

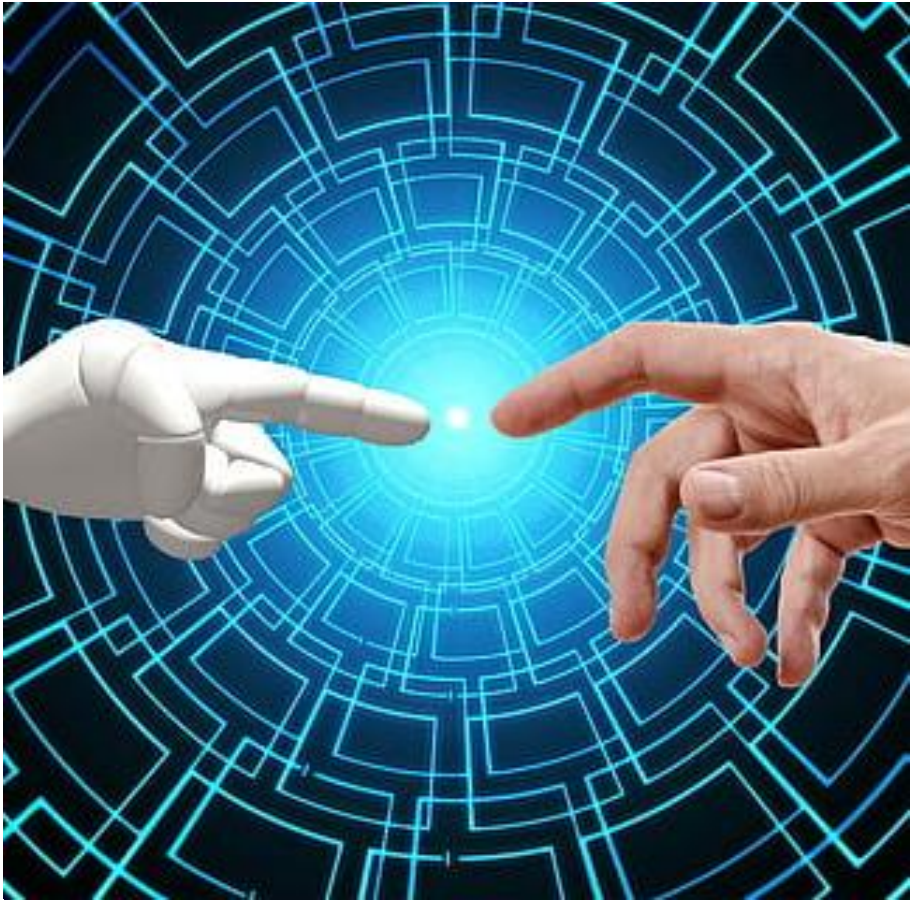
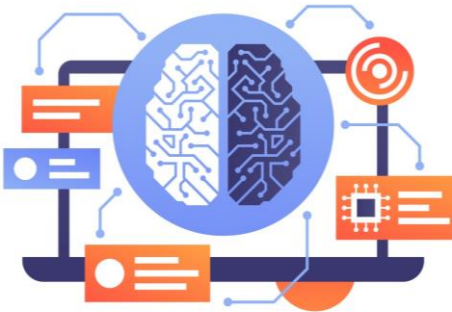
.
.

Artificial Intelligence (AI) for Engineering

COS40007

Dr. Abdur Forkan
Senior Research Fellow, AI and Machine Learning
Digital Innovation Lab

