

Project Topic → "Telecom Churn Modeling" 8-10 AM

Approach → ① Theory and approach  
                  ② Python Implementation }

# ① (Customer Lifecycle)

Customer acquisition

→ Bringing new customers to the business

- i) Marketing Ads
- ii) SMS / whatsapp

Customer early Engagement

customer are encouraged to increase usage of services

[ pre paid → (Financial stability & CIBIL score) ↑  
post paid → ↑ ✕ unimportant ]

# ③

Cross sell and business growth

(Selling additional services and make the customer more engaged)

# ④

Customer retention

assess profitability and retain the customer  
let the cust go

## Customer Acquisition

- i) Come up with features or factors that can be predictive of customer value

e.g. CB12 score

Higher CB12 score, means better post paid customers

- Q. who gets the post paid offer?

## Early Eng

- ii) what can be done to increase usage of services?

- Q. Which set of treatments can improve engagement?

→ digitally → Targeted Ads

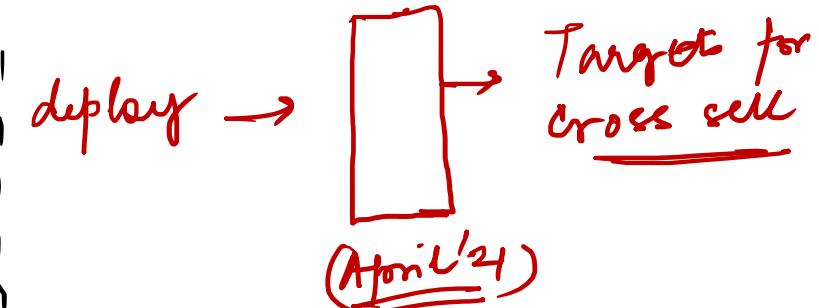
## Cross sell

- iii) what's additional that can be sold?

- ① Build a response model from earlier campaigns (Past campaign) Dec'20



Present day (May'21)



## Customer Retention

Attrition → (Canceling any ongoing relationship)

### Reactive

- i) Customer calls in  
and wants to cancel
- Profitability assessment
- ① Offers will be made

"Pro active"

- i) Predicting which customers are likely  
to attrite

(Our use case for the project)

(Project  
Objectives) → Develop a predictive framework for enabling  
proactive retention strategy for a Telecom company

Project Presentation flow :- (For interviews)

- ① speak about the impact first (Potential Benefits)  
mention what was enabled by your solution. For e.g. 30%  
reduction in cost. OR 20% increase in profit.
- ② which part of lifecycle is addressed by your solution ?
- ③ Justify the approach (supervised / unsupervised)
- ④ steps performed. to solve the problem.

Motivation → "Reduce the churn rate and improve retention"

Outlier analysis }  
Bivariate analysis }

### Missing Value Imputation for Numerical Features

```
from sklearn.impute import SimpleImputer  
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')  
num_1=pd.DataFrame(imputer.fit_transform(num),index=num.index,columns=num.columns)
```

### Missing Value Imputation for Categorical Features

```
from sklearn.impute import SimpleImputer  
imputer = SimpleImputer(missing_values=np.nan, strategy='most frequent')  
char_1=pd.DataFrame(imputer.fit_transform(char),index=char.index,columns=char.columns)
```

Some useful concept to be used in Python Implementation

## Databases

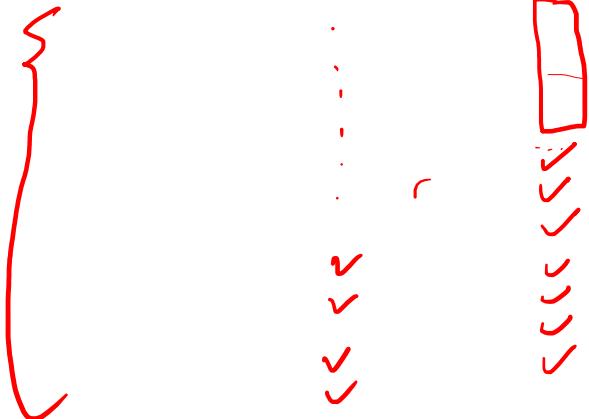
- i) Automated capture of data.
- ii) Tools are used to scan forms, OCR etc

 → CC → self reported / not reported → missing value  
optional.

