algorithm
- 'squared_error', 'huber', 'epsilon_insensitive', or 'squared_epsilon_insensitive' - regression losses

SGDClassifier implements a plain stochastic gradient descent learning routine.

- the gradient of the loss is estimated with one sample at a time and the model is updated along the way with a decreasing learning rate (or strength) schedule.

Advantages:

- Efficiency.
- Ease of implementation

Disadvantages:

- Requires a number of hyperparameters.
- Sensitive to feature scaling.

It is important

- to permute (shuffle) the training data before fitting the model.
- to standardize the features.

→ ShuffleSplit
→ Standardizing

Step 1: Instantiate a SGDClassifier estimator by setting appropriate loss parameter to define classifier of interest. By default it uses hinge loss, which is used for training linear support vector machine.

```
1 from sklearn.linear_model import SGDClassifier
```

```
2 SGD_classifier = SGDClassifier(loss='log')
```

Here we have used 'log' loss that defines a logistic regression classifier.

Step 2: Call fit method on SGD classifier object with training feature matrix and label vector as arguments.

```
1 # Model training with feature matrix X_train and
2 # label vector or matrix y_train
3 SGD_classifier.fit(X_train, y_train)
```

- Least square classification (RidgeClassifier)
- Perceptron (Perceptron)
- Logistic regression (LogisticRegression)

Alternatively we can use SGDClassifier with appropriate loss setting for implementing these classifiers:

- loss = 'log' for logistic regression
- loss = 'perceptron' for perceptron
- loss = 'squared_error' for least square classification

loss = hinge
(default)

similar

Classification estimators implements a few common methods like fit, score, decision_function, and predict.

1) RidgeClassifier is a classifier variant of which regressor?

- Ridge regressor
- Lasso regressor
- NoSQL regressor
- None of the above

2) The first step to apply Ridge classifier should be-

- Call fit method on ridge classifier object with training feature matrix and label vector as arguments
- Instantiate a classification estimator without passing any arguments
- Set alpha to float value
- Use solvers to find the parameter weights that minimize the cost function

3) Which of the following solvers uses the `scipy.linalg.solve` function to obtain the closed-form solution for minimizing the cost function?

- svd
- cholesky
- lsqr
- lbfgs

4) For large scale data, which solver can be used in RidgeClassifier?

- sparse_cg
- n_samples

5) To choose the solver automatically based on the type of data, we can use-

- ridge_classifier = RidgeClassifier(solver=sparse_cg)
- ridge_classifier = RidgeClassifier(solver=auto)
- ridge_classifier = RidgeClassifier(solver=svd)
- ridge_classifier = RidgeClassifier(solver=cholesky)

6) How do we make predictions on data samples?

- Use predict method to predict class labels for samples
- Use densify method to predict class labels for samples
- Use sparsify method to predict class labels for samples
- None of the above

7) Mention TRUE or FALSE : LogisticRegression is a linear model for classification rather than regression.

- TRUE
- FALSE

8) Which of the following codes can be used to set regularization using penalty in Logistic Regression classifier?

- logit_classifier = LogisticRegression(penalty='l2')
- logit_classifier = LogisticClassifier(solver='auto')
- logit_classifier = LogisticClassifier(solver='svd')
- logit_classifier = LogisticRegression(penalty=0.1)

9) Which of the following solvers is used for a Logistic Regression classifier applying on a small dataset?

- liblinear
- sag
- saga
- None of the above

10) Mention TRUE or FALSE : In LogisticRegression classifier, all the solvers support all the penalties.

- TRUE
- FALSE

11) Which of the following statements are true about SGD?

- It is an efficient approach to fitting linear classifiers under convex loss functions
- It is an optimization technique
- It corresponds to a specific family of machine learning models
- None of the above

12) What does a SGDClassifier implement?

- a plain stochastic gradient descent learning routine.

13) The disadvantages of SGDClassifier include:

- Efficiency
- Ease of implementation
- Requires a number of hyperparameters
- Sensitive to feature scaling

14) Which loss functions are supported in SGDClassifier?

- hinge
- modified huber
- squared hinge
- perceptron

15) How to specify maximum number of epochs for SGDClassifier?

- epoch()
- max_iter
- max_iter()
- None of the above

L2: Multi Learning Classification

Basics of multiclass, multilabel and multioutput classification

1 output
no. of labels > 2

> 1 output

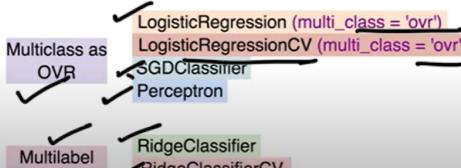
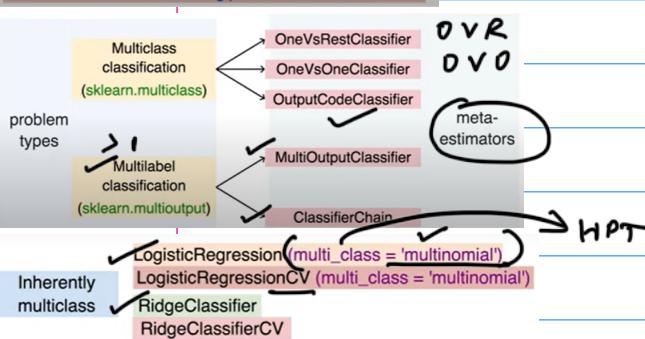
Multiclass classification has exactly one output label and the total number of labels > 2.