Seminar 4: 21st August 2024

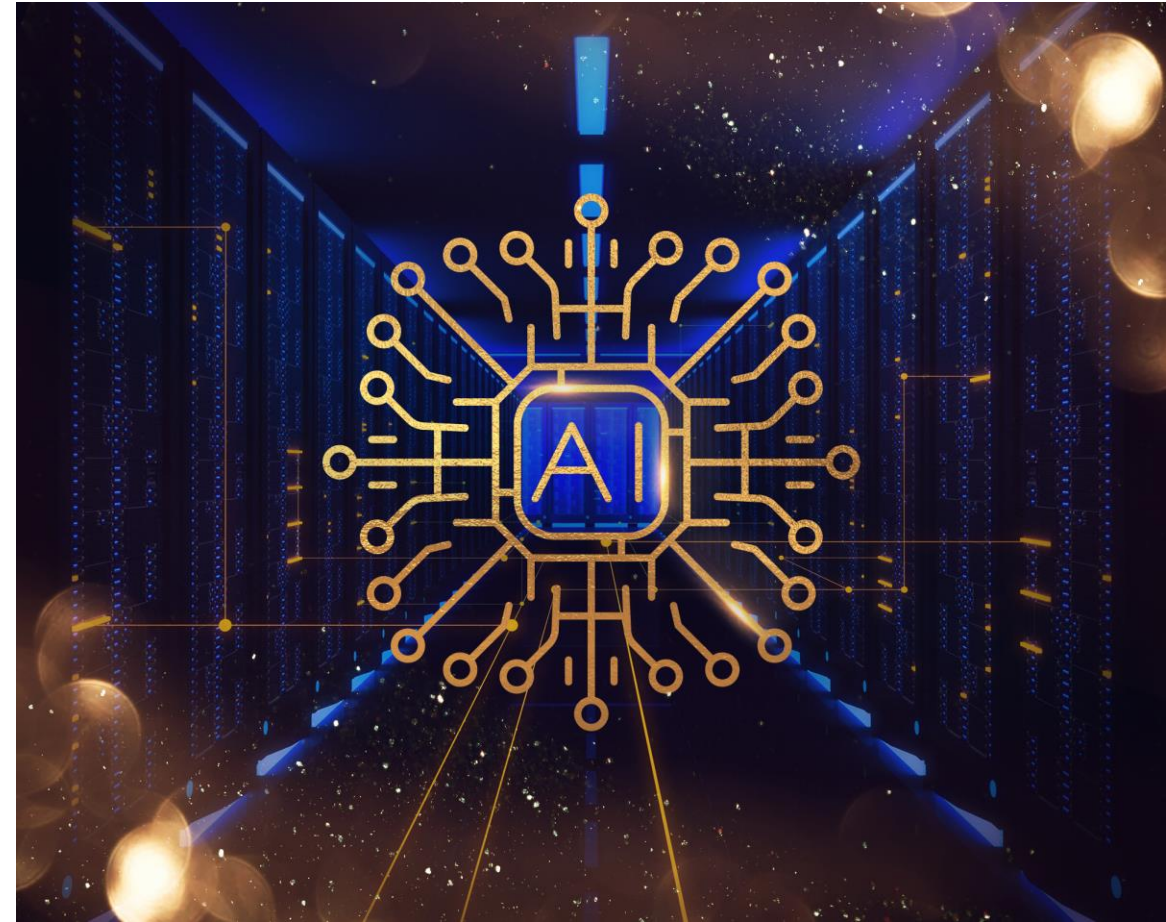


. .
. .

.
.
.
.

