

Logistic Regression2

09 February 2022 20:46

Logistic Regression

EDA :- Find variables that have some relation with 'Churn'

Remove Variables that are having multicollinearity ✓

→ Missing values

→ outlier (input variables)

Observation

Imbalanced dataset

1 - 14.5%

0 - 85.5%

Balanced : 1 : 30% ; 0 : 70%
Yes No

future

Imbalanced dataset

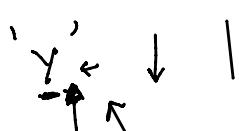
Imbalanced learning

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_m x_m = \log \left(\frac{p}{1-p} \right)$$

$$\text{Odds ratio} = (e^\theta)$$

Target (heart Risk) ← age, BMI, HDL, LDL, Smoker, Active

θ_1 θ_2 θ_3 θ_4 θ_5 θ_6 ↑



Odds ratio (e^θ) 1.5 2.2 0.7 1.8 5.6 ✓ 0.06 ↓

pred

CM

		Pred		
		N	P	
Actual	N	TN	FP	2
	P	FN	TP	20
[100]		80	20	✓

Precision $\frac{20}{22}$ ~ FP moderately
Fancy)

diagnosing Cancer

$$Recall = \frac{20}{100} ? = 0.2$$

FN is NOT acceptable

Relaxing

		N	P	
		20	100	500
Pots		400		

comes with cost

$$\text{Accuracy} = \frac{T.P + T.N}{T.P + F.P + F.N}$$

Precision and Recall $F-1 \text{ Score} = \left(\frac{2 \cdot P \cdot R}{P+R} \right)$

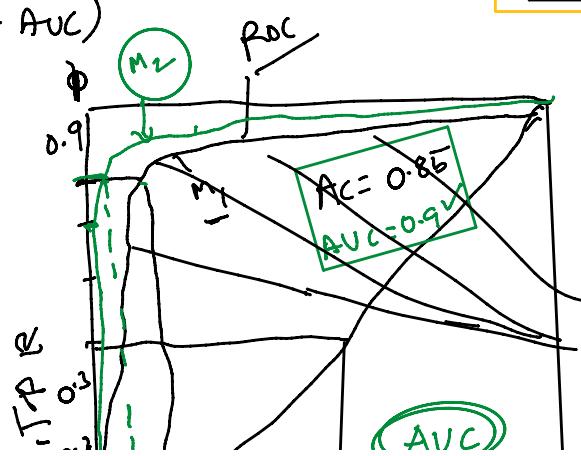
① Why 0.5 cut off?

② How do we compare models?

ROC Curve

		Cutoff	0.1	0.2	0.9
Prob	Pred	0	0	0	0
		0.1	1	0	0
0.2		1	1	0	0
...		1	1	0	0
...		1	1	0	0

(ROC-AUC)