- For more than one output, there are two types of classification models:



We will refer both these models as multi-label classification models, where # of output labels > 1.

Multiclass, multilabel, multioutput problems are referred to as multi-learning problems.



Classification task with more than two classes.
Each example is labeled with exactly one class

In Iris dataset, → > 2 label
• There are three class labels: setosa, versicolor and virginica.
• Each example has exactly one label of the three available class labels.
• Thus, this is an instance of a multi-class classification.

In MNIST digit recognition dataset,
• There are 10 class labels: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9.
• Each example has exactly one label of the 10 available class labels.
• Thus, this is an instance of a multi-class classification.

Each example is marked with a single label out of k labels. The shape of label vector is (n, 1).

Use LabelBinarizer transformation to convert the class label to multi-class format.

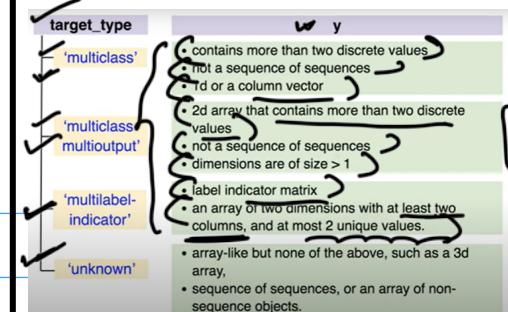
```
1 from sklearn.preprocessing import LabelBinarizer
2 y = np.array(['apple', 'pear', 'apple', 'orange'])
3 y_dense = LabelBinarizer().fit_transform(y)
```

The resulting label vector has shape of (n, k).

Use type_of_target to determine the type of the label.

```
1 from sklearn.utils.multiclass import type_of_target
2 type_of_target(y)
```

In case, y is a vector with more than two discrete values, type_of_target returns multiclass.



Apart from these, there are three more types, type_of_target can determine targets corresponding to regression and binary classification.

- continuous - regression target
- continuous-multioutput - multi-output target

All classifiers in scikit-learn perform multiclass classification out-of-the-box.

Use sklearn.multiclass module only when you want to experiment with different multiclass strategies.

- Using different multi-class strategy than the one implemented by default may affect performance of classifier in terms of either generalization error or computational resource requirement.

OVR - OneVsRestClassifier

- Fits one classifier per class c - c vs not c.
- This approach is computationally efficient and requires only k classifiers.
- The resulting model is interpretable.

```
1 from sklearn.multiclass import OneVsRestClassifier
2 OneVsRestClassifier(LinearSVC(random_state=0)).fit(X, y)
```

- We need to supply estimator as an argument in the constructor.
- Support methods like other classifiers - fit, predict, predict_proba, partial_fit.

OneVsRest classifier also supports multilabel classification.

OVA - OneVsOneClassifier

- Fits one classifier per pair of classes. Total classifiers = $\binom{k}{2}$.
- Predicts class that receives maximum votes.
- The tie among classes is broken by selecting the class with the highest aggregate classification confidence.

```
1 from sklearn.multiclass import OneVsOneClassifier
2 OneVsOneClassifier(LinearSVC(random_state=0)).fit(X, y)
```

- We need to supply estimator as an argument in the constructor.
- Support methods like other classifiers - fit, predict, predict_proba, partial_fit.

OneVsOne classifier processes subset of data at a time and is useful in cases where the classifier does not scale well.

OneVsRestClassifier

- Fits one classifier per class.
- For each classifier, the class is fitted against all the other classes.

OneVsOneClassifier

- Fits one classifier per pair of classes.
- At prediction time, the class which received the most votes is selected.

MultiOutputClassifier

- Able to estimate a series of target functions that are trained on a single predictor matrix to predict a series of responses.

- Allows multiple target variable classifications.

- For a multi-label classification problem with k classes, k binary classifiers are assigned an integer between 0 and k - 1.
- These integers define the order of models in the chain.