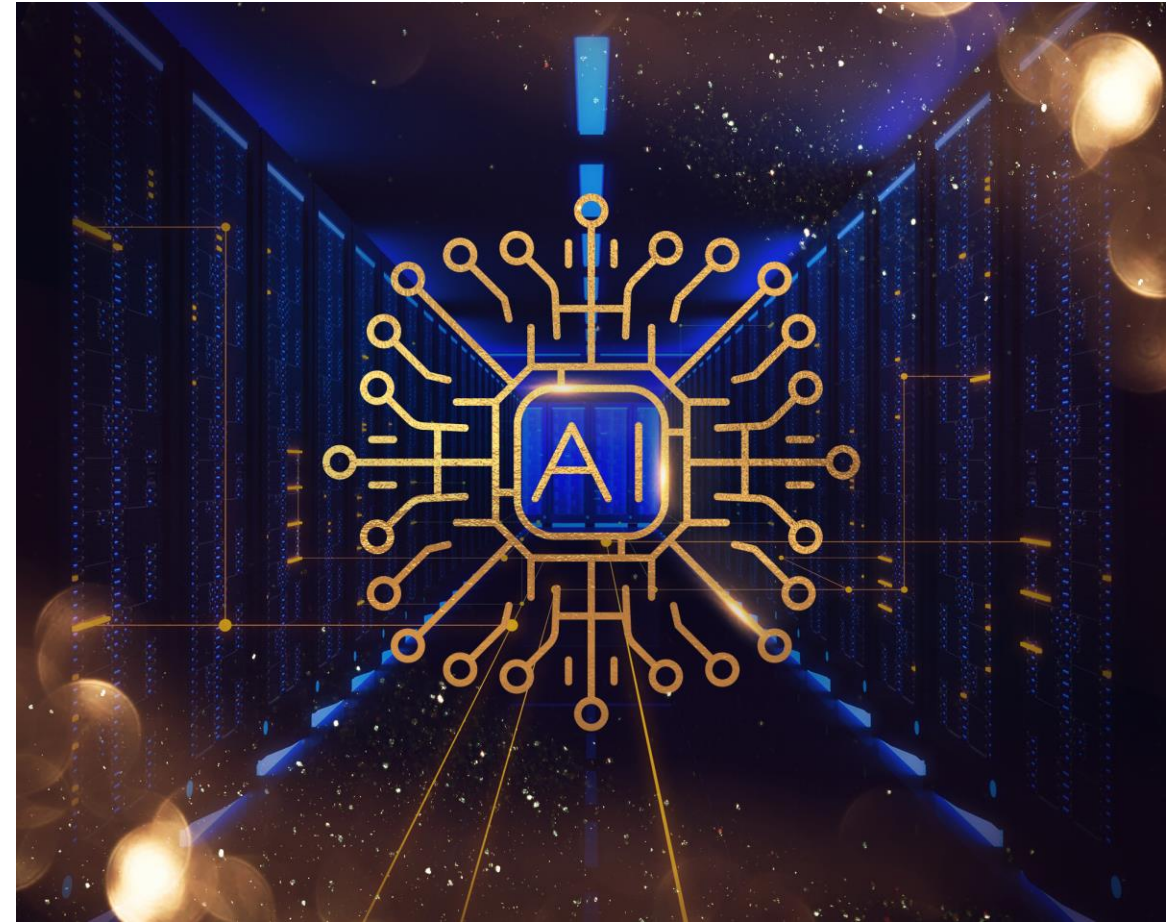
Overview

- ❑ Scikit-learn
- ❑ Examples of scikit-learn for Machine Learning
- ❑ Basics of Clustering



Required Reading

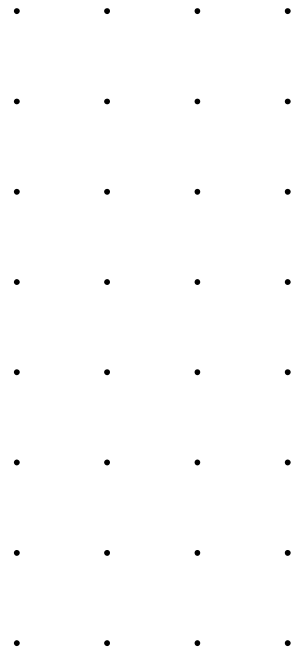
- Chapter 1-6 of “Applied Machine Learning and AI for Engineers”
- Chapter 2-6 of “Machine Learning with Pytorch and Scikit-Learn”



• • • • • • • • • •
• • • • • • • • • •

At the end of this you should be able to

- Understand what functions are available in Scikit-learn
- Understand how to perform training and validation in Scikit-learn
- Understand how to do evaluation using scikit-learn
- Understand how to do clustering using scikit-learn



A diagram consisting of a grid of small black dots. The dots are arranged in 8 horizontal rows and 4 vertical columns, forming a rectangular pattern.

```
print(data[i])
```

scikit-learn: machine Learning with Python

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



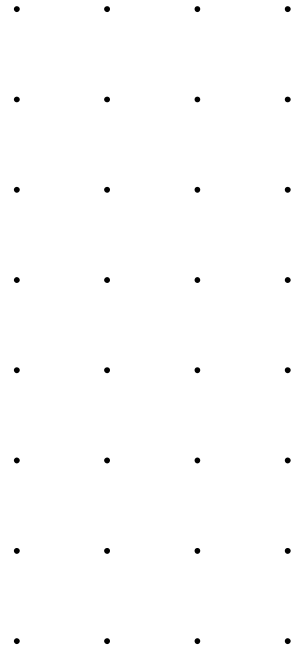
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.

