

Complete Machine Learning

— / —

① Linear Regression

Goal: To find the best fitting line (hyperplane) that minimizes the error b/w (\hat{y} and y)

$$y = w_0 + b \quad / \quad y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Assumptions:

1. Linearity

2. Independence of errors

3. Homoscedasticity: Constant variance of errors

4. Normality: Errors distributed normally

5. No multicollinearity.

* The best fit line will be the one that optimizes the values of m (slope) and b (intercept) so the predicted y values are as close as possible to the actual data points.

To find a best fit line: Least Squares Method

Minimise the sum of squared diff b/w actual & predicted values from the line \rightarrow residuals. (mean should be 0)

$$\{\text{SSE} / \text{RSS}\} : \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

(Residual Sum of Squares)

Ensure the diff b/w predicted & actual values is small as possible

② Multivariate Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

Goal: Estimate β 's that minimises RSS.

- Multicollinearity: Occurs when predictors are highly correlated.
- leads to unstable β estimates.
 - Detect using Variation Inflation Factor (VIF)

$$VIF_i = \frac{1}{1 - R^2_i}$$

High VIF (> 10) \rightarrow strong multicoll...

* Use feature selection / regularization (Ridge / Lasso) to fix this.

(singular)

- $X^T X$ is non invertible, causes:

- Multicollinearity
- $p > n$ (features more than samples)

→ If $X^T X$ singular. Linear Regression cannot be solved by normal equations. (use when $p < 1000$)

FIX ?

⇒ GRADIENT DESCENT

$$\text{Cost function } J = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad [\text{MSE}]$$

[RSS $\propto J$]

GD: Optimization technique \rightarrow used for Linear Regression, by minimizing the prediction error. Repeatedly adjust model parameters to reduce the diff b/w \hat{y}_i and y_i .

1. Initialize Parameters: random values for m and b .

2. Calculate Cost function: $\frac{1}{n} \sum_{i=1}^n (y_i - (mx_i + b))^2$

3. Compute Gradient

$$m: \frac{\partial J}{\partial m} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 x_i$$

$$b: \frac{\partial J}{\partial b} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

4. Update parameters :

$$m = m - \alpha \frac{\partial J}{\partial m} \quad (\text{We move opposite to the gradient direction.})$$
$$b = b - \alpha \frac{\partial J}{\partial b}$$

5. Repeat until convergence.

- * Always converges for convex functions like Linear Regression
- * Gradient Descent updates β values in the direction that reduces the total squared residuals - minimize RSS

→ R^2 : Coefficient of Determination
↳ Indicates how much variation the developed model can explain or capture. Range : $[0, 1]$

$$(R^2 = 1 - \frac{RSS}{TSS})$$

→ Adjusted R^2 : Accounts the number of predictors in the model and penalizes the model for including irrelevant predictors that don't contribute significantly to explain the variance in dependent variables. { Helps prevent overfitting }

$$\text{Adjusted } R^2 = 1 - \left(\frac{(1-R^2) \cdot (n-1)}{n-k-1} \right)$$

→ SUBSET SELECTION

Used when there are many predictors, and we may not want to use all.

1. Forward Selection ($\emptyset \rightarrow \text{select}$)

2. Backward Selection (full $\rightarrow \text{select/deselect}$)

REGULARIZATION / FEATURE SELECTION

③ Shrinkage Methods

Techniques that constrain / regularize the coeffs of a regression model to reduce variance and prevent overfitting

→ Shrink large coeffs (β) to 0.

Coeffs can get very large when:

- Multicollinearity

- Model Overfitting

Shrinkage helps by adding a penalty term to the cost function

1. Ridge Regression (L₂ Regularization)

$$J(\beta) = \text{RSS} + \lambda \sum_{j=1}^n \beta_j^2$$

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^n \beta_j^2$$

{ Coeffs are smaller in magnitude }

- As $\lambda \uparrow$, coeffs shrink towards 0, not zero exactly.
- Reduces variance but introduces small bias.
- Always scale before (mean=0, std-dev=1) fitting.

2. Lasso Regression (L₁ Regularization)

$$J(\beta) = \text{RSS} + \lambda \sum_{j=1}^n |\beta_j|$$

(feature)

- Performs both shrinkage & variable selection
- Some β s exactly becomes 0, eliminate them

3. Elastic Net ($L_1 + L_2$)

- Useful when predictors are highly correlated.
- Ridge → Stability + Lasso → Feature Selection.

(high $R^2 \rightarrow$ high bias)
low var)

DIMENSIONALITY REDUCTION / FEATURE EXTRACTION

④ Principal Components Regression (PCR) (PCA + Regression)

Dimensionality Reduction technique \rightarrow applies PCA to the feature matrix X before regression.

↳ Top k principal components (PCs) of X - reducing variance retaining most information.

