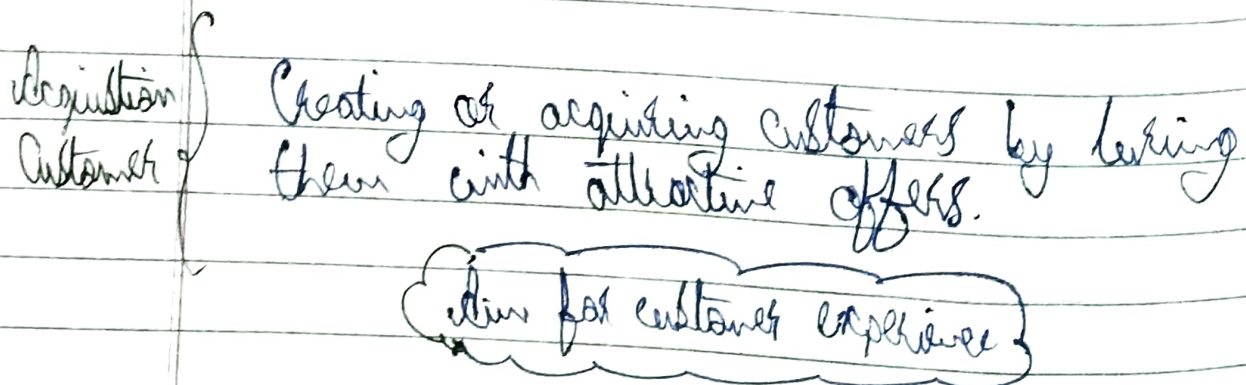# Bank Case Study

* Predicting Bank Loan Default for Bring loyalty

* Sales strategy to influence existing / new customers

|  | Cross-Selling | Up-Selling |
|---|---|---|
| Increasing Revenue | To buy Complementary products / services in addition to org. purchase | To purchase a higher valued product or servicing of same product line. |

*Aim for effectiveness*

| | |
|---|---|
| Minimizing Customer Churn [Attrition] | When customers leave your product / service due to several reasons. |

*Aim to (↑) Promotional Offers*

| | |
|---|---|
| Acquisition Customer | Creating or acquiring customers by luring them with attractive offers. |

*Aim for customer experience*

**✱ Predicting Fraudulent Activity**

Illegal activities in order to receive money, funds, credits from bk. / Institution.

Phishing   Cr./ Db Scam   Check   Accounting   Identity   Money
                            Scans     fraud        fraud    laundering

**✱ Predicting Bank Defaults with Log Regrn.**

Logit model demonstrates a regrn. model where, resp. var. is binary / dichotomous & indept. var. can be binary / continuous / ordial

**✱ Logit types**
- Bi-Nominal = Default / Not-Df.
- Multi-Nominal = Yes / No / Maybe  [No. Res]
- Ordered = Ratings

$$Eg: \ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \beta_n X_n$$

where,
P = Probl. of Binary outcome.

Cost

**\* Odds Ratio** = $\dfrac{\text{Prob. of occurrence of events}}{\text{Prob}^{l}\text{. of occurrence of non events}}$

| Table | D | ND | Total |
|-------|-----|-----|-------|
| U | 50 | 10 | 60 |
| E | 40 | 20 | 60 |

Odds Ratio$_U$ = $\dfrac{\left(\frac{50}{60}\right)}{\left(1-\frac{50}{60}\right)}$ = $\dfrac{(0.8)}{(0.2)}$ = $\boxed{4}$

Odds Ratio$_E$ = $\dfrac{\left(\frac{40}{60}\right)}{\left(1-\frac{40}{60}\right)}$ = $\dfrac{(0.6)}{(0.4)}$ = $\boxed{1.5}$

Odds Ratios = $\dfrac{OR_U}{OR_E}$ = $\dfrac{4}{1.5}$ = $\boxed{2.66}$

**\*\*** The odds for an unemployed person to default on a loan is 2.66 times higher than employed.

**\* Logit Model Curve** = S - Shape

= ⌐╱

✻ <u>Assump.<sup>n</sup></u>

- Resp. Var. must be binary
- logit model estimates $P(Y=1)$, so resp. var. must abide by this assump.<sup>n</sup> & outcome must be in line with this code
- Model should be as perfect fit as possible nor over neither under.
- • Indept. var.'s must not be correlated to each other.
- Indept. var.'s must be linearly correlated to log odds.
- Sample size must be large

✻ <u>Model Evaluation</u>

- Likelihood Ratio test : Used to compare (2) nested glm's & if p-value for the full model $< 0.05$, then its a better fit.

