

UCI Heart Disease Dataset

16 columns, 920 rows

Columns

cp → chest pain type

trestbps → Resting blood pressure

chol → serum cholesterol

fbs → fasting blood sugar

restecg → Resting ecg results

thalch → Max. heart rate achieved

exang → exercise induced Angina

num → disease severity

thal → Thalassemia status (Hereditary disease
(haemoglobin))

Ca → No. of major vessels coloured by fluoroscopy

oldpeak - ST depression induced by exercise
relative to rest

Slope - Trend of ST segment during peak exercise



Linear Regression

Model Eqn :- $\hat{y} = b_0 + b_1 x$ $\Rightarrow \hat{y} \in (-\infty, \infty)$

\hat{y} = model prediction

b_0 = intercept

b_1 = regression coefficient / slope

For multiple features:-

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

Cost Function (MSE)

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

Cost Function

$m \rightarrow$ total training samples

Gradient Descent Update Rule

$$\theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

\vdots

$m =$ dataset size

$\alpha =$ learning rate
in gradient descent

$x_j^{(i)}$
 j^{th} feature
of i^{th} example

Evaluation Metrics

$$MAE = \frac{1}{m} \sum | \hat{y} - y |$$

$$MSE = \frac{1}{n} \sum (\hat{y} - y)^2$$

$$RMSE = \sqrt{MSE}$$

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

Logistic Regression

Sigmoid Function

$$f(z) = \frac{1}{1 + e^{-z}} \Rightarrow \hat{y} \in [0, 1]$$

↳
Output

Logistic Hypothesis

$$\hat{y} = \sigma(b_0 + b_1 x_1 + \dots + b_n x_n)$$

Log Loss (Binary Cross Entropy)

$$J(\theta) = -\frac{1}{m} \sum [y \ln(\hat{y}) + (1-y) \ln(1-\hat{y})]$$

Gradient Descent Update

$$\Theta_j := \Theta_j - \alpha \cdot \frac{1}{m} \sum (\hat{y} - y) x_j$$

Odds & Log-Odds

$$\text{Odds} = \frac{P}{1-P}$$

P = predicted probability

$$\log(\text{Odds}) = \ln\left(\frac{P}{1-P}\right)$$

Decision

Boundary

$$\hat{y} \geq 0.5 \Rightarrow 1, \quad \hat{y} < 0.5 \Rightarrow 0$$

Classification

Metrics

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

TP = True Positive

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

TN = True Negative

$$\text{F1} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

FN = False Negative

Random

Forest

Gini Impurity

$$\text{Gini} = 1 - \sum p_i^2$$

p_i = probability of class i at a node

Entropy

$$\text{Entropy} = -\sum p_i \log_2(p_i)$$

Information Gain

$$T_1 - T_2 - \dots - T_n < \dots < \dots$$

$$-I(G) = \text{Entropy}(\text{parent}) - \sum \frac{n_j}{n} \text{Entropy}(\text{child}_j)$$

Number of Features Per Split

Classification : \sqrt{P}

$P = \text{total features}$

Regression : $\frac{P}{3}$

Bias- Variance Decomposition

Total Error = Bias² + Variance + Irreducible Error

Bagging (Bootstrap Aggregation)

Classification = Mode of all data

Regression = Mean of all data