1) Mention TRUE or FALSE : Perceptron classification shares the same underlying implementation with SGDClassifier

- TRUE

- FALSE

2) Mention TRUE or FALSE : Partial_fit method is a way to plot the iteration vs. loss curve

- TRUE

- FALSE

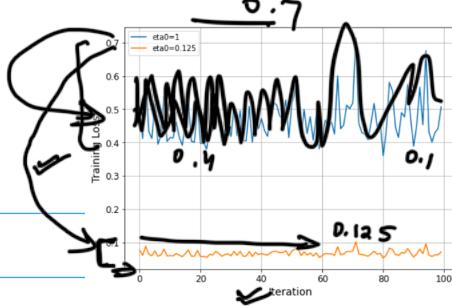


Figure 1: Iteration vs. Training Loss

3) The general trend of variation in training loss with number of iterations-

- Training losses increase exponentially with the number of iterations.
- For $\text{eta}_0 = 1$ and $\text{eta}_0 = 0.125$, the mean of the training loss remains within a particular range of training loss with change in iteration
- For $\text{eta}_0 = 1$, the training loss decreases with number of iteration and for $\text{eta}_0 = 0.125$, the training loss increases with number of iteration.
- For $\text{eta}_0 = 1$, the training loss increases with number of iteration and for $\text{eta}_0 = 0.125$, the training loss decreases with number of iteration.

4) The mean value of training loss for $\text{eta}_0 = 0.125$ is in the range-

- is less than 0.1
- 0.4 to 0.5
- mean is equal to 1
- is greater than or equal to 0.3

5) For what value of eta_0 , the graph of training loss is noisy?

- 0.125
- TRUE
- FALSE

7) Single label multiclass classification have how many output labels?

- 2
- 1
- 3
- 4

8) Mention TRUE or FALSE: OVR classifier supports multilabel classification.

- TRUE
- FALSE

So far we learnt how to train classifiers for binary, multi-class and multi-label/output cases.

We will learn how to evaluate these classifiers with different scoring functions and with cross-validation.

We will also study how to set hyper-parameters for classifiers.

Many cross-validation and HPT methods discussed in the regression context are also applicable in classifiers.

- We will not repeat that discussion in this topic.
- Instead we will focus on only additional methods that are specific to classifiers.

Stratified cross validation iterators

There may be issues like class imbalance in classification, which tend to impact the cross validation folds.

The overall class distribution and the ones in folds may be different and this has implications in effective model training.

sklearn.model_selection module provides three stratified APIs to create folds such that the overall class distribution is replicated in individual folds.

sklearn.model_selection module provides the following three stratified APIs to create folds such that the overall class distribution is replicated in individual folds.

- StratifiedKFold
- RepeatedStratifiedKFold
- StratifiedShuffleSplit

LogisticRegressionCV
Support in-build cross validation for optimizing hyperparameters

- The following are key parameters for HPT and cross validation

$cv = 5$

cv specifies cross validation iterator

$scoring$ specifies scoring function to use for HPT

Cs specifies regularization strengths to experiment with

- Choosing the best hyper-parameters

Scores averaged across folds, values corresponding to the best score are selected and final refit with these parameters

$refit = True$

$refit = False$

automatically fit

sklearn.metrics implements a bunch of classification scoring metrics based on true labels and predicted labels as inputs.

- accuracy_score
- balanced_accuracy_score
- top_k_accuracy_score
- roc_auc_score
- precision_score
- recall_score
- f1_score

`confusion_matrix` evaluates classification accuracy by computing the confusion matrix with each row corresponding to the true class.

Std.

```
from sklearn.metrics import confusion_matrix
```

TP

FP

FN

TN

Example:

array([[2, 0, 0], [0, 0, 1], [1, 0, 2]])

Entry i, j in a confusion matrix

number of observations actually in group i , but predicted to be in group j .

- Confusion matrix

```
ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=clf.classes_)
```

- From estimators

```
ConfusionMatrixDisplay.from_estimator(clf, X_test, y_test)
```

- From predictions

```
ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
```

The `classification_report` function builds a text report showing the main classification metrics.

```
from sklearn.metrics import classification_report
print(classification_report(y_true, y_predicted))
```

```
from sklearn.metrics import precision_recall_curve
precision, recall, thresholds = precision_recall_curve(y_true, y_predicted)
```

```
from sklearn.metrics import roc_curve
fpr, tpr, thresholds = metrics.roc_curve(y_true, y_scores, pos_label=2)
```

Treat data as a collection of binary problems, one for each class.

- Then, average binary metric calculations across the set of classes.

Can be done using `average` parameter.

`macro` calculates the mean of the binary metrics

`weighted` computes the average of binary metrics in which each class's score is weighted by its presence in the true data sample.

`micro` gives each sample-class pair an equal contribution to the overall metric

`samples` calculates the metric over the true and predicted classes for each sample in the evaluation data, and returns their average

`None` returns an array with the score for each class

1) Which stratified APIs are provided by `sklearn.model_selection` module to create folds?

- StratifiedKFold
- RepeatedStratifiedKFold
- StratifiedShuffleSplit

`OVR` → One vs Rest

2) Mention TRUE or FALSE: confusion matrix evaluates classification accuracy by computing the confusion matrix with each row corresponding to the true class.

- TRUE
- FALSE



3) Confusion matrix can be displayed with which API in `sklearn.metrics`?

- ConfusionMatrixDisplay
- RepeatedStratifiedKFold
- ConfusionDisplay
- OVR

Naive Bayes

4) Which NB is used to implement the Gaussian Naive Bayes algorithm for classification?

