

1) Handling Missing Values -

`df.columns.fillna()`

fill with

mean

1) mean/median

Data is numerical.
mean → Normally distributed data.
median → Skewed data.

2) mode

Categorical variable.

`df.groupby('occupation')[['Income']].transform(lambda x: x.mean())`

2) Handle outliers -

Quartiles → The parts of data that each contain 25% of the data, low to high when sorted in ascending order.
Data → Q_1, Q_2, Q_3, Q_4 .

IQR → Inter quartile range → $Q_3 - Q_1$.

If (point > $Q_3 + 1.5 \times \text{IQR}$) outlier.

if (point < $Q_1 - 1.5 \times \text{IQR}$) outlier.

3) Cap data → Replace outliers with Q_1 & Q_3 respectively.

4) log(data) → Compares the data & brings the outliers closer.

2) Scaling The Data -

3) Standardization → Converting the values to Z-scores.

- Good for uniformly distributed data.
- less sensitive to outliers.
- use in models that are sensitive to scale.

$$x_{\text{std}} = \frac{(x - \bar{x})}{\sigma}$$

→ distance away from the mean.

→ 1 unit away from the mean.

→ # of units away from mean.

6) Min-Max Normalization \rightarrow Convert the values b/w 0 & 1.

- Very sensitive to outliers.
- used in NN, Image transformation, feature engineering.
 - \downarrow Convert pixel values from [0-255] \rightarrow [0, 1]

$$\frac{x}{x_{\text{min-max}}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

from sklearn.preprocessing import StandardScaler,
MinMaxScaler

Scaler = StandardScaler()
Merger = MinMaxScaler()

$x_{\text{std}} = \text{Scaler.fit_transform}(x)$

$x_{\text{min-max}} = \text{Merger.fit_transform}(x)$

3) Data Encoding \rightarrow Convert categorical values into numeric values.



a) Label Encoding \rightarrow use when the categories have no inherent order.

from sklearn.preprocessing import LabelEncoder

encoder = LabelEncoder()

data = [e, w, ev, s, n]

encoder.fit_transform(data)

e 4
w 4
ev 3
s 4
n 1 | \rightarrow data = {3, 4, 4, 2, 1}

\downarrow data after encoding.

6) ordinal Encoding → when the categories have inherent order b/w them.

on-2
data = [hs, b, m, phd, hs]
data_order = [hs, bs, ms, phd]

from sklearn.preprocessing import OrdinalEncoder

en = OrdinalEncoder(categories=[data_order])

d_encoded = en.fit_transform(data)

d_encoded → [1, 2, 3, 4, 1]
hs bs ms phd hs

7) One-hot Encoding → convert categorical variables to binary columns where each represents the occurrence of a variable.

on-hot-encoder pd.get_dummies() } → Both do the same thing but with inherent differences.

8) One-Hot Encoder → - use with sklearn if ML pipelines.

from sklearn.preprocessing import OneHotEncoder

encoder = OneHotEncoder(sparse=False)

n = encoder.fit_transform(data_categorical_columns)

(→ returns a numpy matrix.)

9) pd.get_dummies() → use when performing exploratory data analysis.
- returns a pandas data frame.

5) Model Evaluation Strategy → Choosing the parts of dataset that will train & test the model for comparison with other models.

10) Train-Test split → Split the data set into two parts & use each part to train & test the dataset separately.

- high variance → high randomness.

- Split data may not be balanced.

- good for quick & easy iterations.

- 6) Cross-Validation - Data is split into multiple folds, for a certain # of repetitions a different fold is selected as validation set & rest other as train set.
- Repetition reduces randomness.
 - Dataset might not be balanced.
 - General case for model selection.

- 7) K-fold Cross Validation - Data is divided into K sets.

- Repeated K times, for each K
 - a validation set is selected & model trained on rest of the data.
 - Folds can be unbalanced.