   Hyp : Ho : Full model better , against
   Ha : Reduced model better

$$\chi^2 = (-2) \ln (\text{Reduced}) - (-2 \ln (\text{full}_{\text{model}}))$$

★
> Not recommended

- Hosmer Lemeshow Test : Use Pearson-Chisq test to check whether obs. prop. of events are same to predicted prob. in model popl. subgroup (10 groups).

$$H = \sum_{g=1}^{10} \frac{(O_g - E_g)^2}{(E_g)}$$

where,

$O_g$ = No. of obs. events .
$E_g$ = No. of expected events.

Hyp : $H_0$: Model is good , against
$H_a$: Model is poor

★ Sig. of Ind". Indept. Var . [Tests to check]

★ Wald Statistic Test : Helps us determing the importance of an individual dept. var. by logistic regrn. coff.

$$W_j = \left( \frac{B_j}{SE_{B_j}} \right)^2$$

Hyp : $H_0$ : Coff. of Interest $= 0$ , against
$H_a$ : " " " $\neq 0$

1) WST accepts $H_0$ ⇒ Indept. var. will not impact model's fit

Drawback - Large coff. could get removed becoz of ($\uparrow$) SE ; as a result, we could underfit the model.

2 Likelihood Ratio : $\left( G = (-2)\left( \ln (\text{reduced}) - \ln (\text{full}) \right) \right)$

The smaller the deviance b/w both models the better is corel. b/w dept. & Indept. var

* Predictive Value Validation [Model Accuracy]

Few measurements like Confusion Matrix & receiver operating characteristic helps us see, how accurately the model is predicting dept. var.

- Confusion Matrix : Technique used to evaluate predictive accuracy of the model

Predicted class

| Conf. Mat. | | x | y |
|---|---|---|---|
| Actual class | X | TN | FP |
| | Y | FN | TP |

True / false — Actual Outcomes
+ve / -ve — Predicted Outcomes

(+ve) = Predicted yes , (-ve) = Predicted No.

= Ho = Customers would not churn, against

= Ha = Customers would churn

## Predicted

| Actual | | (+ve) | (-ve) |
|---|---|---|---|
| | T | Correct | Correct |
| | F | Type I Error | Type-II Error |

FP =. Not occurring predicted as occurred

FN = Occurring predicted as not occurred

| | | Actual (Disease) | | |
|---|---|---|---|---|
| | | + | — | |
| Predicted (claim) | + | True (+) | False (+) | Type-I Error |
| | — | False (-) | True (-) | |

Type-II Error

★ Accuracy Rate = $\dfrac{\text{Accurately predicted}}{\text{Total Outcomes}}$

★ Error (Misclassification) Rate = $1 - \dfrac{\text{Accuracy}}{\text{Rate}}$

( Sensitivity )

★ ( T(+ve) Rate = $\dfrac{\text{Accurately predicted (+ve)}}{\text{Total actual (+ve)}}$ = $\dfrac{TP}{TP + FN}$ )

★ ( F(+ve) Rate = $\dfrac{\text{Inaccurately predicted (+ve)}}{\text{Total actual (-ve)}}$ = $\dfrac{FP}{FP + TN}$ )

★ ( Specificity = $\dfrac{\text{Accurately predicted (-ve)}}{\text{Total actual (-ve)}}$ = $1 - FPR$ )

★ ( Precision = $\dfrac{\text{Accurately predicted (+ve)}}{\text{Total (+ve)}}$ = $\dfrac{TP}{TP + FP}$ )

★ ( Prevalence = $\dfrac{\text{Actual (+ve)}}{\text{Total Outcomes}}$ = $\dfrac{TP + FN}{TN + FN + TP + FP}$ )

★ Reciever Operating Curve : Used to measure binary classifier performance visually & AUC is used to Quantify model performance.

★ Plots TPR on y axis , against FPR on x axis

* Generally, AUC > 70% = Accurate Model
  So, the more AUC ~ 100%, the more accuracy

* ROC = Above diag = Better
  ROC = Below diagonal = worse
  ROC = $\approx 1$ = Best