Actual

N	P
0	TP
1	FN

+ F.P.R

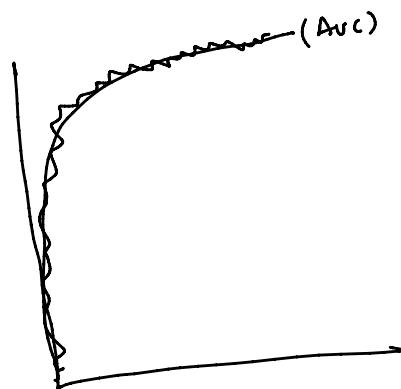
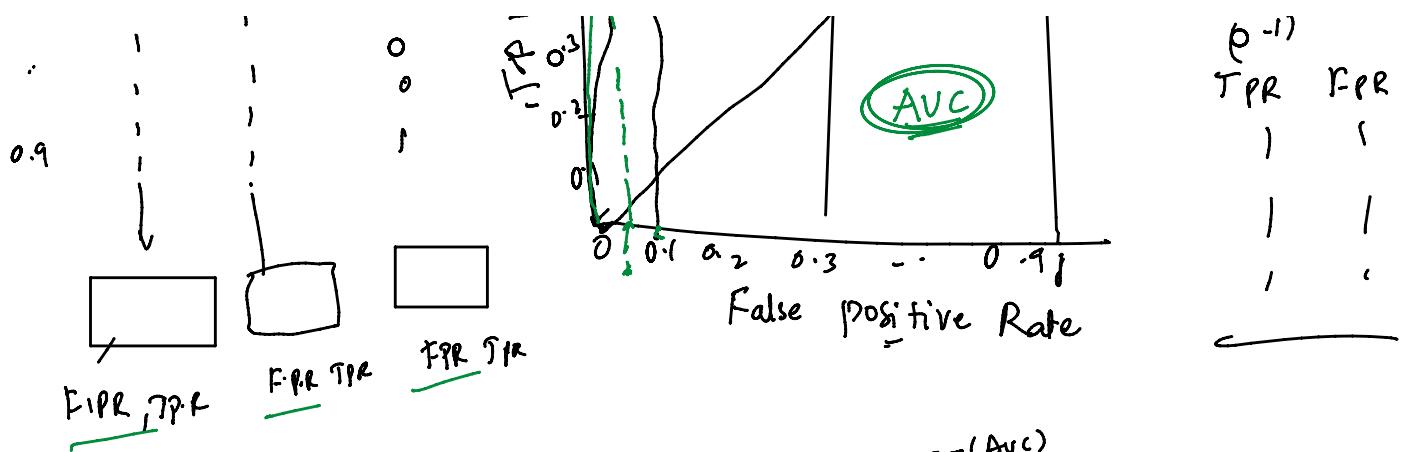
N	P
FN	FP