Classification

Identifying to which category an object belongs to.

Applications: Activity Recognition, Image recognition.

Algorithms: SVM, nearest neighbors, random forest



Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

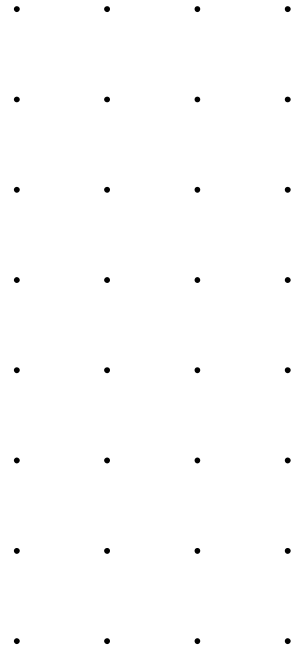
Algorithms: SVR, ridge regression, Lasso

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes.

Algorithms: k-Means, spectral clustering, mean-shift



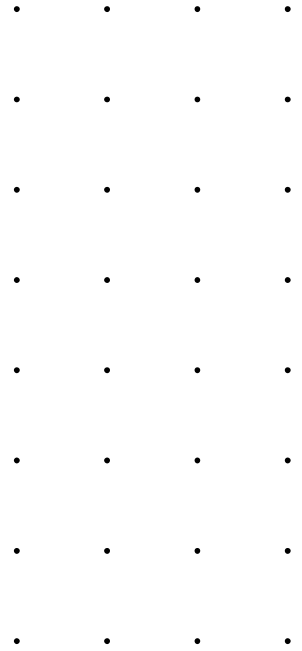
Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix

factorization

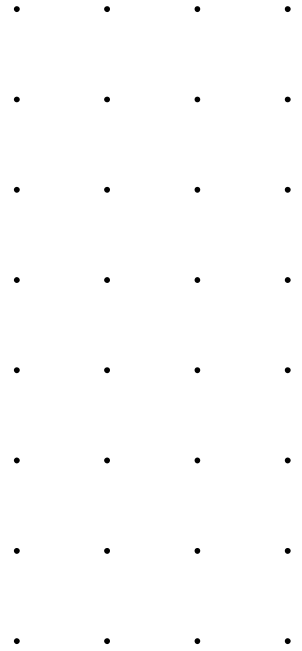


Model selection

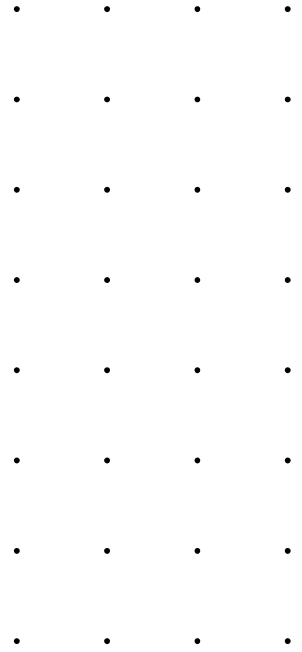
Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning.

Modules: grid search, cross validation, metrics.



Preprocessing



Feature extraction and normalization.

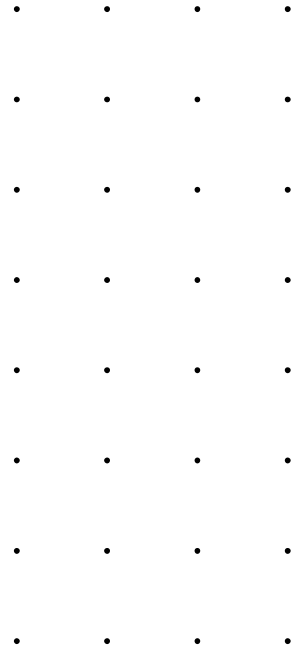
Goal: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction.