- 8) Repeated K-fold - Repeat K-fold validation multiple times.

- Reduced variance.
- Imbalanced splits.
- Increases model robustness.

- 9) Stratified-K-folds - Split the data into K folds & makes sure that each split is balanced. ✓

- Balanced dataset.
 - good for classification tasks.
- Each class in each split has equal # of elements.

- 10) Linear Regression - Model the relationship b/w a set of variables by fitting a straight line into them.
- one of the simplest & widely used ml algorithms.

Procedure -

$$Y = mx + c$$

X Y x^2 y^2 xy need to find these two.

$$\frac{m}{n} \frac{\Sigma y}{\Sigma x} \frac{\Sigma x^2}{n} \frac{\Sigma y^2}{n} \frac{\Sigma xy}{n}$$

$$m = \frac{n \sum xy - (\sum x \sum y)}{n \sum x^2 - (\sum x)^2}$$

\downarrow with an unit change in x, how much does y change.

$$C = \frac{\sum y - n\bar{x}}{n}$$

```

from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
model = LinearRegression()
model.fit(x, y) → Series.
y = model.predict(x-test) → DataFrame
print(mean_squared_error(y, ŷ), r2_score(y, ŷ))

```

Linear Regression Metrics

	Residue $\rightarrow y_i - \hat{y}_i$	
Median	\downarrow	Mean
or Absolute error \rightarrow		
MedianAE $\rightarrow \text{median}(y_i - \hat{y}_i)$		or Mean squared error.
- Less sensitive to outliers.		or Mean absolute error.
- Robust to variance in errors.		or Root mean squared error.
Silversides		

or Mean Absolute Deviation \rightarrow $y \times \sum_{i=1}^n |y_i - \hat{y}_i|$
 - Less sensitive to outliers.
 - All errors should be equally weighted.

6) Mean Squared Error \rightarrow

- Sensitive to outliers.
- Good for model optimization.
- Penalizes larger errors more.

$$Y_n \times \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

7) Root Mean Square Error \rightarrow

- Similar properties to Mean Squared Error.
- Brings error to y 's units.

$$Y_n \times \sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

8) $R^2 \rightarrow$

$$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

- To determine how much variance ($y_i - \hat{y}_i$) does the model explains, when compared with the average model.

→ The lower the better.
↳ The higher the better.

8) Assumptions For Linear Regression \rightarrow

1) Linearity \rightarrow There should a linear relationship between the dependent & independent variable.
if not followed →
- Incorrect predictions.

2) Independence \rightarrow The data points should be independent of each other.
- Misingding tests.

3) Homoscedasticity \rightarrow The variance in errors should be similar across all levels of the input variable.

In \rightarrow In sales data, no higher errors for high priced item compared to normal priced items.

- unreliable predictions at certain ranges.

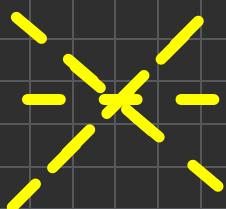
4) Normality of Residuals → The residual should be normally distributed.

$$\epsilon_i(y_i - \hat{y}_i)$$

- Incorrect confidence intervals.

5) No Multicollinearity → The variables in the independent variables set should be independent of each other.

- Less information gained from dependent variables.



Logistic Regression → Simple & computationally easy classification algorithm.

→ Classification → Task of putting/classifying an item into a certain category.

Linear regression → Continuous outcome.

- Simple & efficient.
- Widely used.
- Probabilistic.

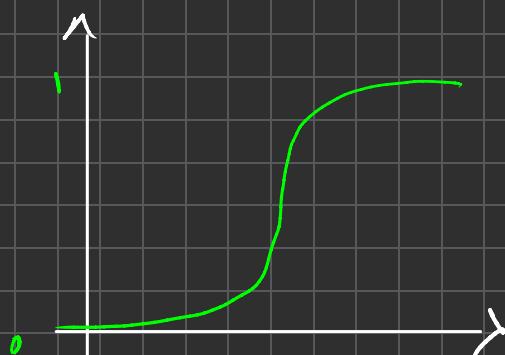
Logistic regression → Discrete outcome.

why? → True classification tasks.

what? → Let $\lambda = \sum_{i=1}^n x_i + w_i$

$$y = \frac{1}{1 + e^{-\lambda}}$$

range(y) = [0, 1] → good for interpreting as probability.



