

Algorithms



Regression : LR

Classification : Logistic Regression

KNN

Naive Bayes

Decision Trees $\begin{cases} C \\ R \end{cases}$

Regression : R² score

Classification Metrics

?

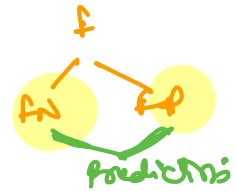
→ way to measure the quality



Confusion Matrix



		Predicted label	
		0	1
(Actual) True label	0	True (TN) Negative	False (FP) Positive
	1	False Negative (FN)	True Positive (TP)



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

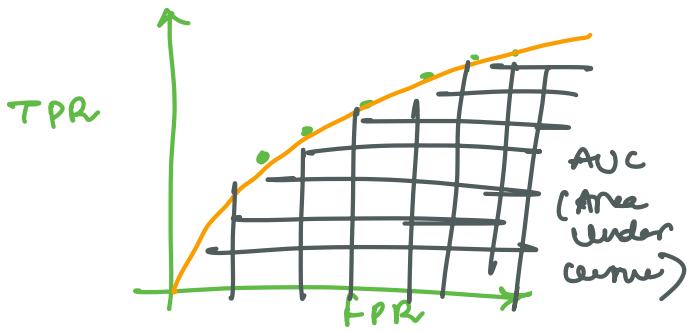
$$\text{Precision} = \frac{TP}{TP + FP} \quad \left. \vphantom{\frac{TP}{TP + FP}} \right\} \begin{array}{l} \text{out of all predicted +ve, how} \\ \text{many were actually correct} \end{array}$$

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

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

↓
HM of Precision & Recall

ROC Curve (Receiver Operating Characteristics)



Different values of threshold

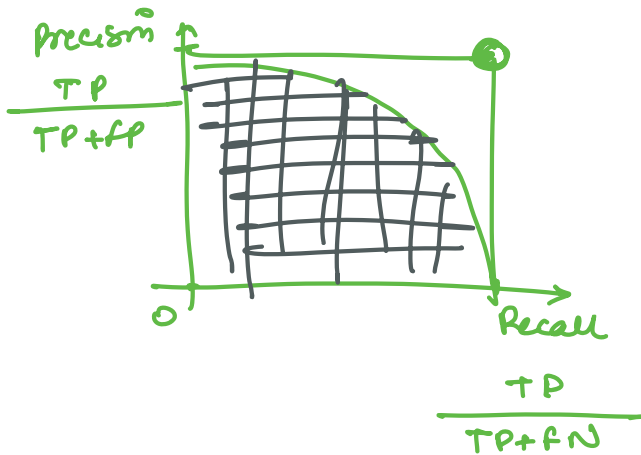
threshold $\rightarrow 0.5$ $\begin{matrix} < 0.5 \rightarrow 0 \\ \geq 0.5 \rightarrow 1 \end{matrix}$

0.8 $\begin{matrix} < 0.8 \rightarrow 0 \\ \geq 0.8 \rightarrow 1 \end{matrix}$

$$\text{True Positive Rate} = \frac{TP}{TP+FN} * 100$$

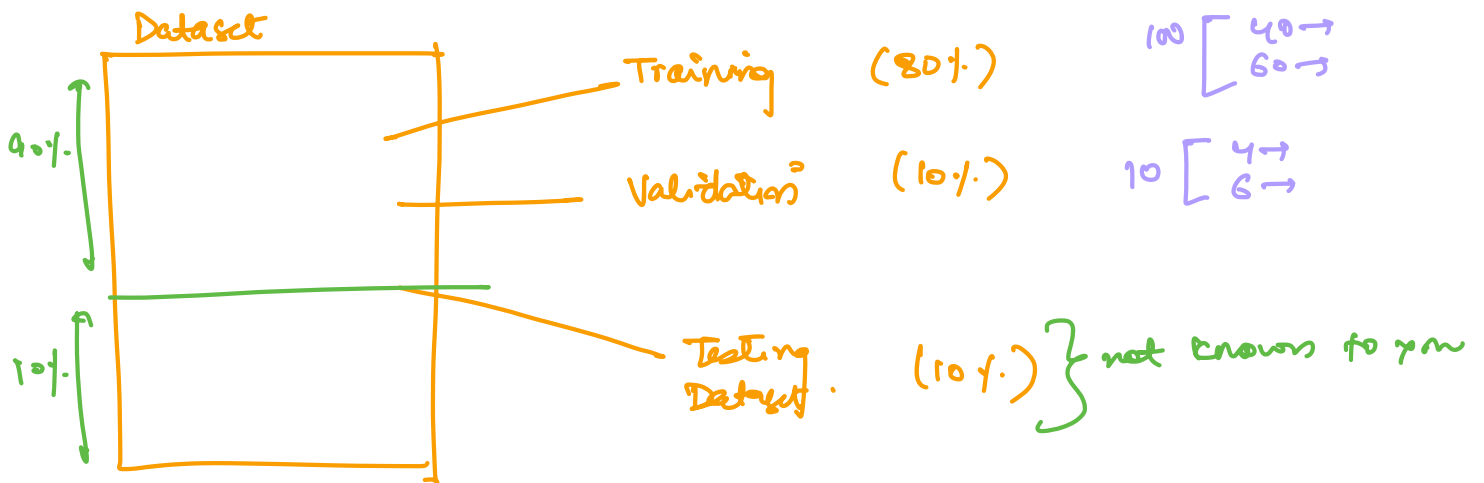
$$\text{False Positive Rate} = \frac{FP}{FP+FN} * 100$$