- GaussianNB
- MultinomialNB
- ComplementNB of these

L1: Naive Bayes Classifier

Naive Bayes classifier

- Naive Bayes classifier applies Bayes' theorem with the "naive" assumption of conditional independence between every pair of features given the value of the class variable.

For a given class variable y and dependent feature vector x_1 through x_m ,

the naive conditional independence assumption is given by:

$$P(x_i|y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_m) = P(x_i|y)$$

Naive Bayes learners and classifiers can be extremely fast compared to more sophisticated methods.

Dep. Reg.
Perceptron
SGD Classifier
Ridge Classifier
Bayes Th.

List of NB Classifiers

- Implemented in `sklearn.naive_bayes` module

GaussianNB	$P, 1 - P$
BernoulliNB	CategoricalNB
MultinomialNB	ComplementNB
Tent	→ data is imbalanced
Implements <code>fit</code> method to estimate parameters of NB classifier with <u>feature matrix</u> and <u>labels</u> as inputs.	
The prediction is performed using <code>predict</code> method.	

1) Which of the following is the correct assumption about the naive Bayes classifier?

- $P(x_i|x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_m) = P(x_i)$
- $P(x_1, x_2, \dots, x_m) = P(x_i)$
- $P(y|x_1, x_2, \dots, x_m) = P(y)$

$P(x_i|y, x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_m) = P(x_i|y)$

2) Which of the following modules of `sklearn` contains the API `GaussianNB`?

- `sklearn.linear_regression`
- `sklearn.model_selection`
- `sklearn.naive_bayes`
- None of the above

3) Consider following code snippet and select the correct option.

```
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
gnb.fit(X, y)
```

where X and y are training data.

- gnb is best suited when features X are imbalanced.
- gnb is best suited when features X are categorical.
- gnb is best suited when X does not have categorical features.
- gnb is best suited for both categorical and numerical features.

4) Which of the following is best suited when the data is imbalanced?

- ComplementNB
- GaussianNB
- BernoulliNB
- None of the above

5) Which of the following Naive Bayes classifiers is suitable for classification of word counts for text classification?

- GaussianNB
- CategoricalNB
- BernoulliNB
- MultinomialNB

feature matrix
does not have
any categorical
data!

$p(1-p)$

MLP Week 7 Notes

L2. K Nearest Neighbours

Nearest neighbor classifier

It is a type of instance-based learning or non-generalizing learning.
 Does not attempt to construct a model.
 Simply stores instances of the training data.

Classification is computed from a simple majority vote of the nearest neighbors of each point.

- Two different implementations of nearest neighbors classifiers are available.

1. KNeighborsClassifier

2. RadiusNeighborsClassifier

How are KNeighborsClassifier and RadiusNeighborsClassifier different?

KNeighborsClassifier

learning based on the k nearest neighbors

most commonly used technique

choice of the value k is highly data-dependent

RadiusNeighborsClassifier

learning based on the number of neighbors within a fixed radius r of each training point

used in cases where the data is not uniformly sampled

fixed value of r is specified, such that points in sparser neighborhoods use fewer nearest neighbors for the classification

Step 1: Instantiate a KNeighborsClassifier estimator without passing any arguments to it to create a classifier object.

```
1 from sklearn.neighbors import KNeighborsClassifier
2 knn_classifier = KNeighborsClassifier()
```

Step 2: Call fit method on KNeighbors classifier object with training feature matrix and label vector as arguments.

```
1 # Model training with feature matrix X_train and
2 # label vector or matrix y_train
3 knn_classifier.fit(X_train, y_train)
```

- Specify the number of nearest neighbors **K** from the training dataset using n_neighbors parameter.

value should be int.

```
1 knn_classifier = KNeighborsClassifier(n_neighbors=5)
```

- It is better to weight the neighbors such that nearer neighbors contribute more to the fit.

weights

- 'uniform': All points in each neighborhood are weighted equally.
- 'distance': weight points by the inverse of their distance.
- closer neighbors of a query point will have a greater influence than neighbors which are further away.

algorithm

- 'ball_tree' will use BallTree
- 'kd_tree' will use KDTree
- 'brute' will use a brute-force search
- 'auto' will attempt to decide the most appropriate algorithm based on the values passed to the fit method.

For 'ball_tree' and 'kd_tree' algorithms, there are some other parameters to be set.

leaf_size

- can affect the speed of the construction and query, as well as the memory required to store the tree
- default = 30

metric

- Distance metric to use for the tree
- It is either string or callable function
- some metrics are listed below:
 - 'euclidean', 'manhattan', 'chebyshev',
 - 'minkowski', 'wminkowski',
 - 'seuclidean', 'mahalanobis'
- default = 'minkowski'

```
1 from sklearn.neighbors import RadiusNeighborsClassifier
2 radius_classifier = RadiusNeighborsClassifier()
```

The number of neighbors is specified within a fixed radius **r** of each training point using radius parameter.

r is a float value.

```
1 radius_classifier = RadiusNeighborsClassifier(radius=1.0)
```

weights

- 'uniform'
- 'distance'
- [callable] function
- default = 'uniform'

algorithm

- 'ball_tree'
- 'kd_tree'
- 'brute'
- 'auto'
- default = 'auto'

1) Which of the following sklearn modules implement relevant functions for the Nearest Neighbors algorithm?

