

# machine learning cheat sheet

## ① Data analysis

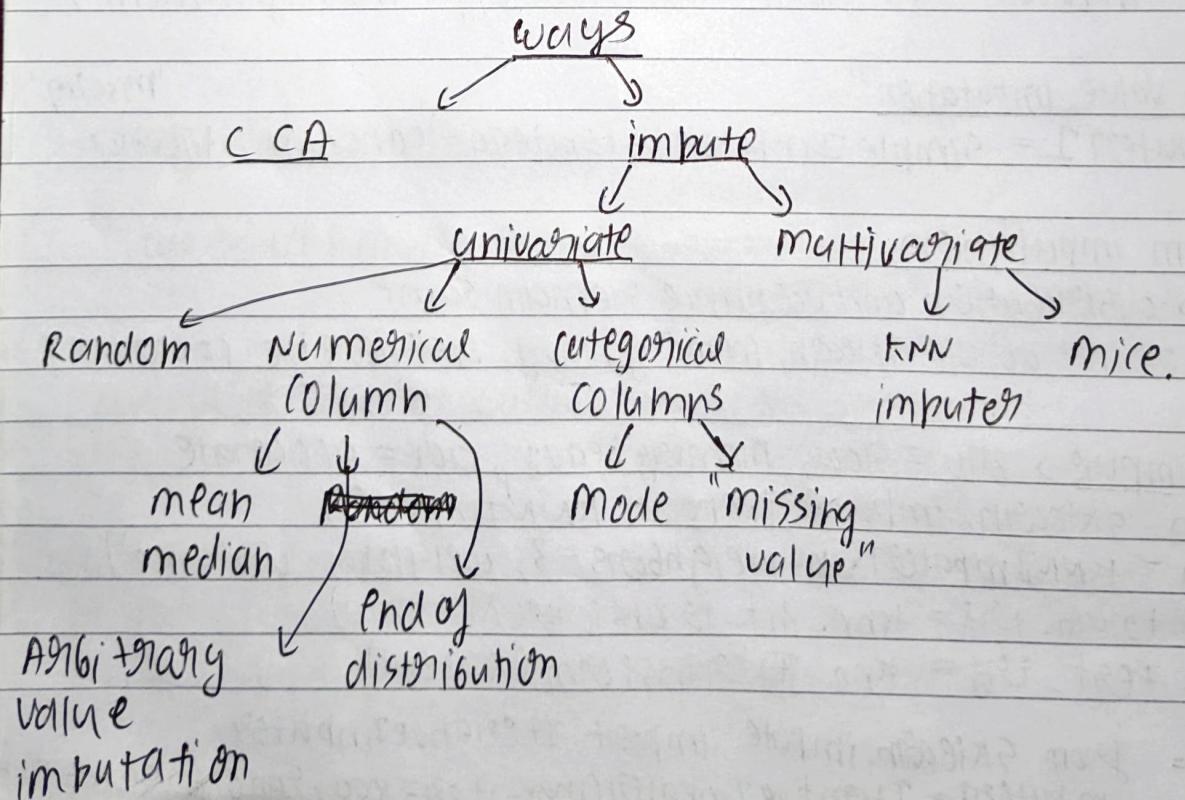
(i) from ydata.profiling import ProfileReport, p=ProfileReport(df),  
p.to\_notebook\_iframe()

(ii) plt.scatter(c1, c2) (iii) lineplot = plt.plot(c1, 12), (iv) plt.bar(c1, 12) (v)  
sns.catplot(c1, 12, data, kind='bar') (vi) kde = sns.kdeplot  
plot(c1) (vii) stats.probplot(data, dist="norm", plot=plt)  
↳ [from scipy import stats]

## ② Feature Engineering

steps = Feature selection → missing value imputation → outlier detection → Feature scaling → encoding of categorical column → encoding of numerical column → Transformation of columns → PCA

## ③ Missing value imputation



$$\rightarrow \pi_{\text{train}}[C_1 + \text{try}] = \pi_{\text{train}}[C_1]$$

→ Now apply transformation on  $C_1 + \text{try}$

we apply CCA and univariate when data missing is less than 5% and we have to do analysis after CCA and univariate imputation on the column.

