

Classification.

Categorizing some unknown items into a discrete set of categories or classes.

• Supervised Learning approach.

Algorithms.

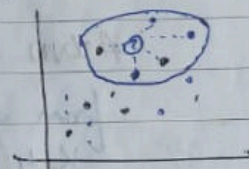
- Decision Trees, • Naive Bayes, • Linear Discriminant Analysis
- K-Nearest Neighbor, • Logistic Regression,
- Neural Network, • Support Vector Machines (SVM)

K-Nearest Neighbor (KNN)

It takes bunch of labeled points and uses them to learn how to label other points.

1. Pick a value for k
2. Cal. the distance of unknown case from all cases.
3. Select the k -observations in the training data that are "nearest" to the unknown data point.
4. Predict the response of the unknown data point using the most popular response value from the k -nearest neighbors

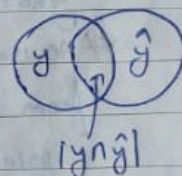
$$Dis(x_1, x_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$



Evaluation Metrics [Accuracy]

Jaccard index

$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|} = \frac{|y \cap \hat{y}|}{|y| + |\hat{y}| - |y \cap \hat{y}|}$$



High Accuracy at 1.0

F1-Score

$$\text{Precision} = TP / (TP + FP)$$

$$\text{Recall} = TP / (TP + FN)$$

$$F1\text{-Score} = 2 \times (\text{prec} \times \text{rec}) / (\text{prec} + \text{rec})$$

Confusion Matrix

	\hat{y}	Churn=1	
		6 TP	9 FN
y	Churn=0	1 FP	24 TN

Churn=1 Churn=0

High Accuracy at 1.0

ML with Python

Page No. _____
Date _____

ML is the subfield of computer science that gives "computers the ability to learn without being explicitly programmed".

- Regression / Estimation
 - Predicting continuous value
- Classification
 - Predicting the item class/category of a case
- Clustering
 - Finding the structure of data; summarization
- Associations
 - Associating frequent co-occurring items/events
- Anomaly detection
 - Discovering abnormal and unusual cases
- Sequence mining
 - Predicting next event
- Dimension Reduction
 - Reducing the size of data (PCA)
- Recommendation systems
 - Recommending items

Python Libraries for ML

Numpy, SciPy, matplotlib, pandas, sklearn

Supervised

- ↳ Regression
- ↳ Classification

Unsupervised

- ↳ Dimension reduction
- ↳ Density estimation
- ↳ Market basket analysis
- ↳ Clustering

Regression

X: Independent 1 or more data.
can be categorical also.

Y: Dependent only 1 and
must be continuous value

Ordinal Regression

• Poisson

• Fast forest quantile

• Linear, Polynomial, Lasso, stepwise, Ridge

• Bayesian linear

• Neural Network

• Decision forest

• Boosted decision tree

• KNN (k-nearest neighbors)

Logloss

$$\text{Log Loss} = -\frac{1}{n} \sum (y \times \log(\hat{y}) + (1-y) \times \log(1-\hat{y}))$$

High Accuracy at 0.0

y	\hat{y}
churn	Predicted churn
1	0.91
1	0.13
0	0.04
0	0.23
0	0.43

```
X = df[['A', 'B', 'C', ...]].values #.astype(float)
```

```
y = df['Z'].values
```

Normalize Data

```
X = preprocessing.StandardScaler().fit(X).transform(X.astype(float))
```

```
from sklearn.model_selection import train_test_split
X_train, y_train, X_test, y_test = train_test_split(X, y, test_size=0.2, random_state=4)
```

KNN

```
from sklearn.neighbors import KNeighborsClassifier
K = 4
```

```
neigh = KNeighborsClassifier(n_neighbors=K).fit(X_train, y_train)
yhat = neigh.predict(X_test)
```

Accuracy

```
from sklearn import metrics
```

```
print("Train set Accuracy: ", metrics.accuracy_score
```

```
(y_train, neigh.predict(X_train))
print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat))
```

For other Ks

```
Ks = 10
```

```
mean_acc = np.zeros((Ks-1))
```

```
std_acc = np.zeros((Ks-1))
```

```
for n in range(1, Ks):
```

```
    neigh = KNeighborsClassifier(n_neighbors=n).fit(X_train, y_train)
    yhat = neigh.predict(X_test)
```