- sklearn.neighbors
- sklearn.mixture
- sklearn.model_selection
- sklearn.ensemble

2) Based on the following statements, select the correct option.

For every dataset, there exists a threshold N such that increasing the value of N beyond this threshold does not significantly affect the accuracy of the model.

If statement I is true, statement II is true. For every dataset, there exists a threshold N such that increasing the value of N beyond this threshold does not significantly affect the accuracy of the model.

Statement I is true, Statement II is false.

Statement I is false, Statement II is true.

Statement I is false, Statement II is true.

3) What does the value k in k-Nearest Neighbors algorithm refer to?

- Number of data points to be classified.
- Number of nearest neighbors to consider for each test data point.
- Number of dimensions of each data point.
- Maximum distance between a test data point and its nearest neighbor.

4) Which of the following statements are true? (Multiple options may be correct.)

- KNN models with low values of k produce complex decision boundaries.
- KNN models with high values of k produce complex decision boundaries.
- KNN models with low values of k produce simple decision boundaries.
- KNN models with high values of k produce simple decision boundaries.

↑ complexity

✓ Based on the following statements, select the correct option.

- A given test data point in KNN can be thought of as being situated at the centroid of k training points nearest to it.
 - If x_1, x_2, \dots, x_k are the k nearest data points of x , then $\hat{y}(x) = \frac{1}{k} \sum_{i=1}^k y(x_i)$.
- Statement I is true, Statement II is true, Statement III is a correct explanation for statement I.
 - Statement I is true, Statement II is true, Statement III is not a correct explanation for statement I.
 - Statement I is true, Statement II is false.
 - Statement I is false, Statement II is true.
- Which of the following are false regarding KNN? (Multiple options may be correct.)
- KNN is a non-parametric supervised learning algorithm.
 - The number k in KNN is a parameter to be estimated from the training data.
 - The time complexity of classifying a test data point using KNN depends only upon the number of training instances, and not their dimension.
 - KNN models suffer from the curse of dimensionality.

L6: SVM in Scikit-Learn

Support Vector Machines

Support Vector Machines (SVM) are a set of supervised learning methods used for classification, regression and outliers detection.

SVM constructs a hyper-plane or set of hyper-planes in a high or infinite dimensional space, which can be used for classification, regression or other tasks.

In sklearn, we have three methods to implement SVM.

SVC

These are similar methods but accept slightly different sets of parameters. Implementation is based on libsvm.

NuSVC

Faster implementation of linear SVM classification with only linear kernel. Implementation is based on liblinear.

LinearSVC

fast

Array X : holding the training samples

```
x = [(0, 0), (1, 1)]
```

shape → (n_samples, n_features)

Array y : holding the class labels (strings or integers)

```
y = [0, 1]
```

Step 1: Instantiate a SVC classifier estimator.

```
1 from sklearn.svm import SVC
2 SVC_classifier = SVC()
```

Step 2: Call fit method on SVC classifier object with training feature matrix and label vector as arguments.

```
1 # Model training with feature matrix X_train and
2 # label vector or matrix y_train
3 SVC_classifier.fit(X_train, y_train)
```

c ↓ Regularization parameter
float value

Default: 1 SVC_classifier = SVC(C=1.0)

c ↑ reg.↑

c ↑ reg.↓

Note:
 • strength of the regularization is inversely proportional to C
 • strictly positive
 • penalty is a squared L2 penalty

kernel

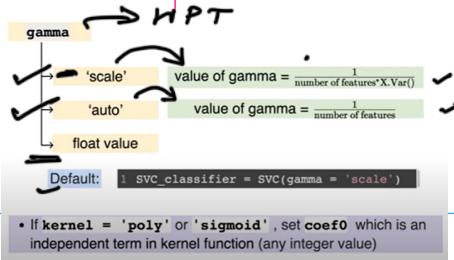
- 'linear'
- 'poly'
- 'rbf'
- 'sigmoid'
- 'precomputed'

thea.

Default: 1 SVC_classifier = SVC(kernel = 'rbf')

→ polynomial features

- If **kernel = poly**, set **degree** (any integer value)
- If **kernel = callable** is given it is used to pre-compute the kernel matrix from data matrices



After the classifier is fit on the training data, there are few attributes which reveal the details of support vectors.

```
1 from sklearn.svm import SVC
2 SVC_classifier = SVC()
3 clf = SVC_classifier.fit(X_train, y_train)
4
5 # to view indices of the support vectors
6 clf.support_
7
8 # to view the support vectors
9 clf.support_vectors_
10
11 # to view the number of support vectors for each class
12 clf.n_support_
```

Step 1: Instantiate a NuSVC classifier estimator.

```
1 from sklearn.svm import NuSVC
```

Instead of C in SVC, nu is introduced in NuSVC to control the number of support vectors and margin errors.

nu is an upper bound on the fraction of margin errors and a lower bound of the fraction of support vectors.