← TPR

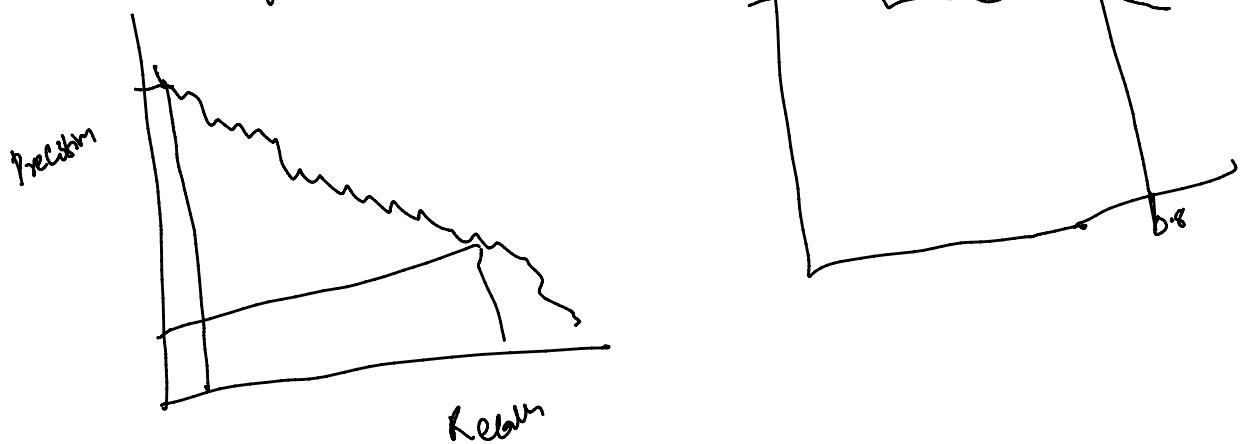
$$T.P.R = \frac{T.P}{T.P + F.N}$$

$$F.P.R = \frac{F.P}{F.P + T.N}$$

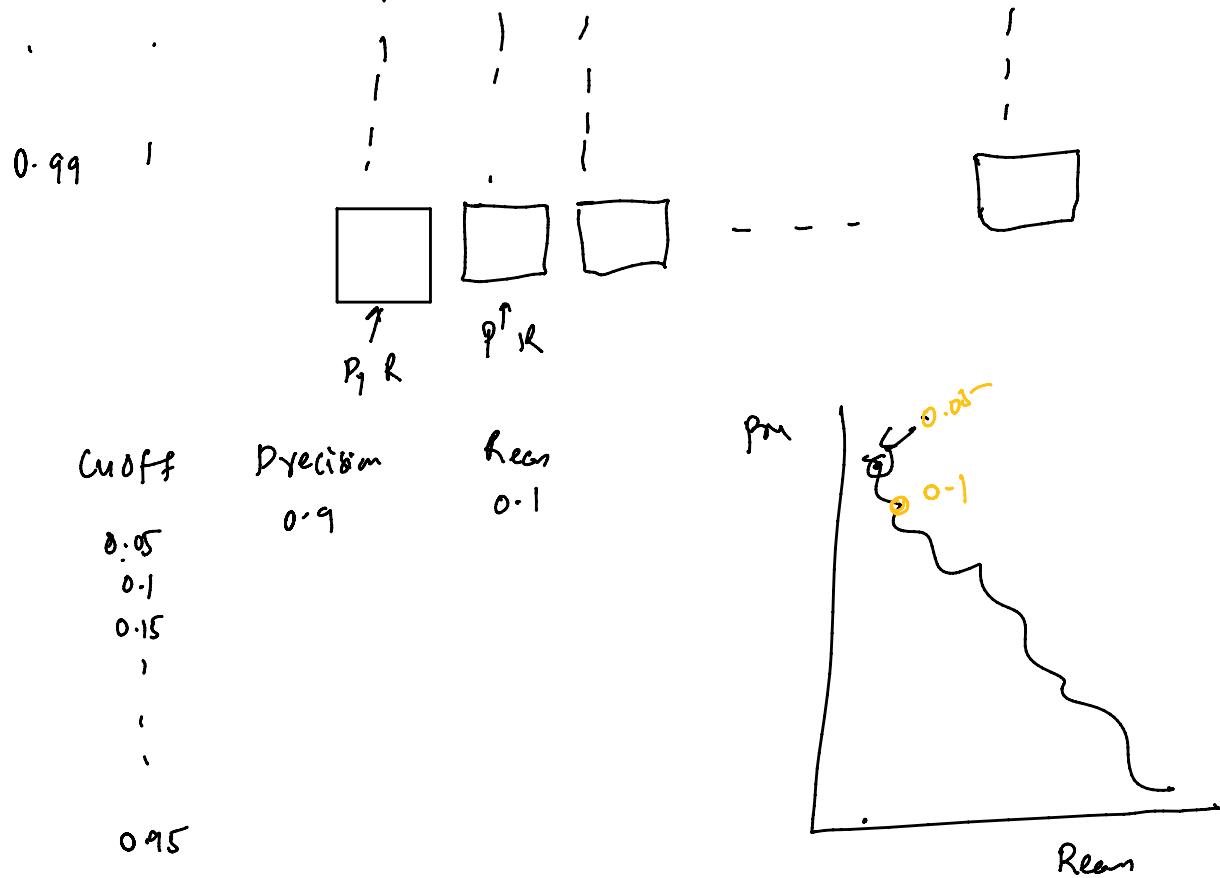
$$\frac{P^{-1}}{T.P.R \quad F.P.R}$$



Precision vs Recall graph (PR Curve)