PR Curve



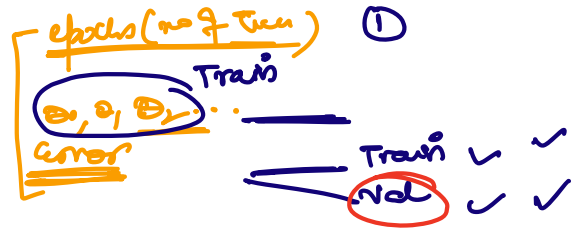
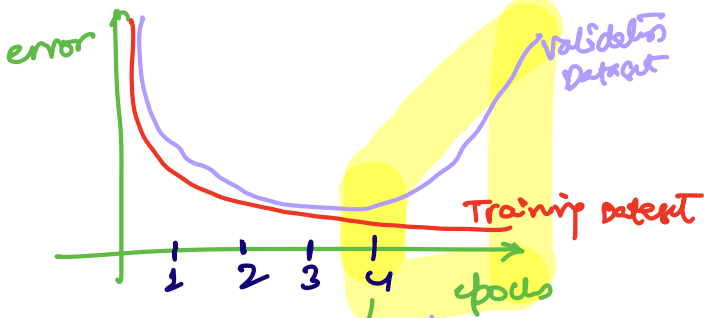
ideally:

high precision $FP=0$
high recall $FN=0$

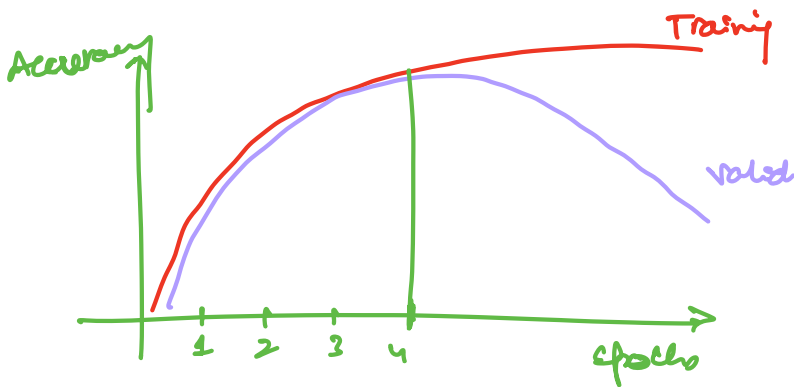


logistic regression

Train Data: Model Train
parameters learn:
 $\theta_0, \theta_1, \theta_2, \dots$

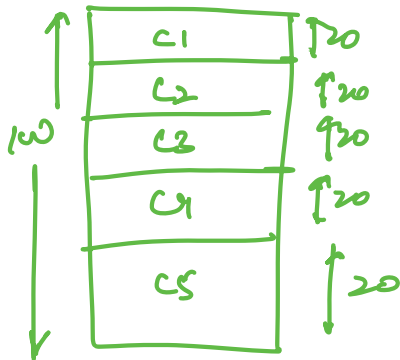


Model overfit
 $\theta_0, \theta_1, \theta_2$ best



k fold cross validation:

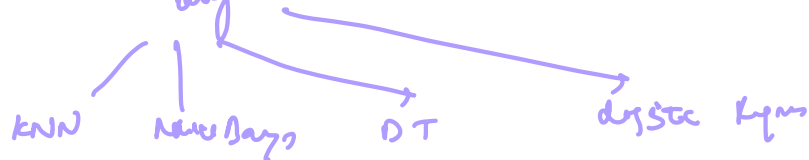
Train + Val



$k = 5, 7, \dots$

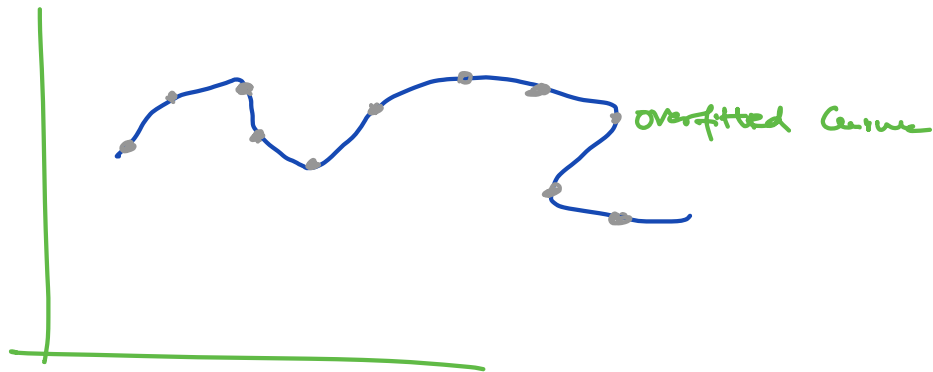


step	TD	VD
1:	C1 C2 C3 C4	C5
2:	C1 C2 C3 C5	C4
3:	C1 C2 C4 C5	C3
4:	C1 C3 C4 C5	C2
5:	C1 C2 C3 C4	C5



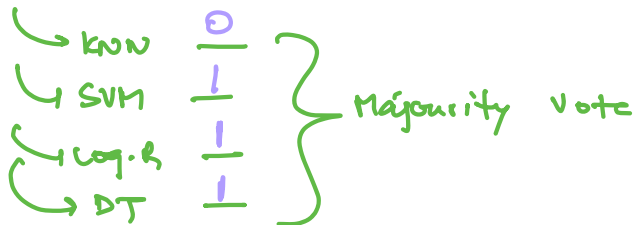
Overfitting:

Model performs very well on training data but does not perform good on test data.



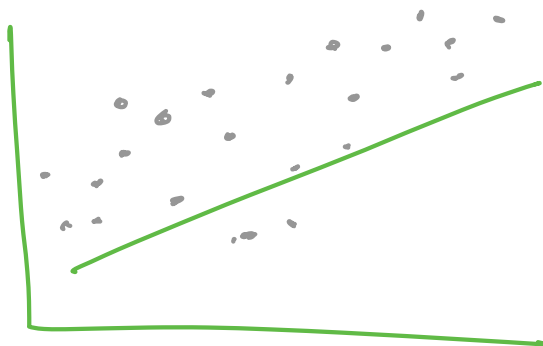
Solution:

- k fold cross validation²
- Sufficient data → examples
- Ensembling technique.



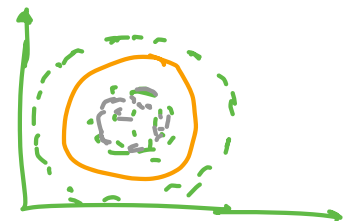
Underfitting:

Model is not going to learn patterns from training data.



underfit model will give poor performance on train as well as test data.

②

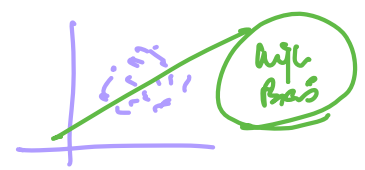


Solutions:

- No. of features increase
- Models Complexity
- Reduce noise
- Increase the duration of training.



BIAS AND VARIANCE:



Bias: wrong assumptions about data — like assuming data is linear in reality is follows a complex form.

It is the inability of the model bcz of that there is diff in predicted value & actual value.

Low Bias: lower assumptions make a model which closely matches the training dataset.
(Simple model)

High Bias: more assumptions model will not match training dataset closely.
(Complex model)

Variance:

↳ measure of spread in data from its mean position.

↳ sensitive to a subset which follows same distribution as your training dataset.

Low variance: model is very less sensitive to changes.

High variance: model is very sensitive to changes and it can result in significant changes if trained on a different subset.