Value of nu should $\in (0, 1]$

Default: nu = 0.5

Other parameters for NuSVC are same as that of SVC.

Step 1: Instantiate a LinearSVC classifier estimator.

```
1 from sklearn.svm import LinearSVC
2 LinearSVC_classifier = LinearSVC()
```

Step 2: Call fit method on SVC classifier object with training feature matrix and label vector as arguments.

```
1 # Model training with feature matrix X_train and
2 # label vector or matrix y_train
3 LinearSVC_classifier.fit(X_train, y_train)
```

penalty

- i1 - adds a L1 penalty term
- i2 - adds a L2 penalty term

i1 leads to coef_ vectors that are sparse.

```
Default: i LinearSVC_classifier = Linear_SVC(penalty = 'l2')
```

loss parameter

- 'hinge' - standard SVM loss
- 'squared_hinge' - square of the hinge loss

Default:

```
1 LinearSVC_classifier = Linear_SVC(loss = 'squared_hinge')
```

c

Regularization parameter

dual

better results

- Select the algorithm to either solve the dual or primal optimization problem.
- When n_samples > n_features, prefer dual=False.

fit_intercept

To calculate the intercept for the model.

- SVC and NuSVC implement the "one-versus-one" approach for multi-class classification.

decision_function_shape

- 'ovo' → one vs one
- 'ovr' → one vs rest

- LinearSVC implements "one-vs-the-rest" approach for multi-class classification.

multi_class

- 'ovr'
- 'crammer_singer'

(out of syllabus)

Advantages of SVM

- Effective in high dimensional spaces.
- Effective in cases where number of dimensions is greater than the number of samples.
- Uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.
- Versatile: different Kernel Functions can be specified for the decision function.

n features > n samples

Disadvantages of SVM

- SVMs do not directly provide probability estimates, these are calculated using an expensive five-fold cross-validation.
- Avoid over-fitting in choosing Kernel functions if the number of features is much greater than the number of samples.

- Which of the following code will throw an error?
 - from sklearn import svm
 - from sklearn import SVC
 - from sklearn import SVR
 - from sklearn.svm import SVC
 - from sklearn import svm.SVC
- Which keyword parameter is used to set the Regularization parameter in svm.SVC?
 - R
 - H
 - RC
 - C
- Which of the following module of sklearn contains the API confusion matrix?
 - sklearn.svm
 - sklearn.datasets
 - sklearn.model_selection
 - sklearn.metrics
- Consider the following code for the classification:


```
smv_clf = Pipeline([
    ('scaler', StandardScaler()),
    ('linear_svc', LinearSVC(C=100, loss='hinge')),
])
```

If this leads to the overfitting, which of the following values may overcome the problem of overfitting?

 - C = 110
 - C = 10
 - C = 150
 - C = 1
- SVM classifiers do not output probabilities for each class.
 - True
 - False
- Consider the following code:


```
smv_clf = Pipeline([
    ('scaler', StandardScaler()),
    ('linear_svc', LinearSVC(C=1, loss='hinge')),
])
```

Which of the following code may be an alternate way to the same code?

 - SVC(C=1)
 - SVC(kernel='linear', C=1)
 - SVC(kernel='linear', C=100)
 - SVC(kernel='rbf', C=1)
- What will the following code implement? (Assume that all the necessary libraries are imported)


```
polynomial_svm_clf = Pipeline([
    ('poly_features', PolynomialFeatures(degree=3)),
    ('scaler', StandardScaler()),
    ('smv_clf', LinearSVC(C=10, loss='hinge'))
])
```

 - It will fit the training data using support vector machine.
 - It will transform the training feature to degree 3 and then will fit the data using SVM.
 - It will normalize the training data using StandardScalar and then will fit the data using SVM.
 - It will transform the training data to degree 3, normalize the transformed data using StandardScalar and then will fit the data using SVM.
- Consider the following code: (Assume that all the necessary libraries are imported)


```
clf = make_pipeline(StandardScaler(), SVC(gamma='auto', kernel='poly'))
```

What will be the degree of the polynomial kernel function?

3
- What is the default value of regularization parameter used in SVC classifier?
 - 1.0
 - 0.1
 - 0.5
 - 0.2
- For a support vector machine model, let z_i be an input instance with label y_i . If $y_i(\mathbf{w}_0 + \mathbf{x}_i^T \mathbf{W}) > 1$ where \mathbf{w}_0 and \mathbf{W} are the estimated parameters of the model, then
 - z_i is a support vector.
 - z_i is not a support vector.
 - z_i is an outlier.
 - Depending upon other data points, z_i may or may not be a support vector.
- Which of the following methods is used to get the indices of support vectors?
 - .support_vectors_
 - .support_
 - .n_support_
 - .nsupport_
- Which of the following methods is used to get the number of support vectors?
 - .support_vectors_
 - .support_
 - .n_support_
 - .nsupport_
- Consider the following code:

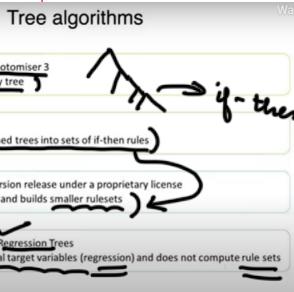

```
from sklearn import svm
clf = svm.NuSVC(gamma='auto')
clf.fit(X, Y)
```

where X and Y are the feature matrix and label vector, respectively. Which kernel is used in the above classification?