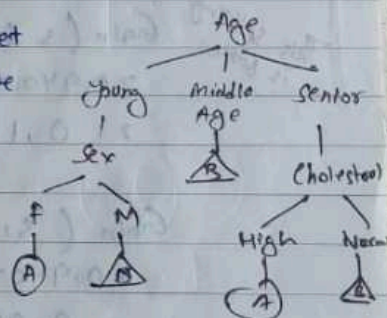
```

mean_acc[n-1] = metrics.accuracy_score(y_test, y_hat)
std_acc[n-1] = np.std(y_hat == y_test) / np.sqrt(y_hat.shape[0])
mean_acc =
# plotting
plt.plot(range(1, ks), mean_acc, 'g')
plt.fill_between(range(1, ks), mean_acc - 1 * std_acc, mean_acc + 1 * std_acc,
alpha = 0.10)
plt.fill_between(range(1, ks), mean_acc - 3 * std_acc, mean_acc + 3 * std_acc,
alpha = 0.10, color = 'g')
plt.legend(['Accuracy', '+/- 1x std', '+/- 3x std'])
plt.ylabel('Accuracy')
plt.xlabel('Number of Neighbors (k)')
plt.tight_layout()
plt.show()
print("The best accuracy was with ", mean_acc.max(), "with k=",
mean_acc.argmax() + 1)

```

Decision Tree.

1. Choose an attribute from your dataset
2. Calculate the significance of attribute in splitting of data.
3. Split data based on the value of the best attribute.
4. Go to step 1.



Entropy : measure of randomness or uncertainty.

1 Drug A

3 Drug A

7 Drug B

5 Drug B

Entropy is low ✓

Entropy is high X

0 Drug A

4 Drug A

8 Drug B

4 Drug B

Entropy is 0 ✓

Entropy = 1 X X

Pure node

Evaluation

```

from sklearn.metrics import jaccard_score
jaccard_score(y_test, y_hat, pos_label=0)

# Confusion Matrix in previous page
print(confusion_matrix(y_test, y_hat, labels=[1, 0]))

cm = confusion_matrix(y_test, y_hat, labels=[1, 0])
np.set_printoptions(precision=2)
plt.figure()
plt.matshow(cm, cmap=plt.cm.Blues, title='Confusion Matrix')
plt.colorbar()
print(classification_report(y_test, y_hat))

```

```

from sklearn.metrics import log_loss
log_loss(y_test, y_hat_prob)

```

General Cost function

$$\sigma(\theta^T X) \rightarrow P(y=1|X)$$

$$\text{cost function: } \text{cost}(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\hat{y}^i, y^i)$$

$$\text{cost}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y=1 \\ -\log(1-\hat{y}) & \text{if } y=0 \end{cases}$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{y}^i) + (1-y^i) \log(1-\hat{y}^i)]$$

Best parameter of our model?

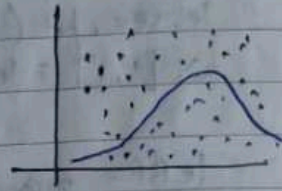
by minimize the cost function

How to min the cost function.

Using Gradient Descent

SVM (Support Vector Machines)

It works by mapping data to high dimensional feature space so that data points can be categorized, even when the data are not otherwise linearly separable.



Mapping data into higher dimensional space is called kernelling.

• Linear, • Polynomial, • RBF • Sigmoid

Dataset:	2D	Clump	UnifSize	UnifShap	BareNuc	class
	-	5			1	2
	-	5			10	2
	-	3			2	2
	-	6			4	2
	-	4			1	2

In class 2 = Benign 4 = malignant

```
ax = cell_df[cell_df['class'] == 4][0:50].plot(kind='scatter', x='Clump',
y='UnifSize', color=' ', label='malignant');
cell_df[cell_df['class'] == 2][0:50].plot(kind='scatter', x='Clump',
y='UnifSize', color=' ', label='benign', ax=ax)
plt.show()
```

```
cell_df.dtypes # Let BareNuc be object
cell_df = cell_df[pd.to_numeric(cell_df['BareNuc'], errors='coerce')]
cell_df['BareNuc'] = cell_df['BareNuc'].astype('int')
```

```
feature_df = cell_df[['Clump', ..., 'mit']]
X = np.asarray(feature_df)
```

```
cell_df['class'] = cell_df['class'].astype('int')
y = np.asarray(cell_df['class'])
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=4)
```


