

→ Underfitting occurs when a model does not give good fit on the training data itself.

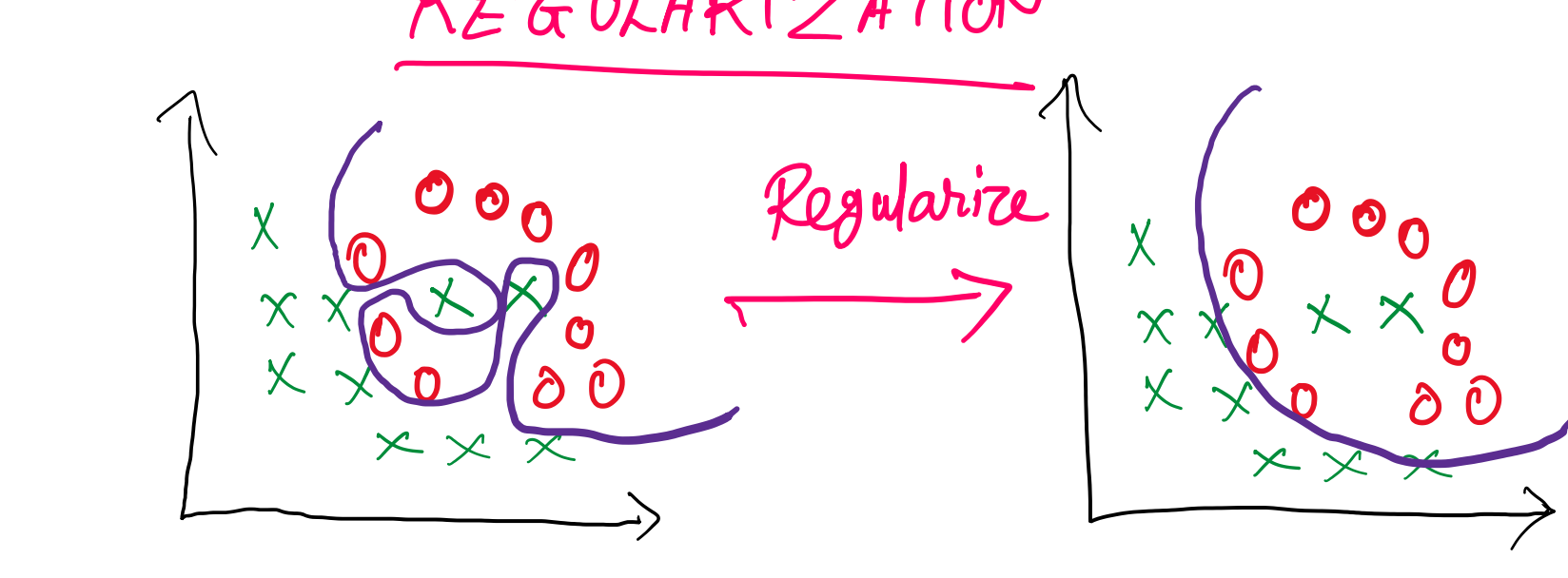
→ Overfitting means the model fits training data very well. However, performs very badly on unseen data.

→ Typically, too simplistic models tend to underfit and too complex models tend to overfit. This is also called bias-variance tradeoff.

→ To prevent underfitting

- Use MORE TRAINING DATA
- Use Slightly More complex Model

→ To prevent Overfitting



Some Popular Types of Regularization

1. L2 Regularization

Eg:  $\|Ax - b\|_2^2 + \lambda \|x\|_2^2$   
 $\lambda > 0$  (Ridge Regression)

2. L1 Regularization

Eg:  $\|Ax - b\|_2^2 + \lambda \|x\|_1$   
 $\lambda > 0$  (LASSO Regression)

3. Elastic Net Regularization

Eg:  $\|Ax - b\|_2^2 + \lambda_1 \|x\|_2^2 + \lambda_2 \|x\|_1$

4. DROPOUT REGULARIZATION

5. Early Stopping.

Model Selection:-

- Use training, Validation and test sets
- Create a model using training data
- Tune it using validation set
- Choose a version that performs best on test data

Evaluation Metrics:- (Regression)

i) Mean Squared Error

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

ii) RMSE

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

iii) MAE

$$= \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

Classification:-

Confusion Matrix

		Actual Values	
		P	N
Predicted Values	P	TP	FP
	N	FN	TN

Labels: Type 1 points to the TP and FN cells; Type 2 points to the FP and TN cells.