 - linear
 - polynomial of degree 3
 - rbf
 - sigmoid
- Which keyword parameter is used in NuSVC to set an upper bound on the fraction of margin errors and a lower bound of the fraction of support vectors?
 - gamma
 - sv
 - nu

MLP Week 8 Notes *Ensemble Methods*

L1: Decision Trees



✓ **scikit-learn** uses an optimized version of the **CART** algorithm; however, it does not support categorical variables for now.

Classification: `sklearn.tree.DecisionTreeClassifier`
Regression: `sklearn.tree.DecisionTreeRegressor`

Both these estimators have the same set of parameters except for `criterion` used for tree splitting.

`splitter`: max_depth, min_samples_split, min_samples_leaf

`splitter`: Strategy for splitting at each node. best, random

`max_depth`: Maximum depth of the tree.

`int`: When `None`, the tree is expanded until all leaves are pure or they contain less than `min_samples_split` samples.

`min_samples_split`: The minimum number of samples required to split an internal node.

`int`, `float`: The minimum number of samples required to be at a leaf node.

`min_samples_leaf`: The minimum number of samples required to be at a leaf node.

`criterion`: Specifies function to measure the quality of a split.

Classification

✓ `gini`
✓ `entropy`

Regression

✓ `squared_error`
✓ `friedman_mse`
✓ `absolute_error`
✓ `poisson`

Visualize a tree →

✓ `decision_tree`: The decision tree to be plotted.
✓ `max_depth`: The maximum depth of the representation. If `None`, the tree is fully generated.
✓ `feature_names`: Names of each of the features.
✓ `class_names`: Names of each of the target classes in ascending numerical order.
✓ `label`: Whether to show informative labels for impurity.

Avoiding overfitting of trees →
Pre-pruning ↗ `tree.fit` for finding the best set of parameters.

Post-pruning
First grows trees without any constraints and then uses `cost_complexity_pruning` with `max_depth` and `min_samples_split`.

Use `min_samples_split` or `min_samples_leaf` to ensure that multiple samples influence every decision in the tree, by controlling which splits will be considered.
A very small number will usually mean the tree will overfit.
A large number will prevent the tree from learning the data.

HPT



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

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

L5: Voting, Bagging and Random Forest

Voting estimators

Class: `sklearn.ensemble.VotingClassifier`

Class: `sklearn.ensemble.VotingRegressor`

Both these estimators take the following common parameters:

✓ `base_estimator`

✓ `weights`

Both these estimators implement the following functions:

✓ `fit`

✓ `predict`

✓ `fit_transform`

✓ `score`

* `VotingClassifier` takes an additional argument:

HPT ↗ voting

hard ↗ hard (default)

Class: `sklearn.ensemble.BaggingClassifier`

Class: `sklearn.ensemble.BaggingRegressor`

Class: `sklearn.ensemble.RandomForestClassifier`

Class: `sklearn.ensemble.RandomForestRegressor`

The parameters can be classified as:

✓ Decision tree parameters

✓ Bagging parameters

• The number of trees are specified by `n_estimators`. ↗ HPT

✓ Default #trees for classification = 10

✓ Default #trees for regression = 100

(bootstrap) specifies whether to use bootstrap samples for training. ↗ boolean

✓ True : bootstrapped samples are used.

✓ False : whole dataset is used.

`oob_score` specifies whether to use out-of-bag samples for estimating generalization error. It is only available when `bootstrap = True`. ↗ test

• `max_samples` specifies the number of samples to be drawn while bootstrapping.

✓ `None` : Use all samples in the training data.

✓ `int` : Use `max_samples` samples from the training data.

✓ `float` : Use `max_samples` * total number of samples from the training data. The value should be between 0 and 1. ↗ 0.6 ↗ 60 %

• `random_state` controls randomness of features and samples selected during bootstrap. ↗ steady

• The number of features to be considered while splitting is specified by `max_features`.

✓ auto, sqrt, log2, int, float

✓ tune ↗

• The depth of the tree is controlled by `max_depth`. The default value is `None`, which means the tree will be grown until all leaf nodes are pure or until leaves contain less than `min_samples_splits` samples.

• estimators_ member variable contains a collection of fitted estimators.

func ↗

• `feature_importances_` member variable contains a list of important features.

✓ `fit` builds forest of trees from the training dataset with the specified parameters.