#1 Convert the dataset into only numeric dataset then normalization
from sklearn.preprocessing import StandardScaler

```
X = df.values[:, 1:]
```

```
X = np.nan_to_num(X)
```

```
clus_dataset = StandardScaler().fit_transform(X)
```

```
clusterNum = 3
```

```
KMeans = KMeans(n_clusters = clusterNum,  
                 n_init = 12)
```

```
cl_means = fit(X)
```

```
labels = cl_means.labels_
```

```
print(labels)
```

→ New column.

```
df["clus_num"] = labels
```

```
df.groupby("clus_num").mean()
```

```
area = np.pi * (X[:, 1])**2
```

```
plt.scatter(X[:, 0], X[:, 3], s=area, c=labels.astype(np.float),  
            alpha=0.5)
```

```
plt.xlabel('Age')
```

```
plt.ylabel('Income')
```

```
plt.show()
```

```
from mpl_toolkits.mplot3d import Axes3D
```

```
fig = plt.figure(figsize=(8,6))
```

```
plt.clf()
```

```
ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=48, azim=134)
```

```
plt.cla()
```

```
ax.set_xlabel('Education')
```

```
ax.set_ylabel('Age')
```

```
ax.set_zlabel('Income')
```

```
ax.scatter(X[:, 1], X[:, 0], X[:, 3], c=labels.astype(np.float))
```



```
import sklearn.cluster import KMeans
import sklearn.datasets as make_blobs
```

Page No.:

Date:

Clustering

K-Means Clustering

On randomly generated dataset

```
np.random.seed(0)
```

```
X, y = make_blobs(n_samples=5000, centers=[[4, 4], [-2, -1], [2, 3], [1, 1]],
                  cluster_std=0.9)
```

```
plt.scatter(X[:, 0], X[:, 1], marker='.'))
```

Setting up K-Means.

```
K_means = KMeans(init="K-means++", n_clusters=4, n_init=12)
```

```
K_means.fit(X)
```

```
K_means.labels = K_means.labels_
```

```
K_means.cluster_centers = K_means.cluster_centers_
```

Plotting.

```
fig = plt.figure(figsize=(6, 4))
```

```
colors = plt.cm.Spectral(np.linspace(0, 1, len(set(K_means.labels_))))
```

```
ax = fig.add_subplot(1, 1, 1)
```

```
for k, col in zip(range(len([[4, 4], [-2, -1], [2, 3], [1, 1]])), colors):
```

```
    my_members = (K_means.labels == k)
```

```
    cluster_center = K_means.cluster_centers[k]
```

```
    ax.plot(X[my_members, 0], X[my_members, 1], 'w',
            markerfacecolor=col, marker='o')
```

```
    ax.plot(cluster_center[0], cluster_center[1], 'o', markerfacecolor=col,
            markeredgecolor='k', markersize=6)
```

```
ax.set_title('K-Means')
```

```
ax.set_xticks(())
```

```
ax.set_yticks(())
```

```
plt.show()
```

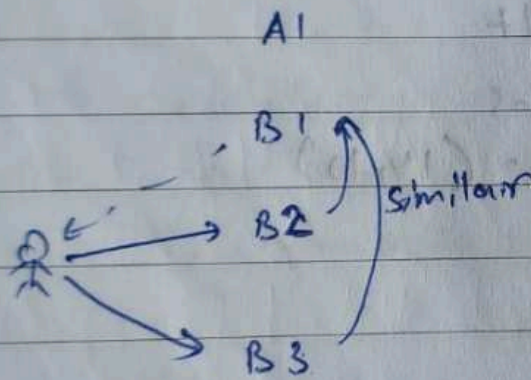

$([1, 0, 0])$, marker 0101, marker size
 $\alpha = 0.75$

plt.show()

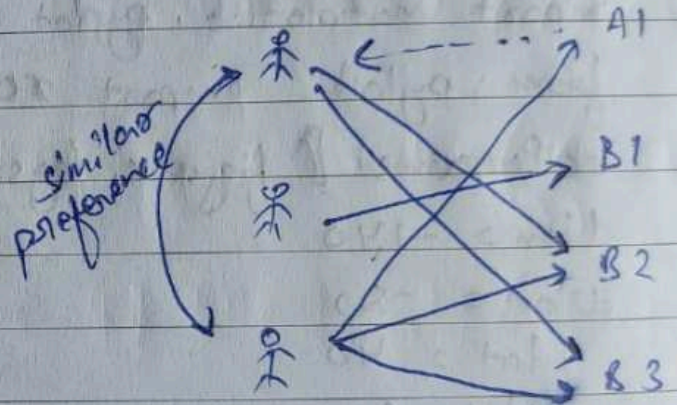
Recommender System.