Dataset

Dataset Loading

Dataset Transformation



Most importantly

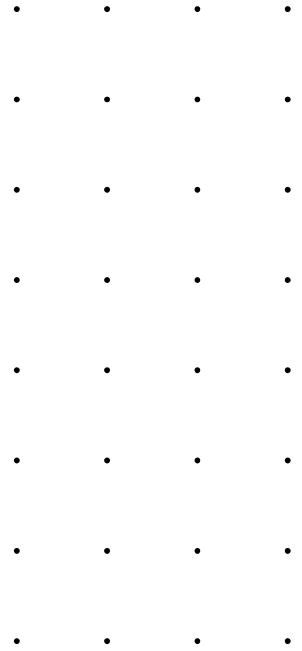
Good Documentation

Good Community Support



Scikit-Learn API

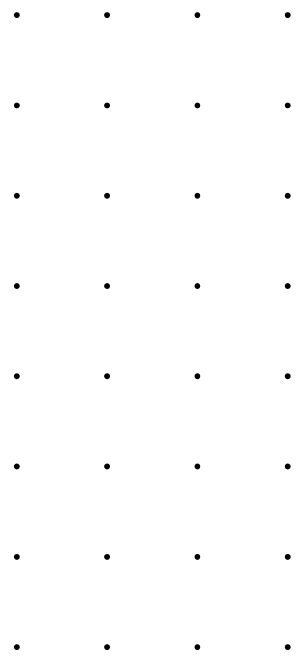
- Object-oriented interface centered around the concept of an Estimator:
- “An estimator is any object that learns from data; it maybe
 - a classification,
 - regression or
 - clustering algorithm or
- a transformer that extracts/filters useful features from
- raw data.”



Scikit-learn

Scikit learn models follow a simple, shared pattern

1. Import the model that you need to use
2. Build the model, setting its hyperparameters
3. Train model parameters on your data: Using the fit method
4. Use the model to make predictions
 - Using the predict/transform methods
 - Sometimes fit and predict/transform are implemented within the same class method



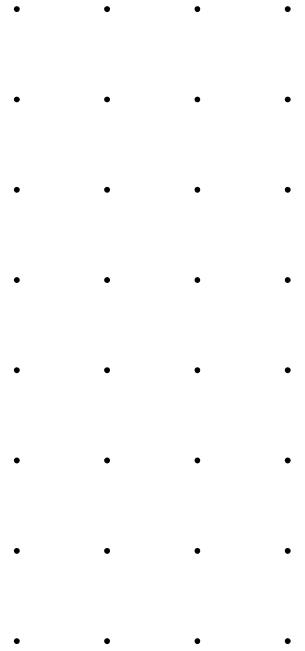
Scikit-learn

- `fit()`: learn model parameters from input data. E.g. train a classifier
- `predict()`: apply model parameters to make predictions on data ▪ E.g. predict class labels
- `fit_predict()`: fit model and make predictions ▪ E.g. apply clustering to data
- `fit_transform()`: fit model and transform data ▪ E.g. apply PCA to transform data



Estimator API

```
class Estimator(object):
    def fit(self, X, y=None):
        """Fits estimator to data. """
        # set state of ``self``
        return self
    def predict(self, X):
        """Predict response of ``X``. """
        # compute predictions ``pred``
        return pre
```



Estimators

- `fit(X,y)` sets the state of the estimator.
- - `X` is usually a 2D numpy array of shape `(num_samples, num_features)`.
- - `y` is a 1D array with shape `(n_samples,)`
- - `predict(X)` returns the class or value
- - `predict_proba()` returns a 2D array of
- shape `(n_samples, n_classes)`

```
from sklearn import svm
estimator = svm.SVC(gamma=0.001)
estimator.fit(X, y)
estimator.predict(x)
```



Transformer

```
class Transformer(Estimator):  
    def transform(self, X):  
        """Transforms the input data. """  
        # transform ``X`` to ``X_prime``  
        return X_prime
```

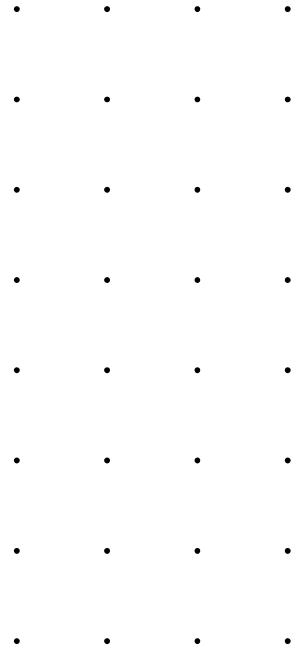
```
from sklearn import preprocessing
```

```
Xt = preprocessing.normalize(X) # Normalizer
```

```
Xt = preprocessing.scale(X) # StandardScaler
```

```
imputer =Imputer(missing_values='Nan',strategy='mean')
```

```
Xt = imputer.fit_transform(X)
```