CCA  $\rightarrow$  df. isnull().mean() \* 100  
 $\hookrightarrow$  new\_df = df[cols].dropna()

$\rightarrow$  production code / model  
 can not handle missing value

mean median

div → change in distribution, correlation + addition of outliers  
 imputer = SimpleImputer(strategy='mean/median')  
 from sklearn.impute import SimpleImputer

Arbitrary value imputation ~~not good for production model as it can't handle missing~~

not used very much

imputer1 = SimpleImputer(strategy='constant', fill\_value=999)

Mode

div → change in distribution, correlation + addition of outliers  
 imputer1 = SimpleImputer(strategy='most\_frequent')

"missing value imputation"

imputer1 = SimpleImputer(strategy='constant', fill\_value=missing)

Random imputation

adv → distribution and variance remain same

div → can not do with sklearn, memory heavy, distribute the covariance

KNN impute → div = slow, memory heavy, adv = accurate

from sklearn.impute import KNNImputer

knn = KNNImputer(n\_neighbors=3, weight='distance')

$X_{\text{train}}_{\text{try}} = \text{knn}.fit\_transform(X_{\text{train}})$

$X_{\text{test}}_{\text{try}} = \text{knn}.transform(X_{\text{test}})$

MICE

from sklearn.impute import IterativeImputer

imputer2 = IterativeImputer(max\_iter=1000, random\_state=42)

## Analysis Part

### Numerical columns:

#### Variance

`print(x -> grain['C1'].var())`

~~`print((x -> grain['C1']).var())`~~

`print((x -> grain['C1'] -> var))`

Same goes for test

#### Distribution

`x -> grain['C1'].plot(kind='kde', qx=ax)`

`x -> grain['C1'].plot(kind='kde', qx=ax, color='green')`

#### COV - covariance

`x -> grain -> cov()`

#### CO - correlation

`x -> grain -> corr()`

#### Outliers

`x -> grain[['C1', 'C1'].boxplot()])`

### Categorical columns:

`temp = pd.concat([`

#### Change in categories

`C`

`x -> grain['C1'].value_counts().len()`

`x -> grain['C1'].diagonal()`

`x -> grain['C1'].value_counts().len(x -> grain)], axis=1)`

`temp.columns = ['original', 'imputed']`

`temp`

### ④ Outlier detection

ways to handle outlier

Trimming

Capping

Other

Treating them as missing

Discretization

ways to detect outlier → Z score treatment

→ Interquartile proximity rule

→ Percentile based approach

#### Z score treatment

→ upper limit  $\Rightarrow \mu + 3\sigma$  < value, lower limit  $\Rightarrow \mu - 3\sigma$  > value

→ used for normally distributed column

$$\text{upper\_limit} = \text{df}['c1'].mean() + 3 * \text{df}['c1'].std()$$

$$\text{lower\_limit} = \text{df}['c1'].mean() - 3 * \text{df}['c1'].std()$$

## Finding outliers

$\text{df}[(\text{df}['c1']) > \text{upper\_limit}) | (\text{df}['c1']) < \text{lower\_limit})]$

## Trimming

new df =  $\text{df}[(\text{df}['c1']) < \text{upper\_limit}) \& (\text{df}['c1']) > \text{lower\_limit}]$

## Capping

$\text{df}['c1'] = \text{np.where}(\text{df}['c1']) > \text{upper\_limit},$   
 $\text{upper\_limit}, \text{np.where}(\text{df}['c1']) < \text{lower\_limit}, \text{lower\_limit},$   
 $\text{df}['c1'])$