It captures the patterns of people's behavior and use to predict what else they might want or like.

Content Based



Collaborative Filtering



Teacher's Signature


```

from scipy import ndimage
from scipy.cluster import hierarchy
from sklearn.spatial import distance
from sklearn.datasets import make_blobs
# on randomly generated data

```

```

X, y = make_blobs (n_samples = 50, centers = [[4, 4], [-2, -1], [1, 1], [10, 4]], cluster_std = 9)

```

```

plt.scatter (X[:, 0], X[:, 1], marker = 'o')

```

```

from sklearn.cluster import AgglomerativeClustering

```

```

agglo = AgglomerativeClustering (n_clusters = 4, linkage = 'average')
agglo.fit (X, y)

```

```

plt.figure (figsize = (6, 4))

```

```

x_min, x_max = np.min (X, axis = 0), np.max (X, axis = 0)

```

```

X = (X - x_min) / (x_max - x_min) # avg. dist from X1

```

```

for i in range (X.shape[0]):

```

```

    plt.text (X[i, 0], X[i, 1], str (y[i]),

```

```

            color = plt.cm.nipy_spectral (agglo.labels_[i] / 10.),

```

```

            fontdict = {'weight': 'bold', 'size': 14})

```

```

plt.xticks ([])

```

```

plt.yticks ([])

```

```

plt.scatter (X[:, 0], X[:, 1], marker = '.')

```

```

plt.show ()

```

```

Z = distance.distance_matrix (X, X)

```

```

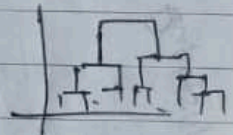
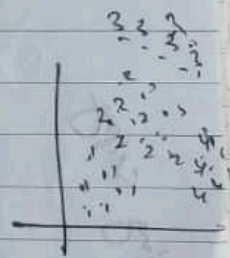
Z = hierarchy.linkage (Z, 'complete')

```

```

dendro = hierarchy.dendrogram (Z)

```



	Action	Supernatural	comedy	Sci-Fi
A _u	1	0	1	1
A _s	1	0	0	0
A _g	1	1	0	1
A _u	0.2	0	0.3	0.3
A _s	0.2	0	0	0
A _g	0.2	0.2	0	0.3

* user profile.

0.8
0.2
0.6

A_u will be above all recommendation.
Suppose Drama come into the list, this movie will never ~~also~~ recommend Drama.

```
import pandas as pd
from math import sqrt
import matplotlib.pyplot as plt
import numpy as np
```

```
movies_df = pd.read_csv('movies.csv')
ratings_df = pd.read_csv('ratings.csv')
movies_df.head()
```

movieId	title	genres
1	A (1995)	Adv Ant comedy fantasy
2	B (1995)	- - -
3	C (1995)	- - -
4	D (1995)	- - -
5	E (1995)	- - -

```
movies_df['year'] = movies_df['title'].str.extract('(\d\d\d\d)')
# expand = first
```

```
movies_df['year'] = movies_df['year'].str.extract('(\d\d\d\d)')
# expand = first
```

```
movies_df['year'] = movies_df['title'].str.replace('(\d\d\d\d)', '1995')
```

```
movies_df['title'] = movies_df['title'].apply(lambda x: x.strip())
movies_df.head()
```



```
userInput = [
    {'title': 'Star Wars R1', 'rating': 5},
    {'title': 'Star Wars R2', 'rating': 2.5},
    ...
]
```

```
inputMovies = pd.DataFrame(userInput)
```

Filtering out the movies by title:

```
inputId = movies_df[movies_df['title'].isin(inputMovies['title'].tolist())]
```

```
inputMovies = pd.merge(inputId, inputMovies)
```

```
inputMovies = inputMovies.drop('genres', 1).drop('year', 1)
```

```
userMovies = moviesWithGenres_df[moviesWithGenres_df['movieId'].isin(
    inputMovies['movieId'].tolist())]
```

```
userMovies.head(1)
```

```
out
movieId
0      1      A
```