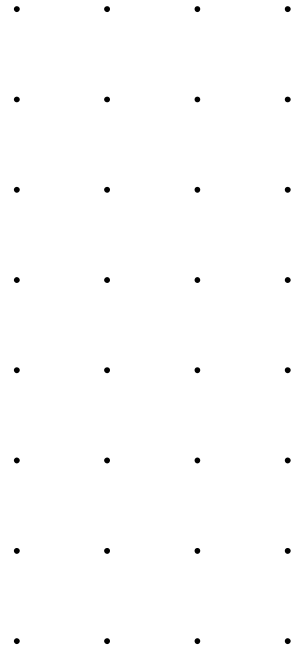
Classification

```
from sklearn import metrics
from sklearn import cross_validation as cv

plits = cv.train_test_split(X, y, test_size=0.2)
X_train, X_test, y_train, y_test = splits

model = ClassifierEstimator()
model.fit(X_train, y_train)
expected = y_test
predicted = model.predict(X_test)

print metrics.classification_report(expected, predicted)
print metrics.confusion_matrix(expected, predicted)
print metrics.f1_score(expected, predicted)
```



MSE and coefficient of determination

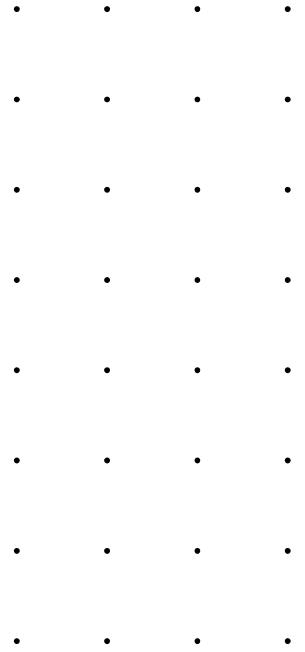
```
from sklearn import metrics
from sklearn import cross_validation as cv

splits = cv.train_test_split(X, y, test_size=0.2)
X_train, X_test, y_train, y_test = splits

model = RegressionEstimator()
model.fit(X_train, y_train)

expected = y_test
predicted = model.predict(y_test)

print metrics.mean_squared_error(expected, predicted)
print metrics.r2_score(expected, predicted)
```



Decision Tree and Random Forest

```
from sklearn import tree
from sklearn.ensemble import RandomForestClassifier
import numpy as np

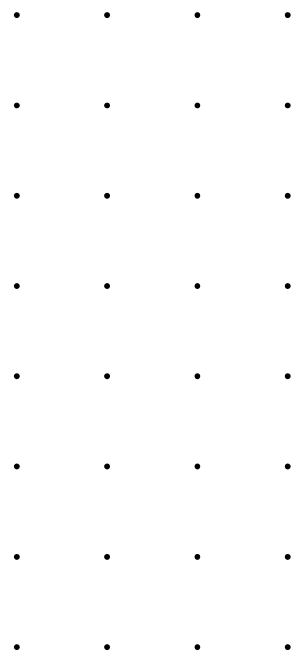
f = open("bc.train.0")
data = np.loadtxt(f)
train = data[:,1:]
trainlabels = data[:,0]

f = open("bc.test.0")
data = np.loadtxt(f)
test = data[:,1:]
testlabels = data[:,0]

clf = tree.DecisionTreeClassifier()
clf.fit(train,trainlabels)
prediction = clf.predict(test)

rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(train,trainlabels)
prediction2 = rfc.predict(test)

err = 0
err2 = 0
for i in range(0, len(prediction), 1):
    if(prediction[i] != testlabels[i]):
        err += 1
    if(prediction2[i] != testlabels[i]):
        err2 += 1
err = err/len(testlabels)
err2 = err2/len(testlabels)
print(err,err2)
```