## Inter Quartile Proximity Rule (Box Plot)

$$Q1 = 25\text{ Percentile} \quad Q3 = 75\text{ Percentile} \quad IQR = Q3 - Q1$$

$$\text{upper\_limit} = Q3 + 1.5 * IQR \geq \text{Value}$$

$$\text{lower\_limit} = Q1 - 1.5 * IQR \leq \text{Value}$$

## Finding outliers

## Normal ways

$$Q1 = \text{df}['c1'].quantile(0.25) \quad | \quad IQR = Q3 - Q1$$

$$Q3 = \text{df}['c1'].quantile(0.75)$$

$$\text{upper\_limit} = Q3 + 1.5 * IQR$$

$$\text{lower\_limit} = Q1 - 1.5 * IQR$$

## Finding outliers

→ (a) ~~outlier detection~~ same as above  
(b) Sns - box plot ( $\text{df}['c1']$ )

## Trimming + capping // same as above

Percentile based approach → here we do not use box plot  
In this method we assume percentile like  $\text{max}=99, \text{min}=1$

$$\text{upper\_limit} = \text{df}['c1'].quantile(0.99)$$

$$\text{lower\_limit} = \text{df}['c1'].quantile(0.01)$$

Finding outliers  
Trimming  
Capping  
same

## ⑤ Feature Scaling

### (i) Standardization

(ii) Normalization → minmax scaling  
Robust scaling → mean normalization  
Max Abs Scaling

### (i) Standardization

→ outliers remain same

from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

### (ii) Min Max Scaling

→ used when you know min value and max value

→ outliers squeeze, Note → cause change in data

from sklearn.preprocessing import MinMaxScaler,  
MaxAbsScaler, RobustScaler

### (iii) Mean Normalization → not supported by sklearn,

### (iv) RobustScaling → used if data has a lot of outliers

### (v) MaxAbsScaling → used if data is sparse (a lot of zeros)

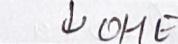
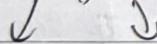
## ⑥ Encoding of categorical columns

No. of unique categories = df['c1'].nunique()

No. of occurrence of unique categories = df['c1'].value\_counts()

### Types

Ordinal categorical data ↴ Nominal categorical data



⊗ ~~ordinal~~

~~nominal~~ ⊗

ordinal encoder

label encoder

### Ordinal Encoder - from sklearn.preprocessing import

OrdinalEncoder, LabelEncoder, OneHotEncoder

oe=OrdinalEncoder(categories=[['Poor', 'Avg', 'Good'], ['D', 'C', 'B', 'A']])

label encoder → le = LabelEncoder() // Here we do not mention the categories

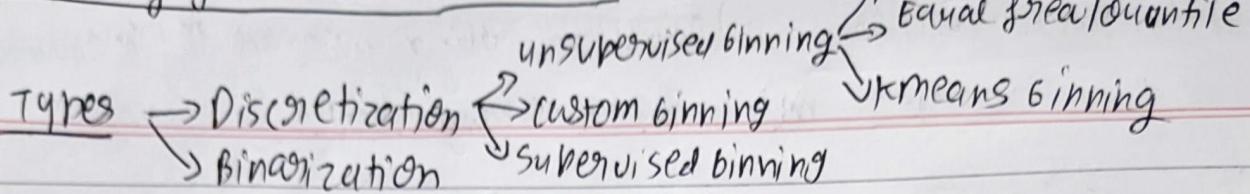
OHE → ohe = OneHotEncoder(drop='first', sparse=False, dtype=np.int32)

Note → if so many categories then ⇒ threshold = 100,

counts = df['c1'].value\_counts(), gtepl = counts[counts <= threshold].index

pd.get\_dummies(df['~~c1~~'].replace(gtepl, 'uncommon'))

