

Cheat Sheet :

importing the dataset

```
dataset = pd.read_csv(' .csv' ) [dataset.iloc[:, 2:4].values]
```

Missing Data values

from sklearn.preprocessing import Imputer

```
imputer = Imputer(missing_values = 'NaN', strategy = 'mean',  
                  axis = 0)
```

↳ Impute column, 1. means row

```
imputer = imputer.fit(x)
```

```
x = imputer.transform(x)
```

Categorical Data

from sklearn.preprocessing import LabelEncoder, OneHotEncoder

```
x labelencoder_x = LabelEncoder()
```

```
x = labelencoder_x.fit_transform(x) [x[:, 0]]
```

```
onehotencoder = OneHotEncoder(categorical_features = [0])
```

```
x = onehotencoder.fit_transform(x).toarray()
```

Splitting the dataset into TrainingSet and TestSet :

from sklearn.cross_validation import train_test_split

x_train, x_test, y_train, y_test = train_test_split

(x, y, test_size = 0.2, random_state = 0)

Feature Scaling : Two ways < Standardisation Normalisation

from sklearn.preprocessing import StandardScaler

sc_x = StandardScaler()

x_train = sc_x.fit_transform(x_train)

x_test = sc_x.transform(x_test)

Linear Regression :

```
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, Y_train)
y_pred = regressor.predict(X_test)
```

Multiple Linear Regression :

encode the data if string exists (Label, OneHot-Encoding)
 # Avoiding Dummy Variable Trap

$$X = X[:, 1:]$$

Fitting Multiple Linear Regression

```
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, Y_train)
y_pred = regressor.predict(X_test)
```

Matrix of Features ≥ 2 :

$[:, 1]$ → create vector

But our Matrix of Features should always be in matrices, not in vector. So $[:, 1:2]$ → upper bound excluded.

Polynomial Regression :

from sklearn.preprocessing import PolynomialFeatures

poly-reg = PolynomialFeatures(degree = 2) → creates a column of squared values.

x-poly = poly-reg.fit_transform(~~x-test~~)

lin-reg-2 = LinearRegression()

lin-reg-2.fit(x-poly, y)

Visualising the results :

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

plt.plot(x, lin-reg-2.predict(poly-reg.fit_transform(x)), color = 'blue')

plt.title(), plt.xlabel(), plt.ylabel(), plt.show()

*** Instead of writing x-poly, we wrote the whole poly-reg.fit-trans(x) because our ~~poly-reg~~ may x-poly may test-set of x i.e x-test.

Clear Results : [polynomial, SVR, Decision Tree, Random Forest]

x-grid = np.arange(min(x), max(x), 0.1)

x-grid = ~~np~~ x-grid.reshape((len(x-grid), 1))

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

plt.plot(x-grid, lin-reg-2.predict(poly-reg.fit_transform(x-grid)),
color = 'blue')

plt.show()

Visualising Test Set Results : [Linear, Multiple linear]

```
plt.scatter(x_test, y_test, color='red')
```

```
plt.plot(x_train, regressor.predict(x_train), color='blue')
```

Visualising Training Set Results : [Linear, Multiple linear]

```
plt.scatter(x_train, y_train, color='red')
```

```
plt.plot(x_train, regressor.predict(x_train), color='blue')
```

SVR Model : || sc-x ~ StandardScaler()

```
from sklearn.svm import SVR
```

```
regressor = SVR(kernel='rbf')
```

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

```
y_pred = regressor.predict(sc_x.transform([6.5])) — X
```

Because the input to predict must be an array.

```
y_pred = regressor.predict(sc_x.transform(np.array([6.5])))
```

↓
Why 2nd bracket?

If 1 bracket — vector

If 2 brackets — Array

```
y_pred = sc_y.inverse_transform(reg.predict(sc_x.transform([6.5])))
```

↓
Because we are predicting y value. Inverse Transform to inverse scaled value.