We only need genre table.

```
userMovies = userMovies.reset_index(drop=True)
```

```
userGenreTable = userMovies.drop('movieId', 1).drop('title', 1).
    drop('genres', 1).drop('year', 1)
```

```
inputMovies['rating']
```

```
out
(1, 2.5)
```

product.

```
userProfile = userGenreTable.transpose().dot(inputMovies['rating'])
userProfile
```

```
out
Adv    10.0
Ad1     .
/       .
/       .
/       .
```


cell movie

```
genreTable = moviesWithGenres_df.set_index(moviesWithGenres_df['movieId'])
```

```
genreTable = genreTable.drop('movieId', 1).div
```

```
recommendationTable_df = (genreTable * userProfile).sum(axis=1)  
/ (userProfile.sum())
```

```
recommendationTable_df = r--_df.sort_values(ascending=False)  
.head()
```

out

movieId

5018

0.74

:

0.73

:

The recommendation table.

```
movies_df.loc[movies_df['movieId'].isin(recommendationTable_df  
.head(20).keys())]
```


out:

movieId	title	genres	year
	A		1999
	B		"
	C		"
	D		"
	E		"

movies_df[genres] = movies_df.genres.str.split('|')

movies_df.head(1)

out

movieId	title	genres	year
0	1	A [Adv, Ani, comedy, fantasy]	1999

moviesWithGenres_df = movies_df.copy()

from index, row in movies_df.iterrows():

for genre in row['genres']:

moviesWithGenres_df.at[index, genre] = 1

moviesWithGenres_df = moviesWithGenres_df.fillna(0)

moviesWithGenres_df.head(1)

out

movieId	title	genres	year	Adv	Ani	com	dram	fantasy	horr	musi	rom	sci	thriller	war	western
0	1	A	1999	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

ratings_df.head(1)

out

userId	movieId	rating	timestamp
0	1	2.5	1207927694

ratings_df = ratings_df.drop('timestamp', 1)

Implementing recommender systems

• Memory-based

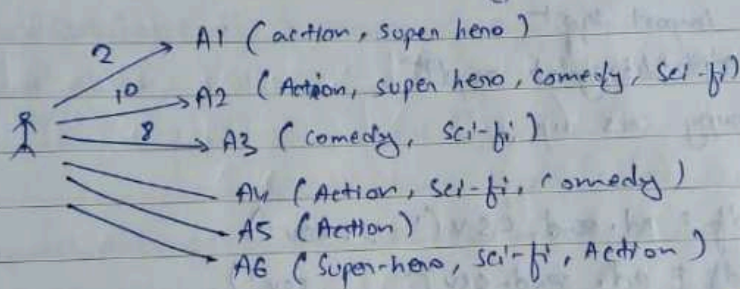
- Uses the entire user item dataset to generate a recommendation
- Uses statistical tech.

• Model based

- Develops a model of users in an attempt to learn their preferences.
- Created using ML tech. like regression, clustering, classification.

Disadvantage of content based.

Model only recommends that users already have seen, related to that. But when a new thing comes, the model don't recommend it.



Input User matrix		movies matrix				
	User		Action	Superhero	comedy	Sci-fi
A1	2	A1	1	1	0	0
A2	10	A2	1	1	1	1
A3	8	A3	0	0	1	1

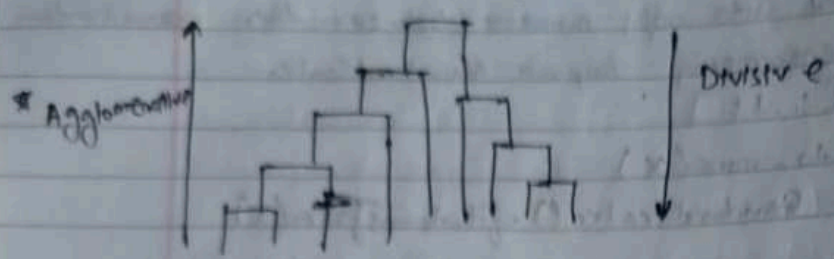
	Action	Superhero	comedy	Sci-fi
A1	2	2	0	0
A2	10	10	10	10
A3	0	0	8	8

(Weighted Genre matrix)