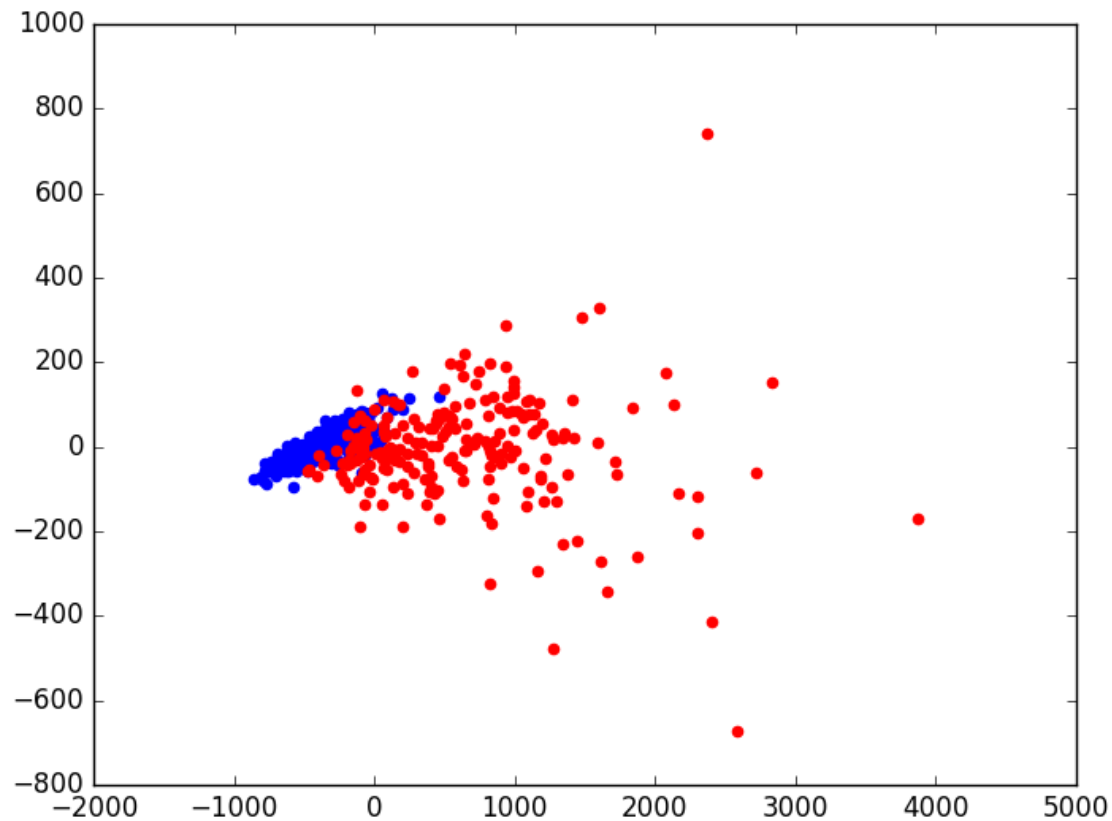


.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.

The logo for Swinburne University of Technology. It consists of a black square on the left containing the word "SWINBURNE" in white, with a small white asterisk on either side of the word "NE". To the right of the black square is a red rectangle containing the text "SWINBURNE UNIVERSITY OF TECHNOLOGY" in white, stacked in three lines.

Dimensionality reduction and visualization with PCA

PCA plot of breast cancer data (output of program in previous slide)