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

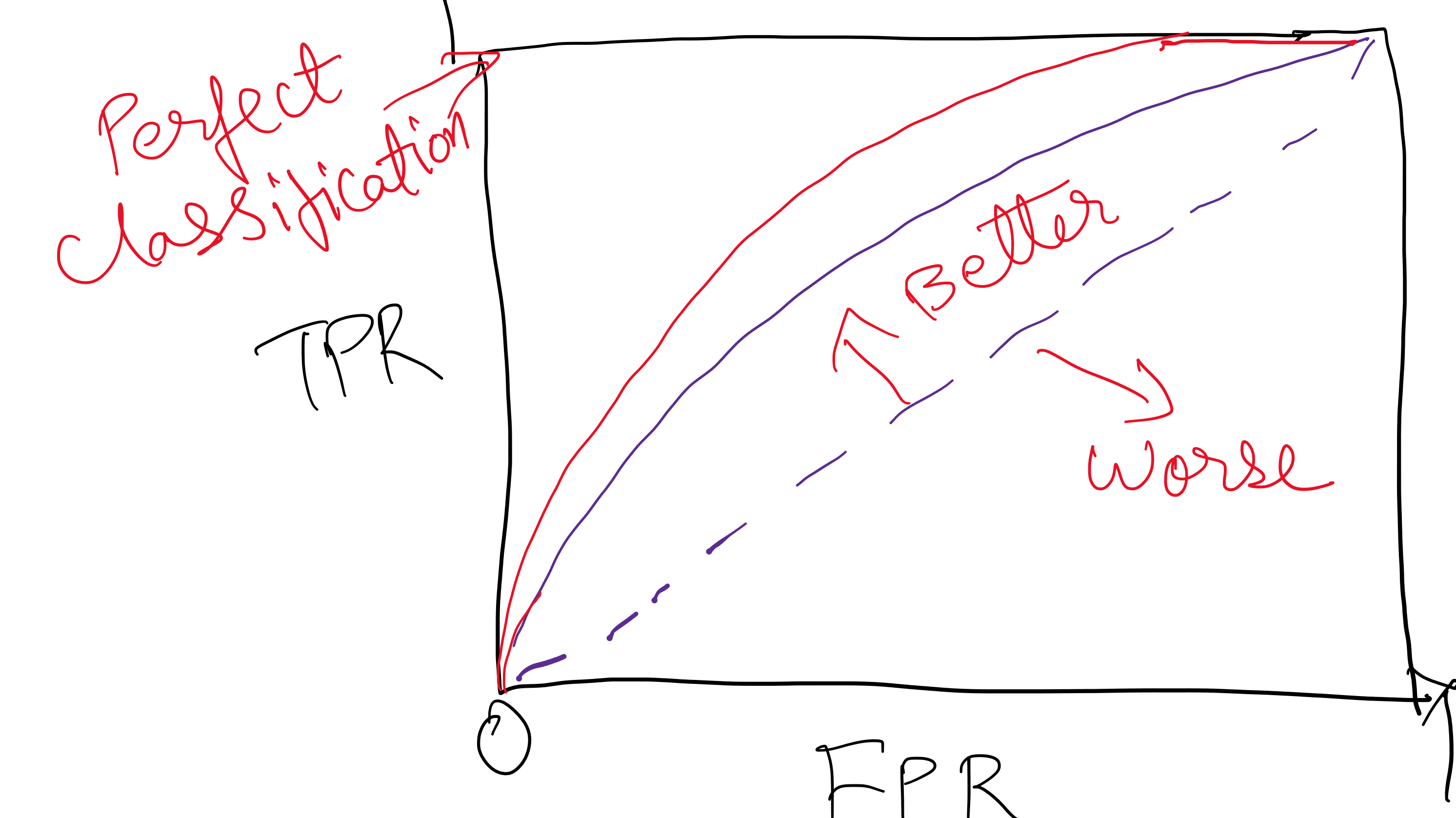
$$\text{Precision} = \frac{TP}{TP + FP}$$

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

$$\begin{aligned} \text{F1 Score} &= \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \\ &= \frac{2 (\text{Precision})(\text{Recall})}{\text{Precision} + \text{Recall}} \\ &= \frac{TP}{TP + \frac{1}{2}(FP + FN)} \end{aligned}$$

ROC curve

$$\text{TPR} = \frac{TP}{TP + FN} \quad \text{FPR} = \frac{FP}{FP + TN}$$



AUC = Area Under Curve  
Higher the better