

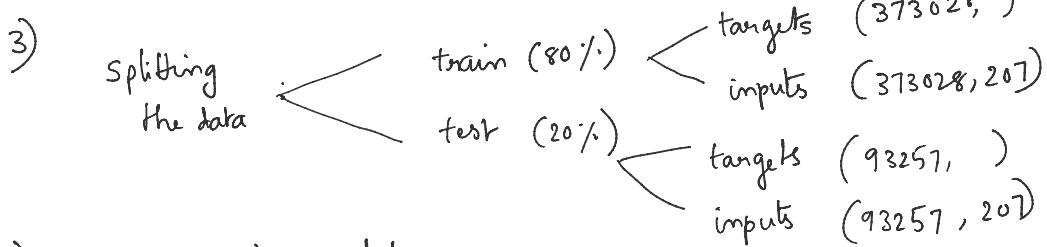
A MODEL DEVELOPMENT :



2) Dependent Variable / Target :

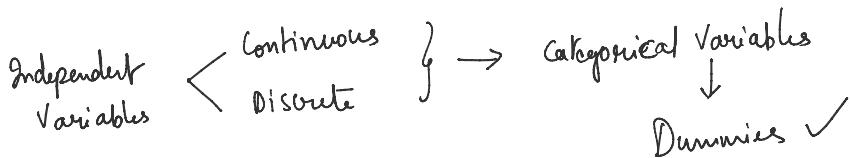
Good → Non-default → 1

Bad → Default → 0



4) More processing of data

↳ which is not needed for test data but needed
for train data for modeling



* Dummy Variables

* Binning on the basis of WOE

WOE = How much evidence does the independent variable
has w.r.t Variation in the dependent variable

$$\begin{aligned}
 WOE_i &= \ln \left(\frac{\% (y=1)_i}{\% (y=0)_i} \right) = \ln \left(\frac{\% \text{ Good}}{\% \text{ bad}} \right) \\
 &= \ln \left(\frac{\% \text{ Non-default}}{\% \text{ default}} \right)
 \end{aligned}$$

line

clustering
Coarse

$$IV = \sum_{i=1}^K \left[(\%.\text{good} - \%.\text{bad}) \times \ln \left(\frac{\%.\text{good}}{\%.\text{bad}} \right) \right]$$

(K-categories)

IV → can be used to identify the variables which have better predictive power w.r.t dependent variable.

5)

MODEL BUILDING :

<u>Feature / inputs</u>	<u>Co-efficients</u>	<u>p-values</u>
↓	↓	↓

filtering out the significant features

based on p-value ($< 0.05 \rightarrow$ Rejecting the null hypothesis)

(Co-efficient \rightarrow significant)

6)

Higher co-efficients \rightarrow greater the odds of being a good borrower

observation \rightarrow odds of categories $\xrightarrow{\text{can be compared with the help of}}$ Reference Categories

7)

As we have trained the model (using train)

Now we have to test (using test)

Now we have to draw one more (ROC curve)

Now we have to test (using test)

VALIDATION:

i) probabilities will be estimated

Decide on the cut-off (0, 1)
Range

ii) Confusion Matrix

		Predicted	
		1	0
Actual	1	TP	FN
	0	FP	TN

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

out of total predictions

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

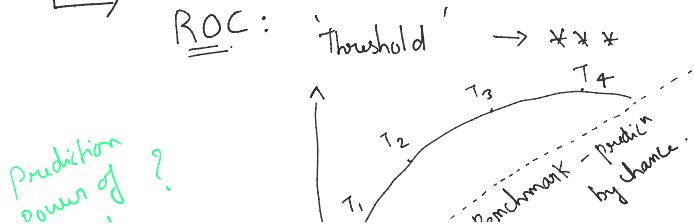
P^3 - precisely predicted positives

$$(TPR) \quad \text{Recall} = \frac{TP}{TP + FN}$$

Rightly Identified positives out of total positive

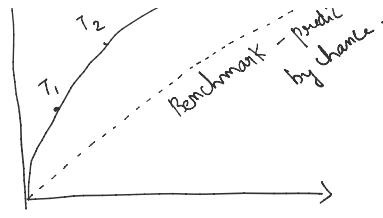
$$4) \text{ Specificity / TNR} = \frac{TN}{TN + FP}$$

$$5) F1 = \text{Harmonic Mean of precision \& Recall} \rightarrow \left[\frac{2 PR}{P + R} \right]$$



AUROC - Area Under ROC

Prediction power of Model ?
TPR (Recall)