1) Confusion Matrix → A matrix used for evaluating model accuracy.

- Displays model predictions for categories.
- Shows how many predictions are correct & how many are wrong.

		Predicted	
		True	False
Actual	True	True - TP	False - FN
	False	False - FP	True - TN

TP → Correctly classify a true item.

FP → Incorrectly classify a -ve item.

FN → Incorrectly classify a true item.

TN → Correctly classify a -ve item.

→ Metrics →

1) Accuracy → $\frac{\# \text{ of correct predictions}}{\text{Total # of predictions}}$

$$= \frac{TP + TN}{TP + FP + FN + TN}$$

2) Precision → $\frac{\# \text{ of correct True predictions}}{\text{Total # of True predictions}}$

$$= \frac{TP}{TP + FP}$$

3) Recall →

$\frac{\# \text{ of correct True predictions}}{\text{Total # of true items}}$

$$= \frac{TP}{TP + FN}$$

4) F-1 Score → Balances Precision & Recall.

$$\text{F-1 Score} = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Regression Advanced -4

↳ Adjusted R² →

↳ Problem with r^2 → Adding an useless column to data may increase r^2 , thus giving inaccurate results.

going from $\hat{y} = w_0 x_0 + \dots + w_n x_n$
to

$\hat{y} = w_0 x_0 + \dots + \frac{w_n x_n}{n+1} + \frac{w_{n+1}}{n+1}$
might change r^2 value.

↳ Adjusted r^2 → Idea → Account for change in dimensionality.

↳ # of features.

$$\text{adjusted } r^2 = 1 - \frac{(1 - r^2)(n - 1)}{(n - m - 1)}$$

range ∈ $(-\infty, 1]$

n → # of rows.
 m → # of columns.

high adjusted r^2 → Model explains more variance in dependent variable.

↳ OLS Method of (statsmodel & library)

↳ ordinary Least Square → estimating model parameters by minimizing sum of square errors.

- Fisher model statistics than sklearn.

↳ offers additional features.

import statsmodel.api as sm

$x_sm = sm.add_constant(x_train)$ → adds a column for the bias co-eff.

model = sm.OLS(y_train, x_sm)

results = model.fit()

results.summary(), results.predict(x_sm)

Overfitting → Model performs well on training data but fails to generalize well on test data.

high variance & low bias.

Underfit → The data is too complex for the model & fails to perform good even on train data.

low variance & high bias.

→ ways to reduce overfitting →

1) use simpler model.

2) Regularization.

3) Early stopping & stop once the validation error starts to increase.

4) Ensemble Methods.

5) Cross-validation.

6) Drop useless features.

7) Data augmentation.

Regularization → A form/variant of regression with a penalty term to reduce overfitting.

1) Ridge ($\lambda=2$) → Loss = sum of sq of errors +
sum of sq of weights
Ridge-regression term.

- useful for limiting/minimizing the weights.

- Prevents the weights from increasing drastically.

2) Lasso ($\lambda=1$) → Loss = sum of sq errors +

sum of absolute weights.
- useful to drop useless features.

→ Lasso term

$$3) \text{ Elastic Net} \rightarrow \text{Loss} = \text{sum of } \ell_2 \text{ errors} + \lambda (\alpha(\ell_2) + (1-\alpha)\ell_1)$$

- shrinks the weights smoothly. Controls shrinkage.
- Removes less helpful variables.
- More stable than Ridge & Lasso alone.

Balanced regularization.