S1: PCA - transforms correlated predictors into uncorrelated components. $[\text{Cov}(PC_1, PC_2) = 0]$

PC_1 (max Variance) $> PC_2$ (2nd max Variance) $> PC_3 \dots$

- * - Principal components are linear combinations of many features
 - Almost all variables contribute to PCs (even if indirectly)
 - Removes multicollinearity completely - PCs are orthogonal
- * Use PCR when features are highly correlated, but we want to use all info \rightarrow Dimensionality Reduction \rightarrow new uncorrelated features

PCA

1. Center and scale X (input - features contribute equally)
2. Compute Cov Matrix $X^T X$ (Capture feature relationships)
3. Find Eigenvalues & Eigenvectors (Variance explained)
4. Sort by largest Eigenvalues (Pick PCs with max variance)
5. Project X onto M eigenvectors
6. Choose top M PCs
7. Regress y on these PCs

\Rightarrow Need to use ' y ' while forming components? Use:
Partial Least Square (PLS)

⑤ Bias - Variance Tradeoff

In supervised learning, we want to reduce the expected prediction error on unseen data.

Total Error can be decomposed to:

$$E[(y - \hat{f}(x))^2] = \text{Bias}^2 + \text{Variance} + \sigma^2$$

Bias: error from wrong assumption

Variance: sensitivity of data

σ : irreducible error ~ noise inherent in data.

1. Bias: how far model the avg model prediction is from the true function $f(x)$.

↑ High Bias: Model too simple, underfitting

$$\text{Bias } (\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

2. Variance: how much the model's prediction changes with different training data.

↑ High Variance: model too sensitive to noise, overfitting

$$\text{Variance } (\hat{f}(x)) = E[\hat{f}(x)] - E[\hat{f}(x)]^2$$

→ Signs

- Training acc \gg Validation acc : overfitting
- Training & Validation acc low : underfitting
- Train error \downarrow Val Err \uparrow : overfitting

⑥ Cross Validation

A technique used to check how well a ml model performs on unseen data while preventing overfitting.

- Splitting dataset into diff parts
- Training on some parts and testing on remaining.
- Averaging the results from multiple validation step to get the final performance.

(MP)

1. k-fold CV

- split data into k roughly equal parts (folds)
- For each fold,
 - train on $k-1$ folds, test on the remaining fold.
- Compute mean test error.

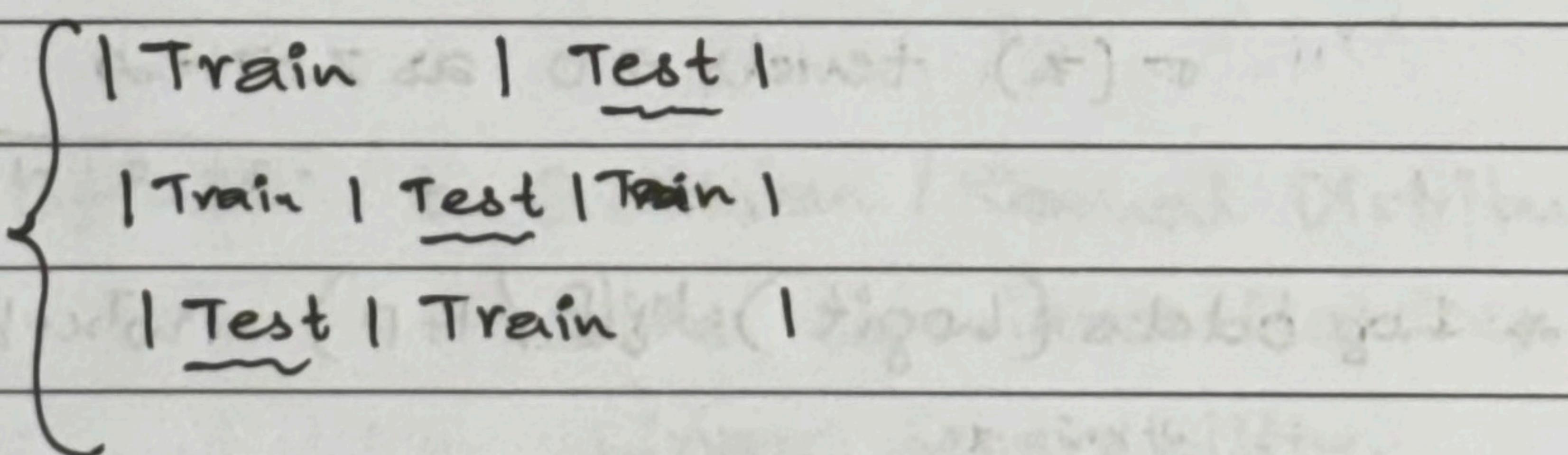
2. Leave one out CV (LOOCV)

- Special case of k-fold where $k=n$
- Train on $n-1$ samples, test on 1, repeat for all

3. Repeated / Stratified CV

Example of

K fold :



Discriminative Model predicts the boundary between classes while generative models learn the underlying data distribution to create new data.