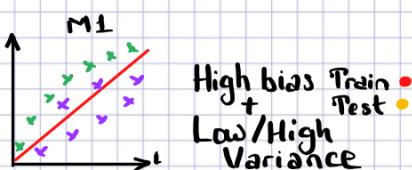


## ① Ridge and Lasso Regression {Overfitting, Bias, Variance, Underfitting}

## ② ElasticNet Regression

## ③ Logistic Regression

### Underfitted Model

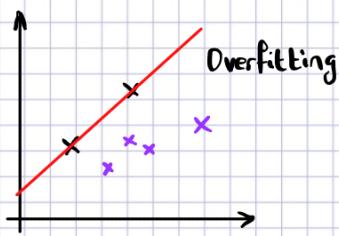


Bias : Training Data

Variance : Test Data

### Ridge and Lasso Regression

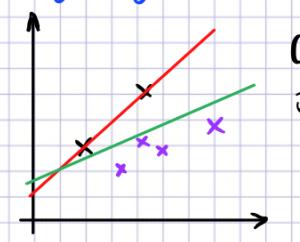
Train Error ↑ | Test Error ↑ { Low bias and high variance }



### Cost function

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

### Ridge Regression (L2- Regularization)

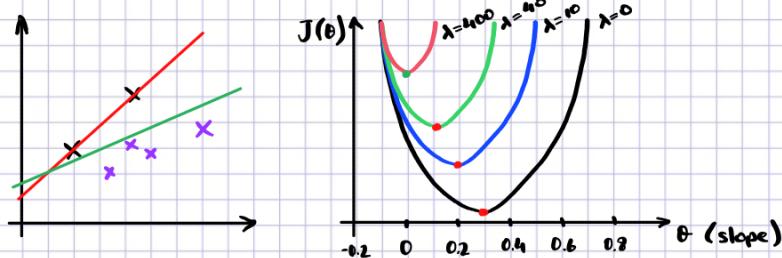


### Cost function (modified)

$$\begin{aligned} J(\theta) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda (\text{slope})^2 \\ &= \text{small number} + \lambda \cdot \frac{(\text{slope})^2}{\text{slope}^2} \\ &\approx 0.71 \end{aligned}$$

↓  
hyperparameter

### Relationship between $\lambda$ and $(\text{slope})^2$



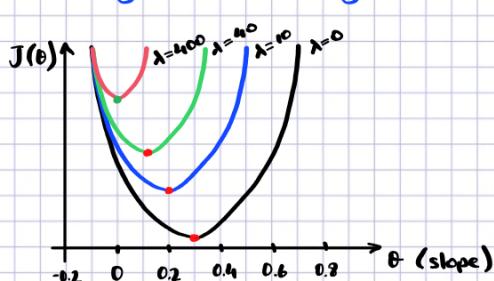
$\lambda \uparrow$  slope ↓ Reduce Overfitting

Global Minima is shifting, while overfitting is being reduced, and if our hypothesis is in form of

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

by reducing the slope ( $\theta_1$ ) we are neglecting feature  $x_1$ .

### Lasso Regression (L1-Regularization)



$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$\theta_i \rightarrow 0$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \cdot |\text{slope}|$$

feature not correlating at all  $\approx 0 \Rightarrow$  Feature neglected

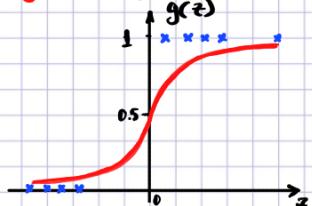
$\lambda \uparrow$  slope ↓ Feature Selection

### ElasticNet Regression

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda_1 (\text{slope})^2 + \lambda_2 |\text{slope}|$$

Reducing Overfitting feature selection

## Logistic Regression [Binary Classification]



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1)$$

$$z := \theta_0 + \theta_1 x_1$$

$$h_{\theta}(x) = g(z) \quad [\text{activation function}]$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Training set:

$$\{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$$

$y = \{0, 1\} \rightarrow 2 \text{ outputs} \rightarrow \text{Binary Classification}$

$$h_{\theta}(z) = \frac{1}{1 + e^{-z}}, z = \theta_0 + \theta_1 x_1, \theta_0 = 0 \Rightarrow \text{intercept} = 0$$

Main Aim: Change  $\theta_1$ , so that our sigmoid function classifies data point

Logistic Regression Cost Function (Log loss)

$$J(\theta) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$J(\theta_0, \theta_1) = y \cdot \log(h_{\theta}(x^{(i)})) + (1-y) \cdot \log(1-h_{\theta}(x^{(i)}))$$

$$J(\theta) = -\frac{1}{2m} \sum_{i=1}^m [y^{(i)} \cdot \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \cdot \log(1-h_{\theta}(x^{(i)}))]$$

$$h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_i)}} \Rightarrow \text{hypothesis} \quad \text{repeat until convergence} \quad \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} (J(\theta)) \right\}$$

## Performance Metrics (Binary Classification)

		Real Label	
		Positive	Negative
Predicted Label	Positive	True Positive TP	False Positive FP
	Negative	True Negative TN	False Negative FN

→ Precision =  $\frac{TP}{TP+FP}$

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

↓

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

We can choose between Precision and Recall after clarifying what our main aim is, whether it's to reduce FN's or FP's.

E.g. Spam classification  $\{FP\} \Rightarrow$  Precision

Patient's cancer (Y/N)  $\{FN\} \Rightarrow$  Recall

But say we have a problem where we need minimal values possible of both FP and FN. In this case we can use F-Beta Score using  $\beta$  as  $\beta$ -value (F $\beta$  score)

$$F\text{-Beta Score} = (1+\beta^2) \frac{\text{Precision} \times \text{Recall}}{\beta^2 \times (\text{Precision} + \text{Recall})}$$

$$F\text{-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$

[Harmonic Mean]

We can adjust this metric by using different values as  $\beta$  depending on what is more important for us whether it is to minimize FP's or FN's

Case 1: minimizing FP's is more important : FP > FN

⇒ reduce  $\beta$  Ex:  $\beta=0.5$

$$\Rightarrow F\text{-Beta Score} = (1+0.25) \times \frac{\text{Precision} \times \text{Recall}}{0.25 \times (\text{Precision} + \text{Recall})} = 5 \times \frac{\text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$

Case 2: minimizing FN's is more important : FN > FP

⇒ reduce  $\beta$  Ex:  $\beta=2$

$$\Rightarrow F\text{-Beta Score} = (1+4) \times \frac{\text{Precision} \times \text{Recall}}{4 \times (\text{Precision} + \text{Recall})} = 1.25 \times \frac{\text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$