User profile	12	12	18	18
--------------	----	----	----	----

Normalized	0.2	0.2	0.3	0.3
------------	-----	-----	-----	-----

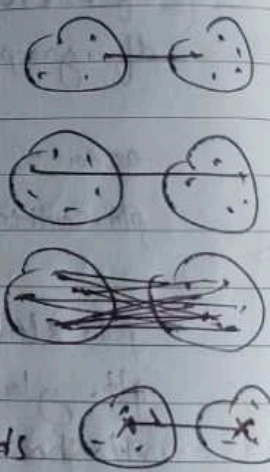
Hierarchical clustering



- algo-
1. Create n clusters, one for each data point
 2. Compute the proximity matrix
 3. Repeat
 - i. Merge the 2 closest clusters
 - ii. Update the proximity matrix
 4. Until only a single cluster remains

Distance b/w clusters

- Single - Linkage Clustering
 - Min distance b/w clus.
- Complete - Linkage
 - Max dis.
- Average Li
 - Avg.
- Centroid Li
 - Distance b/w cluster centroids



Adv.

- Doesn't required no. of clusters to be specified
- Easy to implement
- Produces a dendrogram, which helps with understanding the data

Disadv.

- Can never undo any previous steps throughout the algo
- Generally has long runtimes
- Sometimes difficult to identify the no. of clusters by the dendrogram.

`cnf_matrix = confusion_matrix(y_test, y_hat, label=[2, 4])`
`np.set_printoptions(precision=2)`
`print(classification_report(y_test, y_hat))`
`plt.figure()`
`plt.confusion_matrix(cnf_matrix, classes=['Benign(2)', 'Malignant(4)'],`
`normalize=False, title='Confusion matrix')`

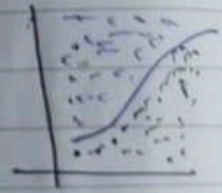
Errors

`from sklearn.metrics import f1_score`
`f1_score(y_test, y_hat, average='weighted')`

`from sklearn.metrics import jaccard_score`
`jaccard_score(y_test, y_hat, pos_label=2)`

Logistic Regression

`churn_df = churn_df[['tenure', 'churn']]`
`churn_df['churn'] = churn_df['churn'].astype(int)`



`X = np.asarray(churn_df[['tenure']])`
`y = np.asarray(churn_df['churn'])`

Normalization of Dataset

`from sklearn import preprocessing`
`X = preprocessing.StandardScaler().fit(X).transform(X)`

`from sklearn.model_selection import train_test_split`
`X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)`

`from sklearn.linear_model import LogisticRegression`
`from sklearn.metrics import confusion_matrix`
`LR = LogisticRegression(C=0.01, solver='liblinear').fit(X_train, y_train)`
`proba = LR.predict_proba(X_test)`
`y_hat = LR.predict(X_test)`

main code

Page No.:

Date:

```
from sklearn import svm
clf = svm.SVC(kernel='rbf')
clf.fit(X_train, y_train)
```

```
y_hat = clf.predict(X_test)
```

#Evaluation.

```
from sklearn.metrics import classification_report, confusion_matrix
import itertools
```

```
def plot_confusion_matrix(cm, classes, normalize=False, title='Confusion
Matrix', cmap=plt.cm.Blues):
```

```
    if normalize:
```

```
        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
```

```
        print("Normalized Confusion Matrix")
```

```
    else:
```

```
        print("Confusion Matrix without normalization")
```

```
    print(cm)
```

```
    plt.imshow(cm, interpolation='nearest', cmap=cmap)
```

```
    plt.title(title)
```

```
    plt.colorbar()
```

```
    tick_marks = np.arange(len(classes))
```

```
    plt.xticks(tick_marks, classes, rotation=45)
```

```
    plt.yticks(tick_marks, classes)
```

```
    fmt = '.2f' if normalize else 'd'
```

```
    thresh = cm.max() / 2.
```

```
    for i, j in itertools.product(range(cm.shape[0]), range(cm
    shape[1])):
```

```
        plt.text(j, i, format(cm[i, j], fmt),
```

```
                  horizontalalignment='center')
```

```
        color = "white" if cm[i, j] > thresh else "black"
```

```
    plt.tight_layout()
```

```
    plt.ylabel("True label")
```

```
    plt.xlabel("Predicted label")
```