## ⑦ Encoding of numerical columns



Equal width → You have to decide number of bins → Create equal size bin  
 $(\text{max value} - \text{min value}) / \text{Number of bins}$

handle outliers and no change in spread of data

→ More we decide bins

Equal freq → handle outliers and normally dist the spread of data

K-mean binning → when data is present in clusters

from sklearn.preprocessing import KBinsDiscretizer, Binarizer

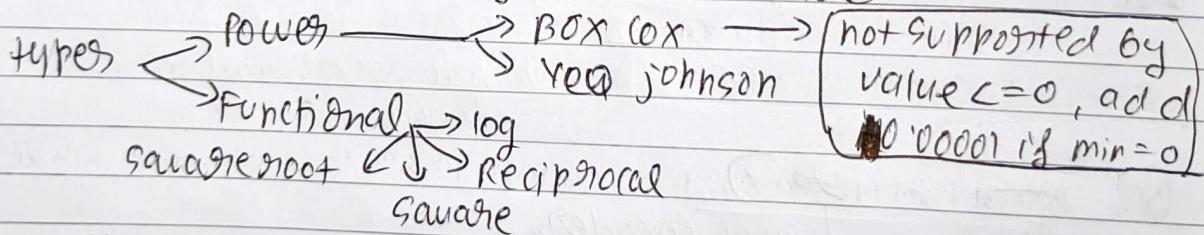
a = KBinsDiscretizer(n\_bins=15, encode='ordinal/onehot/onehot-dense', strategy='uniform/quantile/kmeans')

Binarization → 0 or 1 Based on threshold ⇒

a = Binarizer(copy=False, threshold=4.0)

## ⑧ Transformation of numerical columns

if Q-Q plot = 45% = Normally distributed + tde



from sklearn.preprocessing import PowerTransformer, FunctionTransformer

pt = PowerTransformer() → 'Yeo-Johnson'

pt = PowerTransformer(method='boxcox')

(Q-Q Trans form) → Right skewed (does not work on -ve data)

Reciprocal transform → square root transform → left skew

square root transform → not same as Q-Q trans

Best → Yeo Johnson

## ⑨ PCA

from sklearn.decomposition import PCA

pca = PCA(n\_components=5)

⑩ Column Transformer

from sklearn.compose import ColumnTransformer

$$+g_{\text{tf}} = \text{ColumnTransformer}(\text{transformers}=[$$

$('t_{\text{tf1}}')$	$\text{, } [6, 8])$
$('t_{\text{tf2}}')$	$\text{, } [3])], \text{ remainder} = \text{'passthrough'}$

$$\begin{array}{l} T1 = FP \\ T2 = FN \end{array}$$

## ⑪ Pipeline

from sklearn.pipeline import Pipeline

pipe = Pipeline([('+g\_{\text{tf1}}', +g\_{\text{tf1}}), ('+g\_{\text{tf2}}', +g\_{\text{tf2}}) ... ()])

pipe.fit(X-train, Y-train)

$(\text{cv} = \text{cross-validation}, \text{int})$

binary	macro = nominal
multi	weighted = ordinal
	precision = type
	recall = type

GridSearchCV

params = { 't\_{\text{tf2}}--max-depth': [None, 3] }  $\beta_1 = \text{type1} + 2$

from sklearn.model\_selection import GridSearchCV

grid = GridSearchCV(pipe, params, cv=5, scoring='accuracy')

/precision/precision-macro/precision-micro/precision-weighted/

recall/recall-macro/recall-micro/recall-weighted/ $f_1/f_1-\cancel{\text{macro}}/f_1-\text{micro}/f_1-\text{weighted}/\text{roc-auc}/\text{roc-auc-avg}/$

$\text{roc-auc-ovo}$ ,  $\text{roc-auc-ovo}$ -weighted,  $\text{roc-auc-ovo}$ -weighted)

classification metrics (neg-mean-absolute-error/neg-mean-squared-error/neg-root-mean-squared-error/ $R^2$  / explained-variance)

Gross validation

grid.best\_params\_

grid.best\_score\_

from sklearn.model\_selection import cross\_val\_score

cross\_val\_score(pipe, X-train, Y-train, cv=5, scoring='').mean()