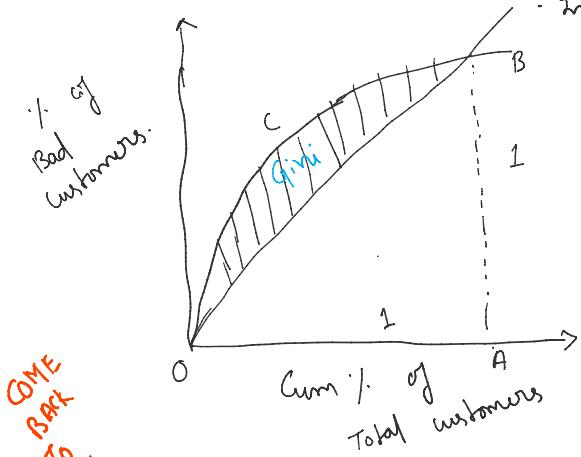


FPR (sensitivity) - Not specificity.



Gini:

(Measure of inequality) (originally economics - Income inequality)



COME BACK
THIS → CHECK TO

$$\text{Area } \Delta OAB = \frac{1}{2} \times 1 \times 1 = \frac{1}{2}$$

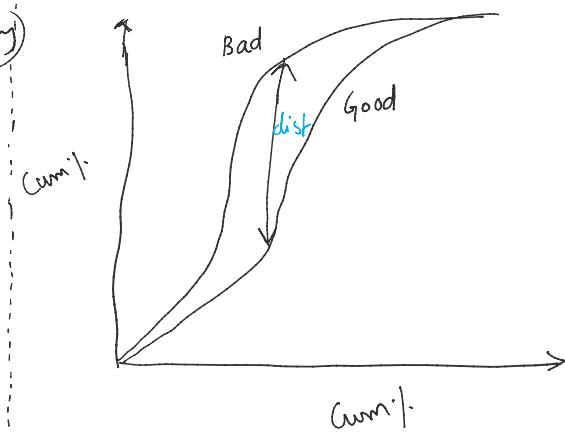
Gini = Shaded reg'

AUROC = \int_0^1 area under

$$\text{Gini} = \text{AUROC} - \frac{1}{2}.$$

Discriminatory power of model

Kolmogorov - Smirnov:



(Very high)

Perfect model \rightarrow Max. dist $\rightarrow KS = 1$
pred'n by chance \rightarrow Max. dist $\rightarrow KS = 0$

(Very low)

→ Applying Model:

Suppose new customer

falls under what all dummy variables

→ Consider only those co-efficients

(if my_cat \rightarrow then co-eff = 0)

Sum of all co-eff \rightarrow prob of non-default

$(1 - \text{prob of non-default}) \rightarrow \text{prob of default}$



Scorecard Development: (including my-categories) \rightarrow for easier interpretations.

(Min = ? \rightarrow a borrower falls into lowest co-eff categories (zeros/neutrals) = -1.49

Scorecard - mapping: (including n/a categories) → for easier interpretations.

Scores	Min = ? → a borrower falls into lowest co-eff categories (zeros/negatives)	(zeros/negatives) = -1.49
	Max = ? → a borrower falls into highest co-eff categories	= 5.6332

Rescaling co-efficients → Scores

$$\text{Variable Score} = \text{Variables-Coff} \times \frac{(\text{max-Score} - \text{min Score})}{(\text{max-SumCoff} - \text{min SumCoff})}$$

scaling intercept to around min-Score

$$\begin{array}{l|l} \text{Min Score} = 360 & \text{scaling done ✓ scoring ✓} \\ \text{Max Score} = 850 & \end{array}$$

↳ Calculating Credit Scores: ✓

↳ Calculating PD from credit scores: (reverse the process)

Banks → more business → lower cut-off of PD

less " → higher " .. "

Setting the Cut-offs:

from ROC calculations

FPR
TPR
thresholds
1)
2)
3)

Derive 'Good' proba from scores

TABLE: thresholds FPR TPR Score

Now lets do the

- 1) Approval Rate
- 2) Rejection Rate

- 3) No. of Approved Applicants ✓
- 4) No. of Rejected Applicants ✓

Based on own risk appetite:

Export scorecard to CSV



Model monitoring:

check whether the model is still relevant

Suppose → now data is 2007 to 2014

Suppose \rightarrow now data is 2007 to 2014

we got new data \rightarrow 2015

whether our model still relevant for this 2015 data

for that we need to calculate the Population Stability Index PSI

PSI = formula

Interpretation of values

'PSI' calculated for all features. $= 0$ No difference in population

$PSI \in [0, 1]$ $= 1$ Absolute difference.

Generally $PSI > 0.25$ \rightarrow indicating to rebuild the model / changes to be made

(check score
airly)

(action to be taken)

tion changed significantly.