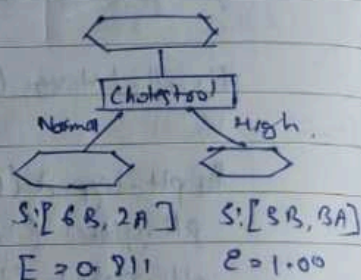
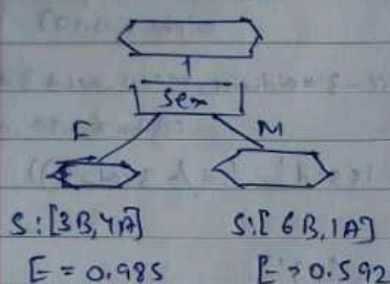

$$\text{Entropy} = -P(A) \log_2(P(A)) - P(B) \log_2(P(B))$$

↑
Total A / Total (A+B)

$S = [9 \text{ (mg B)}, 5 \text{ (mg A)}]$ Before splitting.

$$E = - (9/14) \log_2(9/14) - (5/14) \log_2(5/14)$$

$$= 0.940$$



which one is better choice.

The tree with higher Information Gain after splitting

Information Gain = (Entropy before split) - (weighted entropy after split)

Entropy decreases Information gain Increases.

This splitting is better

Gain (S, Sex)

$$= 0.940 - [(7/14)0.985 + (7/14)0.592]$$

$$= 0.151$$

Gain (S, Cholesterol)

$$= 0.940 - [(8/14)0.911 + (6/14)1.0]$$

$$= 0.048$$

Code

```
X = my_data[['Age', 'Sex', 'BP', 'Cholesterol', 'Na-to-K Ratio']]
```


Convert categorical variables into dummy (one hot encoding)

```
from sklearn import preprocessing  
le-sex = preprocessing.LabelEncoder()  
le-sex.fit(['F', 'M'])  
X[:, 1] = le-sex.transform(X[:, 1])
```

```
le-BP = preprocessing.LabelEncoder()  
le-BP.fit(['low', 'Normal', 'High'])  
X[:, 2] = le-BP.transform(X[:, 2])
```

```
le-chol = preprocessing.LabelEncoder()  
le-chol.fit(['Normal', 'High'])  
X[:, 3] = le-chol.transform(X[:, 3])  
print(X[0:5])
```

Code

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,  
                                                    random_state=3)
```

```
from sklearn.tree import DecisionTreeClassifier  
import sklearn.tree as tree
```

```
drugTree = DecisionTreeClassifier(criterion="entropy", max_depth=4)  
drugTree.fit(X_train, y_train)
```

```
predTree = drugTree.predict(X_test)
```

Accuracy.

```
from sklearn import metrics
```

```
import matplotlib.pyplot as plt
```

```
print("Decision Tree's Accuracy: ", metrics.accuracy_score(y_test, predTree))
```

It is accuracy near 1.0 ✓

```
tree.plot_tree(drugTree)
```

```
plt.show()
```


Multiple LR

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

How to estimate θ ?

- Ordinary Least Square

Tries to estimate the value of θ by minimizing MSE.

- An Optimization algo

• Gradient Descent for large dataset.

```
from sklearn import linear_model
```

```
reg = linear_model.LinearRegression()
```

```
x = np.asarray(train[['S', 'C', 'FC']])
```

```
y = np.asarray(train[['CO2']])
```

```
reg.fit(x, y)
```

```
print(reg.coef_)
```

Evaluation

```
y_hat = reg.predict(test[['S', 'C', 'FC']])
```

```
x = np.asarray(test[['S', 'C', 'FC']])
```

```
y = np.asarray(test[['CO2']])
```

```
print("Residual sum of squares: %.2f"
```

```
      % np.mean((y_hat - y) ** 2))
```

```
print("Variance score: %.2f" % reg.score(x, y))
```

Polynomial Regression

```
from sklearn.preprocessing import PolynomialFeatures
```

```
from sklearn import linear_model
```

```
train_x = np.asarray(train[['S']])
```

```
train_y = np.asarray(train[['CO2']])
```

```
test_x = np.asarray(test[['S']])
```

```
test_y = np.asarray(test[['CO2']])
```

```
poly = PolynomialFeatures(degree=2)
```

```
train_x_poly = poly.fit_transform(train_x)
```

```
clf = linear_model.LinearRegression()
```

```
train_y_ = clf.fit(train_x_poly, train_y)
```



```
plt.scatter (train.E5, train.CO2, color = 'b')
```

```
xx = np.arange (0.0, 10.0, 0.1)
```

```
yy = clf.intercept_ [0] + clf.coef_ [0] [1] * xx + clf.coef_ [0] [2]
      * np.power (xx, 2)
```

```
plt.plot (xx, yy, '-r')
```

```
plt.xlabel ("Engine Size")
```

```
plt.ylabel ("Emission")
```

Evaluation .