## Linear Regression

from sklearn.linear\_model import LinearRegression

lr = LinearRegression(), lr.fit(X-train, Y-train), lr.predict(X-test)

## Polynomial

from sklearn.preprocessing import PolynomialFeatures

poly = PolynomialFeatures(degree=2, include=True)

$X\_train\_trans = \text{poly}.fit\_transform(X\_train)$

$X\_test\_trans = \text{poly}.transform(X\_test)$

→ Apply linear regression

Gross validate → for

training and test data

→ condition  $\leftarrow$  overfitting

→ then use Regularization

## Ridge Regularization (L2)

\* Apply if all columns are ~~not~~ imp

from sklearn.linear\_model import Ridge

## Lasso Regularization (L1)

\* Apply if all columns are not imp

from sklearn.linear\_model import Lasso

## Elastic Net Regularization

\* Apply if you know all columns are imp or not imp

from sklearn.linear\_model import ElasticNet

## Logistic Regression

from sklearn.linear\_model import LogisticRegression

for binary class =  $\text{log} = \text{LogisticRegression}()$

for multi class =  $\text{log} = \text{LogisticRegression(multi\_class='ovr')}$

Polynomial  $\rightarrow$  same as linear regression

Hyperparameters  $\rightarrow$  penalty = 'l1/l2/elastinet'  $\rightarrow$  Default = l2  
range = [0.01 ... 100] = C2 = fint small value  $\rightarrow$  high regularization  
solver = 5 different default = 'liblinear'

## Decision Tree

from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor

## Hyperparameters

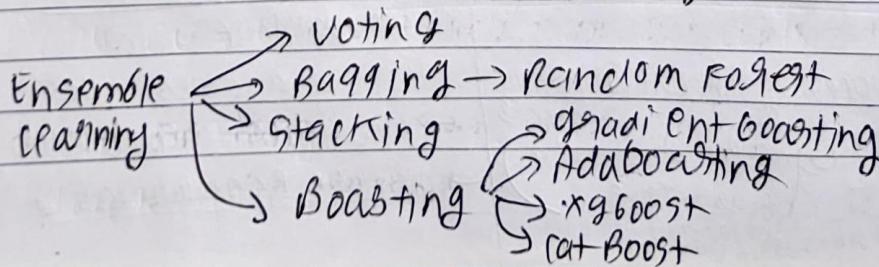
Regression  $\rightarrow$  criterion = squared\_error, Poisson, absolute\_error, mse

Classification  $\rightarrow$  criterion = gini, entropy, log\_loss

max\_depth = 10,

min\_samples\_split = jnt default = 2

## Ensemble Learning



## Voting (use different algo)

classification → Hard  
Regression → Soft

from sklearn.ensemble import VotingClassifier, VotingRegressor

lr = LinearRegression()

dt = DecisionTreeRegressor(Classifier())

svr = SVR()

estimators = [(lr, lr), (dt, dt), (svr, svr)]

Hard ⇒ vc = VotingClassifier(estimators=estimators, voting='hard')  
Soft

⇒ vr = VotingRegressor(estimators)

Bagging [use only 1 algo] DecisionTree

Types → Bagging = row sampling + replacement

→ Parting = row sampling + hot replacement

→ Random Subspace = column Sampling + with or without replacement

→ Random Patches = column + row Sampling + with or without replacement

from sklearn.ensemble import BaggingClassifier, BaggingRegressor

Bagging → bag = BaggingClassifier(

base\_estimator=DT(),

n\_estimators=500

max\_samples=0.5 (each model get 0.5 sized)

bootstrap=True

random\_state=42)

Parting → bag = BaggingClassifier(base\_estimator=n\_estimators=500,  
max\_samples=0.5, bootstrap=False, random\_state=42)