K-means PCA plot in scikit-learn

```
from sklearn.cluster import KMeans
from matplotlib import pyplot as plt
import numpy as np
from sklearn.decomposition import PCA

f = open("bc")
data = np.loadtxt(f)
X = data[:,1:]
Y = data[:,0]

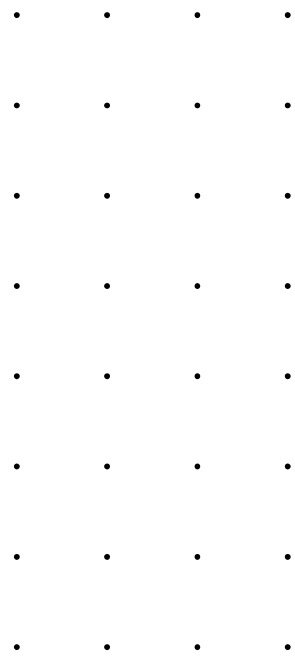
pca = PCA(n_components=2)
pca.fit(X)
newdata = pca.transform(X)
clustering = KMeans(n_clusters=2,init='random').fit(X)

x = []
y = []
for i in range(0, len(newdata), 1):
    if(clustering.labels_[i] == 0):
        x.append(newdata[i][0])
        y.append(newdata[i][1])

plt.scatter(x, y, color='blue')

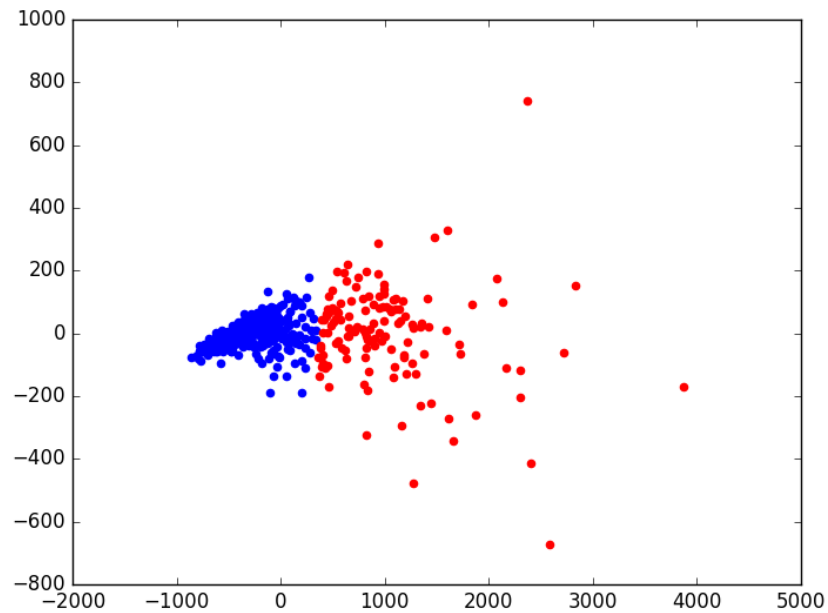
x = []
y = []
for i in range(0, len(newdata), 1):
    if(clustering.labels_[i] == 1):
        x.append(newdata[i][0])
        y.append(newdata[i][1])

plt.scatter(x, y, color='red')
plt.show()
```

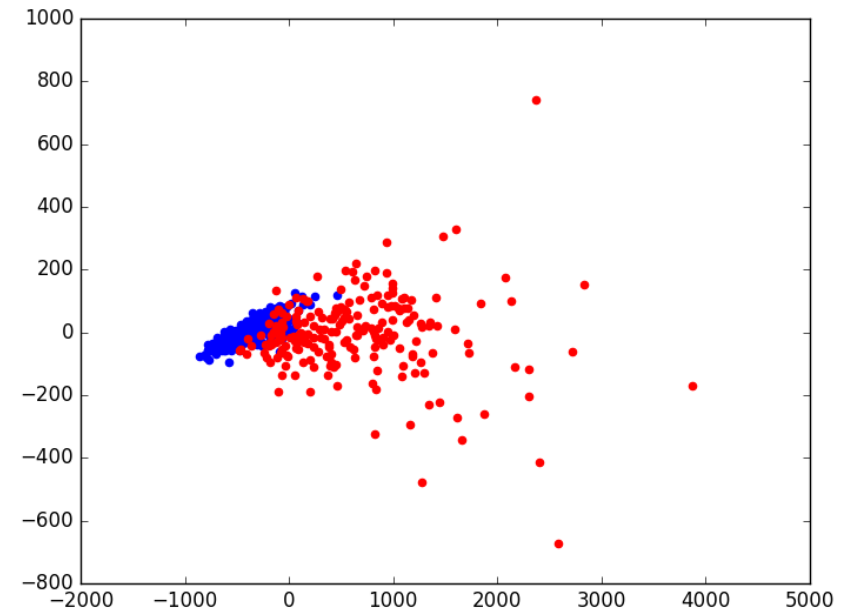


K-means PCA plot in scikit-learn

PCA plot of breast cancer data
colored by k-means labels



PCA plot of breast cancer data
colored by true labels



Learn, Practice and Enjoy the AI journey

