

HOW TO DEAL WITH A ML PROJECT

① IMPORT LIBRARIES

5 import pandas as pd
" " numpy as np

```
import matplotlib.pyplot as plt ]- For visualisation  
"      seaborn as sns
```

10 import warning
warnings.filterwarnings('ignore') } To remove warnings

from sklearn.decomposition import PCA } To convert
15 " sklearn.manifold import TSNE } high dimensional
import matplotlib as plt data to low
dimensional

from sklearn.model_selection import RandomizedSearchCV
For cross validation and
improve accuracy

```
from sklearn.linear_model import LogisticRegression  
" " , svm import SVC  
" " , tree import DecisionTreeClassifier  
" " , ensemble import RandomForestClassifier  
  
25  
" " , ensemble import RandomForestClassifier  
  
ML algo.
```

30 from sklearn.metrics import accuracy_score
To check accuracy

② LOAD DATA

train = pd.read_csv ('train.csv')
 test = " " " ("test.csv")

5 → Check data

↳ train.head()

↳ test.head()

↳ train.info()

↳ train.shape

↳ train.describe()

③ DATA PRE PROCESSING

15 → Checking null values

* train.dropna()

* train.isnull().sum()

* train.drop(index, inplace=True)

OR

train.drop(col.name, axis=0/1)

* train[col] = train[col].fillna('To fill')

* train[col] = np.where(np.isnan, tofill, train[col])

25 → Checking duplicates

* sum (train.duplicated())

* train.drop_duplicates (keep='first', inplace
= True)

④ ANALYSING DATA

1 DATA WRANGLING

* Labeling of data

5 eg Sex: M | F

↓
0 1 1

use

10 from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
train['Sex'] = le.fit_transform(train['Sex'])

* Use graphs for analysing

→ Check distribution (should be normal)

15 sns.distplot (x=col, data=train, kde=True)

20 * To plot multiple graphs in a single

features = ['Col 1', 'Col 2', 'Col 3', ...]

for i in features:

25 sns.distplot (x=i, data=train, kde=True)

OR

features = ['Col 1', 'Col 2', 'Col 3', ...]

l = list (enumerate(features)) # It will give index

30 * for i in l:
plt.subplot(3, 3, i[0]+1)

sns.scatterplot (x=i[1], y='Age', data=train)

Here $i[0]$ = index
 $i[1]$ = features

DATE : ___/___/
PAGE :

* To change figure size use

`plt.figure(figsize=(13, 12))`

5

* Use Boxplot to check outliers

`sns.boxplot(x=col, data=train)`

10



15

`sns.boxplot(data=train)`

* Use z-score or quantile(stand. dev.) to remove outliers

$$z\text{-score} = \frac{(df.\text{columnname} - df.\text{col.mean})}{df.\text{col.std}}$$

std → standard deviation

then remove all z-score having value less than -3
more than +3

25

OR

`q = train['col'].quantile(0.97) #change accordingly
train = train[train['col'] < q]`

30

* Check correlation of parameters

sns.heatmap (train.corr(), annot=True)

↳ Check values w.r.t. output column

5

* cd

10

⑤ FEATURE ENGINEERING

↳ Scaling

from sklearn.preprocessing import
StandardScaler

15

OR

min. MaxScaler

x = pd.DataFrame (MinMaxScaler).fit_transform(x)

20

↳ Correlation

↳ heatmap

↳ Using PCA / TSNE

25

x = train.drop(['Output Col'], axis=1)

pca = PCA (n_components=2, random_state=0)
• fit_transform(x)

30

sns.scatterplot (x=pca[:,0], y=pca[:,1])

hue = train['Outcome'])

• Fly with TSNE

⑥ ML Modelling

* Splitting of data

5
 $x_{train} = train.drop(['Output col'], axis=1)$

$y_{train} = train['Output col']$

10
 $x_{test} = test.drop(['output (if any')], axis=1)$

OR .

remove any unwanted columns

↳ check by correlations

15 * Apply models using hyperparameters, Across Validation

↳

20 * For logistic regression

parameters = {'max_iter': [100, 200, 300]}

lr = LogisticRegression()

lr_rs = RandomizedSearchCV(lr, param_dist

= parameters,

cv=5, n_jobs=-1

= 42)

lr_rs.fit(x_train, y_train)

30
 $y_{pred} = lr_rs.predict(x_{test})$

* For Kernel SVM

parameters = {
 'kernel': ['linear', 'rbf', 'poly', 'sigmoid'],
 5 'C': [100, 5],
 }
 # same

* for Decision tree

10 parameters = {
 'max_depth': np.arange(2, 10, 2),
 }
 # same

* For Randomforest

15 parameters = {
 'n_estimators': np.arange(20, 100, 10),
 'max_depth': np.arange(2, 16, 2),
 }
 # same

OR To increase accuracy use Catboost model

train-data = cb.Pool(x_train, y_train)
 25 model = cb.CatBoostRegressor(loss function='RMSE')

grid search

30 grid = { 'iterations': [100, 150],
 'learning_rate': [0.03, 0.1],
 'depth': [2, 4, 6, 8]}

model = grid - Search(grid, train-data)

pred = model.predict(x-test)

part (RMSB) (A.O) ✓ good model

Extra Info

5

* Date Time application

eg Date → 9/26/2009

↓
splitting
}

df['Year'] = df['Date'].apply(lambda
x: x.split('/')[2])

use for month, day

15

* also you can use

df['Year'] = pd.to_datetime(df['Date']).dt.year

use for time, hour, min, sec,
day, month etc

25

30