Decision Tree Regression:

```
from sklearn.tree import DecisionTreeRegressor  
regressor = DecisionTreeRegressor(random_state = 0)  
regressor.fit(x, y)
```

Random Forest Regression: Ensemble learning (or) Multiple Deep Forest Trees.

```
from sklearn.ensemble import RandomForestRegressor  
regressor = RandomForestRegressor(n_estimators = 10, random_state = 0)  
regressor.fit(x, y)  
y_pred = regressor.predict(6.5)
```

CLASSIFICATION

(4)

Logistic Regression :

```
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)
```

Confusion Matrix :

```
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
>> cm
```

K - Nearest Neighbours :

```
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 5,
                                metric = 'minkowski', p = 2)
metric = 'minkowski' → To choose Euclidean distance
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)
```

Support Vector Machine (SVM):

```
-from sklearn.svm import SVC  
classifier = SVC(kernel = 'linear', random_state = 0)  
classifier.fit(x_train, y_train)  
y_predict = classifier.predict(x_test)
```

Kernel SVM :

```
-from sklearn.svm import SVC  
classifier = SVC(kernel = 'rbf', random_state = 0)  
classifier.fit(x_train, y_train)  
y_pred = classifier.predict(x_test)
```

Naive - Bayes Theorem : $P(A/B) = \frac{P(B/A) * P(A)}{P(B)}$

*** See concept in the notebook (mandatory).

```
-from sklearn.naive-bayes import GaussianNB  
classifier = GaussianNB()  
classifier.fit(x_train, y_train)  
y_pred = classifier.predict(x_test)
```


Decision Tree Classification:

```

from sklearn.tree import DecisionTreeClassifier
classifier = DecisionTreeClassifier(criterion = 'entropy',
                                   random_state = 0)
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)

```

Random Forest Classification:

```

from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 10,
                                   criterion = 'entropy', random_state = 0)
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)

```

CLUSTERING

K-Means clustering: from sklearn.cluster import KMeans

#using elbow method

wcss = []

for i in range(1, 11):

kmeans = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300,
n_init = 10, random_state = 0)

```
kmeans.fit(x)
```

```
wcss.append(kmeans.inertia_)
```

method to calculate wcss.

```
plt.plot(range(1, 11), wcss)
```

```
plt.title(' '), plt.xlabel(' '), plt.ylabel(' '), plt.show()
```

Applying k-means to dataset

```
kmeans = KMeans(n_clusters = 5, init = 'k-means++', max_iter = 300,  
                 n_init = 10, random_state = 0)
```

```
y_kmeans = kmeans.fit_predict(x)
```

Hierarchical Clustering : Agglomerative Approach

```
import scipy.cluster.hierarchy as sch
```

```
dendrogram = sch.dendrogram(sch.linkage(dataset, method = 'ward'))
```

→ Agglomerative
Algorithm

ward method in
Algorithm

```
plt.plot(' ')
```

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

```
:
```

```
plt.show(' ')
```

Optimum clusters : 5

Fitting Hierarchical Clustering to the dataset

```
from sklearn.cluster import AgglomerativeClustering
```

```
hc = AgglomerativeClustering(n_clusters = 5, affinity = 'euclidean',
```

```
                             linkage = 'ward')  
y_hc = hc.fit_predict(x)
```

Association Rule Learning

6

Apriori : SEE concept in notebook (Mandatory)

In Python, Apriori expects input as lists of lists.

```
dataset = pd.read_csv('')
```

```
transactions = []
```

```
for i in range(0, len(dataset)7501):
```

```
    transactions.append([str(dataset.values[i, j])
```

```
                        for j in range(0, 20)])
```

↓
To make list of lists.

```
# training apriori with dataset  
from apyori import apriori
```

```
rules = apriori(transactions, min_support = 0.003,
```

```
                min_confidence = 0.2, min_lift = 3, min_length = 2)
```

```
# visualising results  
results = list(rules)
```