Random Subspaces → bag = BaggingClassifier(base\_estimator=n\_estimators=500,

max\_samples=1, bootstrap=False, max\_features=0.5,  
bootstrap\_features=True, random\_state=42)

Random Patches → bag = BaggingClassifier(base\_estimator=n\_estimators=500,

max\_samples=0.5, bootstrap=True, max\_features=0.5,

bootstrap\_features=True, random\_state=42)

## Random Forest

### Bagging + DT

from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor

Hyperparameters → n\_estimators, max\_features, ~~max\_depth~~,  
 bootstrap = True, max\_samples, random\_state,  
 oob\_score = True + inf. oob\_score

## AdaBoost

from sklearn.ensemble import AdaBoostClassifier, AdaBoostRegressor

Hyperparameters → ~~n\_estimators~~ = DecisionTree, n\_estimators  
 algorithm = SAMME, SAMME-R, learning\_rate =  
 generally less than 1,

## K-mean clustering

from sklearn.cluster import KMeans  
 wcss []

for i in range(1, 11)

Km = KMeans(n\_clusters = i)  
 Km.fit\_predict(X)

wcss.append(Km.inertia\_)

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

## Gradient Boosting

from sklearn.ensemble import GradientBoostingRegressor  
 GradientBoostingClassifier

Hyperparameters = n\_estimators, learning\_rate = 0.1, max\_depth,  
 random\_state = 42

## Stacking and Blending Ensemble use diff API

Merge CV = Cross-validation act as  $K = K\text{-fold mean}$

from ~~stacking~~ sklearn.ensemble import StackingClassifier,  
 StackingRegressor

```
base_models = [  
    ('91', LinearRegression()),  
    ('98', RandomForestRegressor(n_estimators=10, random_state=42)),  
    ('96', GradientBoostingRegressor(n_estimators=10, random_state=42))  
]
```

stack = StackingRegressor( estimators=base\_models,  
final\_estimator=LinearRegression(), cv=10)

### Hierarchical clustering

import scipy.cluster.hierarchy as shc  
dend = shc.dendrogram(shc.linkage(data, method='ward'))

from sklearn.cluster import AgglomerativeClustering  
cluster = AgglomerativeClustering(n\_clusters=5, affinity='euclidean', linkage='ward')

(x-test, y-pred)

### K Nearest Neighbors

choosing the best model for i in range(1,16):

knn = KNeighborsClassifier(  
n\_neighbors=i)

knn.fit(x-train, y-train)

y-pred = knn.predict(x-test)

scores.append(accuracy\_score)

~~import mat~~

plt.plot(range(1,16), scores) // choose ~~not~~ the k value  
with highest accuracy

Now make model → with that k value

### SVM Naive Bayes

SVM Regressor → not used very much

from sklearn.svm import SVC

Hyperparameters → C = regularization parameter, kernel = 'rbf', 'linear' ... , degree, ~~gamma~~=auto, random\_state=42

import xgboost as xgb

### XGBoost

mean 2 class with  
y ↴

① xgb\_clf = xgb.XGBClassifier(max\_depth=6, learning\_rate=0.1, n\_estimators=120, objective='multi-class' num\_class=2)

② xgb\_regr = xgb.XGBRegressor() → These hyperparameters

### Regression Objective

- 'neg: squarederror'
- 'neg: squaredlogerror'
- 'neg: logistic'
- 'neg: pseudohubererror'

### Classification Objective

- 'binary: logistic'
- 'binary: logitraw'
- 'multi: softmax'
- 'multi: softprob'

### CatBoost

import catboost as cb

① cat\_clf = cb.CatBoostClassifier(iterations=100, depth=6, learning\_rate=0.1, loss\_function='MultiClass')

② cat\_regr = cb.CatBoostRegressor()

### Regression loss function

- RMSE, MAE, Quantile, LogLoss, Quantile, Poisson

### Classification loss function

- LogLoss, CrossEntropy, Multiclass, MulticlassOneVsAll

### LightGBM