✓ `decision_path` returns decision path in the forest.

✓ `predict` returns class label in classification and output value in regression.

✓ ↗ RF etc. ↗ extra complex models ↗

✓ `predict_proba` and `predict_log_proba` returns probabilities and their logs for classification set up.

complex models

L8: Boosting: Adaboost, GradientBoost, XGBoost

Class: `sklearn.ensemble.AdaBoostClassifier`

XGBoost

↓
boost
models

base_estimator_ ↴
base-estimator

base_estimator_ ↴
• Default estimator is `DecisionTreeClassifier` with depth = 1.

n_estimators_ ↴
Maximum number of estimators where boosting is terminated. The default value is 50.

learning_rate_ ↴
Weight applied to each classifier during boosting.
Higher value here would increase contribution of individual classifiers.
• There is a trade-off between `n_estimators` and `learning_rate`.

base_estimator_ ↴
Base estimator of ensemble.

estimators_ ↴
Collection of fitted sub-estimators.

estimator_weights_ ↴
Weights for each estimator in ensemble.

estimator_errors_ ↴
Errors for each estimator in ensemble.

n-estimator × 1

learning
-rate

Estimator
func ↴

There are two most important parameters of these estimators:

- `n_estimators`
- `learning_rate`

sklearn.ensemble.GradientBoostingClassifier supports both binary and multiclass classification.

MLP Week 9 Notes

L3: Neural Networks: MLP

Multilayer Perceptron (MLP)

- It is a supervised learning algorithm.
- MLP learns a non-linear function approximator for either classification or regression depending on the given dataset.
- In sklearn, we implement MLP using:
 - MLPClassifier for classification
 - MLPRegressor for regression

MLPClassifier supports multi-class classification by applying Softmax as the output function.

It also supports multi-label classification in which a sample can belong to more than one class.

MLPRegressor also supports multi-output regression, in which a sample can have more than one target.

Step 1: Instantiate a MLP classifier estimator.

```
1 from sklearn.neural_network import MLPClassifier
2 MLP_clf = MLPClassifier()
```

Step 2: Call fit method on MLP classifier object with training feature matrix and label vector as arguments.

Step 3: After fitting (training), the model can make predictions for new samples (X_{test}) using two methods:

```
1 MLP_clf.predict(X_test)
2 MLP_clf.predict_proba(X_test)
```

- gives labels for new samples
- for example: array([1, 0])

- gives vector of probability estimates per sample
- for example: array([1.967...e-04, 9.998...e-01])

hidden_layer_sizes

- The parameter sets the number of layers and the number of neurons in each layer.

It is a tuple where each element in the tuple represents the number of neurons at the i th position where i is the index of the tuple.

- The length of tuple denotes the total number of hidden layers in the network.

To create a 3 hidden layer neural network with 15 neurons in first layer, 10 neurons in second layer and 5 neurons in third layer:

```
1 MLPClassifier(hidden_layer_sizes=(15,10,5)) = (15,10,5)
```

no-op activation returns $f(x) = x$

logistic sigmoid function returns $f(x) = \frac{1}{(1+exp(-x))}$

'identity'

'logistic'

'tanh'

'relu'

Default

hyperbolic tan function returns $f(x) = \tanh(x)$

rectified linear unit function returns $f(x) = max(0,x)$

- MLPClassifier optimizes the log-loss function using LBFGS or stochastic gradient descent

solver: lbfgs, sgd, adam

If the solver is 'lbfgs', the classifier will not use minibatch.

- Size of minibatches can be set to other stochastic optimizers: batch_size (int)

- default batch_size is 'auto'.

```
1 batch_size=min(200, n_samples)
```

intercepts_ \rightarrow estimator of MLP

- It is a list of shape $(n_layers - 1)$
- The i th element in the list bias vector corresponding to layer $i + 1$.

Example:

- "Bias values for first hidden layer:"

```
1 print(MLP_clf.intercepts_[0])
```

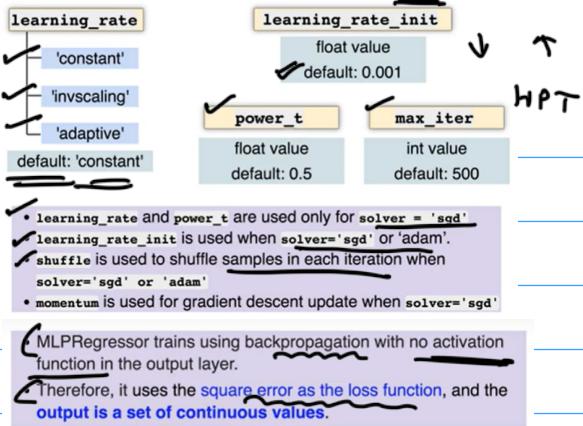
- "Bias values for second hidden layer:"

```
1 print(MLP_clf.intercepts_[1])
```

third

[2]

MLP Week 9 Notes



The parameters of MLPRegressor are the same as that of MLPClassifier.