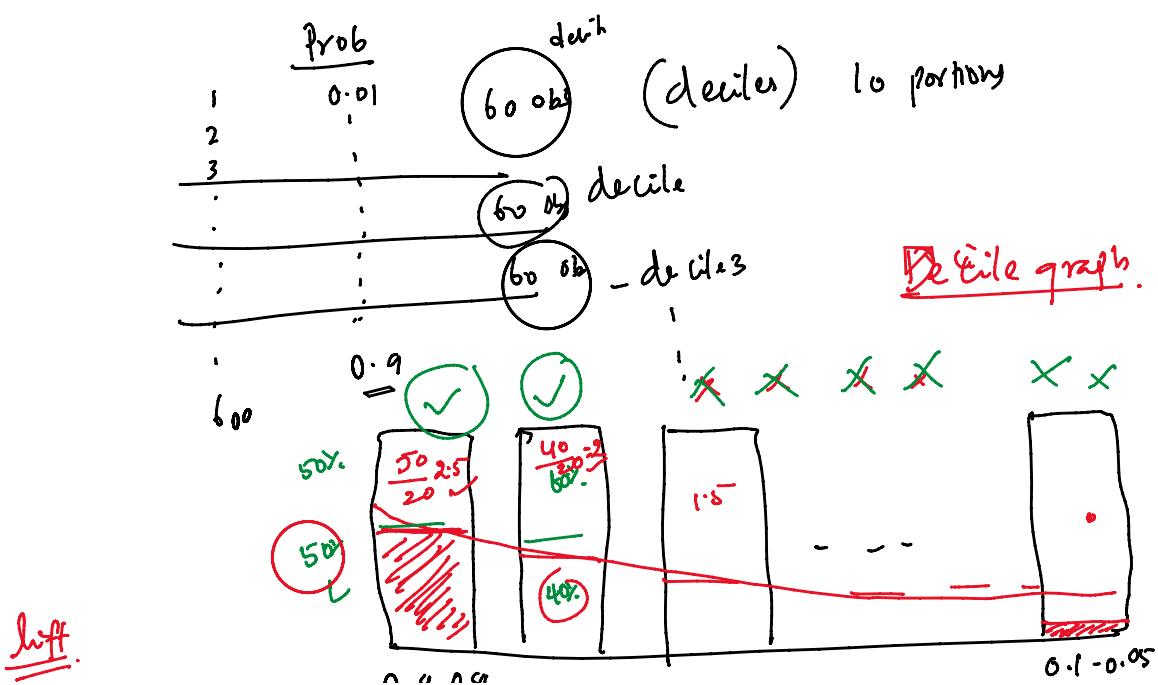
\hat{y}	Actuals	Cutoff 0.05	0.0	0.15	...	0.95
\hat{y}_{prob}	y	$\overline{P_C}$	P_C	$\overline{P_C}$		
0.1	0	1	0	0		0
0.11	0	1	1	0		0
0.12	1	1	1	0		0
.	0	1	1	1		0
.	.	1	1	1		1
.	.	1	1	1		1



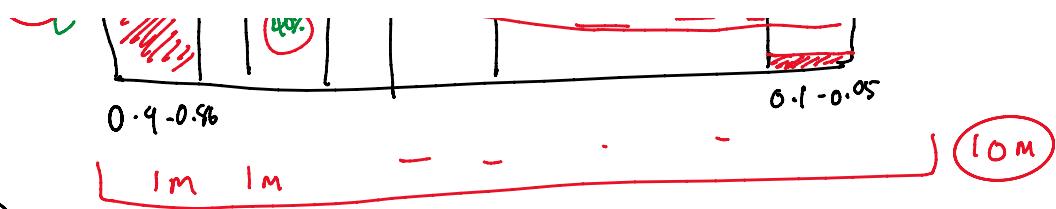
Analysis / charts

(Benchmark)

20% first 16.8/call
≥ 5



buff.

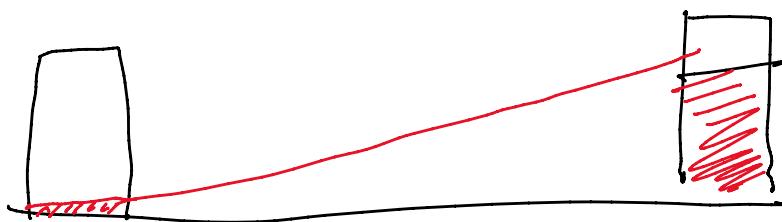


Randomly (100) \rightarrow (20%)

Top decile (high prob) Model 2.5

Decile 100 \rightarrow (50%) 20%

Decile 2 \rightarrow 40% 20%



	M	T	W	T	F	S	S
Pred	N	N	R	R	N	N	N
	N	N	R	N	R	N	N

A green arrow points from the first 'R' in the 'Pred' row to the first 'R' in the second row. A yellow arrow points from the second 'R' in the 'Pred' row to the second 'R' in the second row. A green arrow points from the 'P-' label to the second 'R' in the second row.