```
import from sklearn.metrics import r2_score
```

```
test-x-poly = poly.transform (test-x)
```

```
test-y- = clf.predict (test-x-poly)
```

```
print ("MAE : %.2f " % np.mean (np.absolute (test-y- - testy)))
```

```
print ("MSE : %.2f " % np.mean ((test-y- - test-y) ** 2))
```

```
print ("R2-score : %.2f " % r2_score (test-y, testy-))
```

[Logistic graph. $Y = a + \frac{b}{1 + e^{(x-d)}}$]

```
x = np.arange (-5.0, 5.0, 0.1)
```

```
Y = 1 - 4 / (1 + np.power (3, x-2))
```

```
plt.plot (x, Y) ]
```


* RSE = Relative Square Error,
$$\frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

$$R^2 = 1 - RSE$$

Lab for SLR

Suppose we need specific column from a df

`cdf = df[['ES', 'C', 'FC', 'CO2']]`

Hist

`viz = cdf[['ES', 'C', 'FC', 'CO2']]`

`viz.hist()`

Scatter plot.

`plt.scatter(cdf.FC, cdf.CO2, color='b')`

Splitting data

`mask = np.random.rand(len(cdf)) < 0.8`

`train = cdf[mask]`

`test = cdf[~mask]`

Training data distribution in scatter plot

`plt.scatter(train.ES, train.CO2)`

Model

`from sklearn import linear_model`

`regr = linear_model.LinearRegression()`

`train_x = np.asarray(train[['ES']])` # convert a list into array.

`train_y = np.asarray(train[['CO2']])`

`regr.fit(train_x, train_y)`

plotting LR.

`plt.scatter(train.ES, train.CO2)`

`plt.plot(train_x, regr.coef_[0][0]*train_x + regr.intercept_[0], '-r')`

Error

`from sklearn.metrics import r2_score`

`test_x = np.asarray(test[['ES']])`

`test_y = np.asarray(test[['CO2']])`

`test_y_ = regr.predict(test_x)`

`print("MAE: %.2f" % np.mean(np.absolute(test_y_ - test_y)))`

`print("MSE: %.2f" % np.mean((test_y_ - test_y)**2))`

`print("R2-score: %.2f" % r2_score(test_y, test_y_))`

Simple LR

$$\text{Error} = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

$$\hat{y} = \theta_0 + \theta_1 x$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{it should be minimum.}$$

$$\theta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x} \quad \left[\begin{array}{l} \bar{x} = \frac{\sum \text{of every } x}{\text{Total}} \\ \bar{y} = \frac{\sum \text{of every } y}{\text{Total}} \end{array} \right]$$

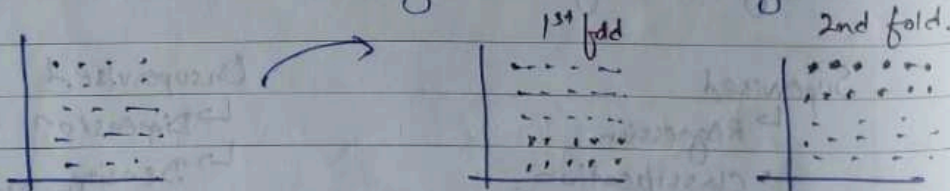
It is fast, no parameter tuning, easy to understand

Model evaluation approaches.

- Train and Test on the same dataset → may result into overfit
high training accuracy
Low out of sample accuracy
- Train/Test split → mutually exclusive, more accurate
But Highly dependent on which datasets the data is trained and tested

- K-fold cross-validation [It solves all issues]

It split some data set into various folds means splitting a table into training and testing set and where each testing set varies from one another and train with the rest of the data set and calculate the accuracy, and at last we average all the accuracy.



MAE: Mean Average Error

$$\frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

MSE: Mean Square Error

$$\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2$$

RMSE: Root Mean Square Error

$$\sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

Relative Absolute Error RAE = $\frac{\sum_{j=1}^n |y_j - \hat{y}_j|}{\sum_{j=1}^n |y_j - \bar{y}|}$

$$\sum_{j=1}^n |y_j - \bar{y}|$$