→ The business records a non activity as missing value. (missing value carries a meaning)

$x_1$	$x_2$	$x_3$	$x_n$	$x_5$
✓	✗	✓	✓	✓
				→ ✗

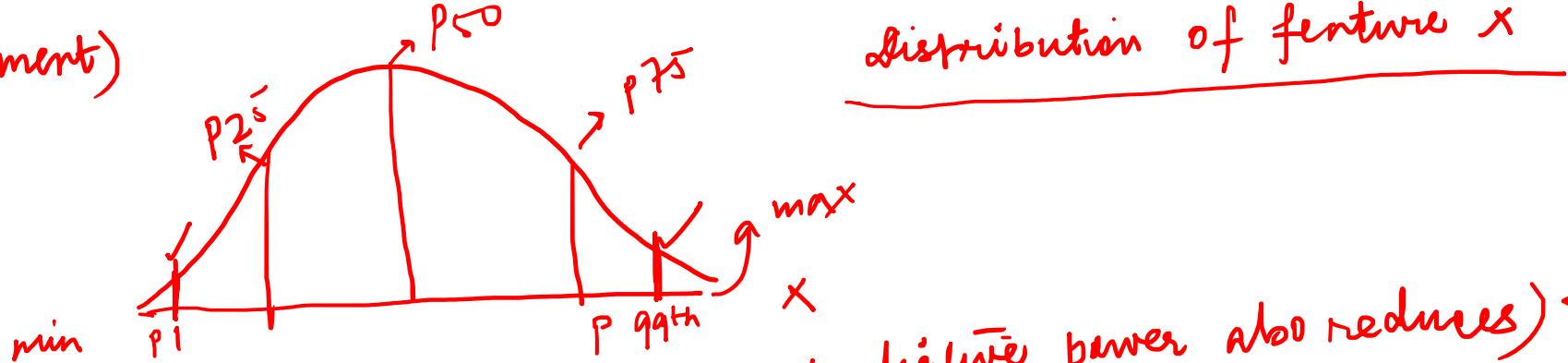
$x_1$  60%  $x_2$  15%.  $x_3$  2%.  $x_4$  0%  $\rightarrow$  (missing value percentages)



(MISSING VALUE  
TREATMENT)

- {
- As a data scientist I do not know what data capture problem led to missing values.
- If a feature records more  $\frac{1}{4}$  of its rows as missing, then it is best to use such a feature as a lost goal.
- $E(\bar{x}) = \mu$

(Outlier treatment)

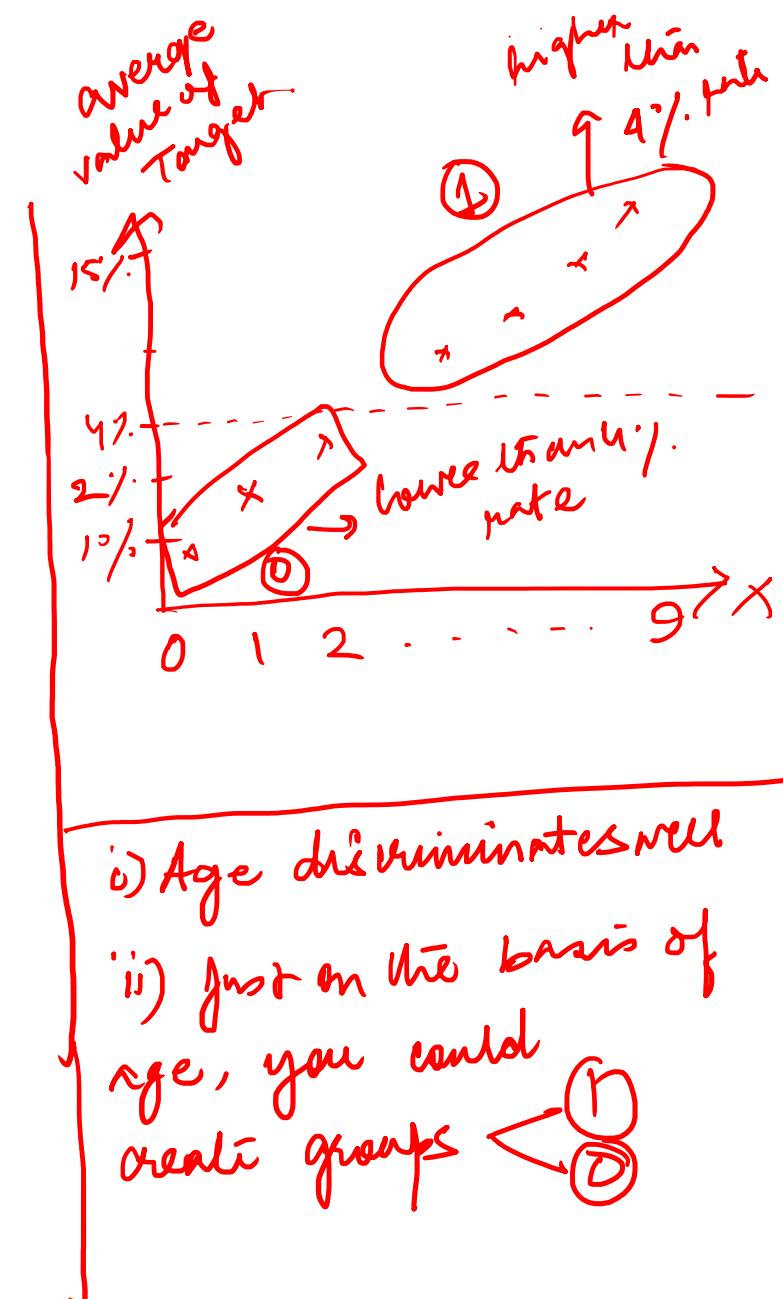
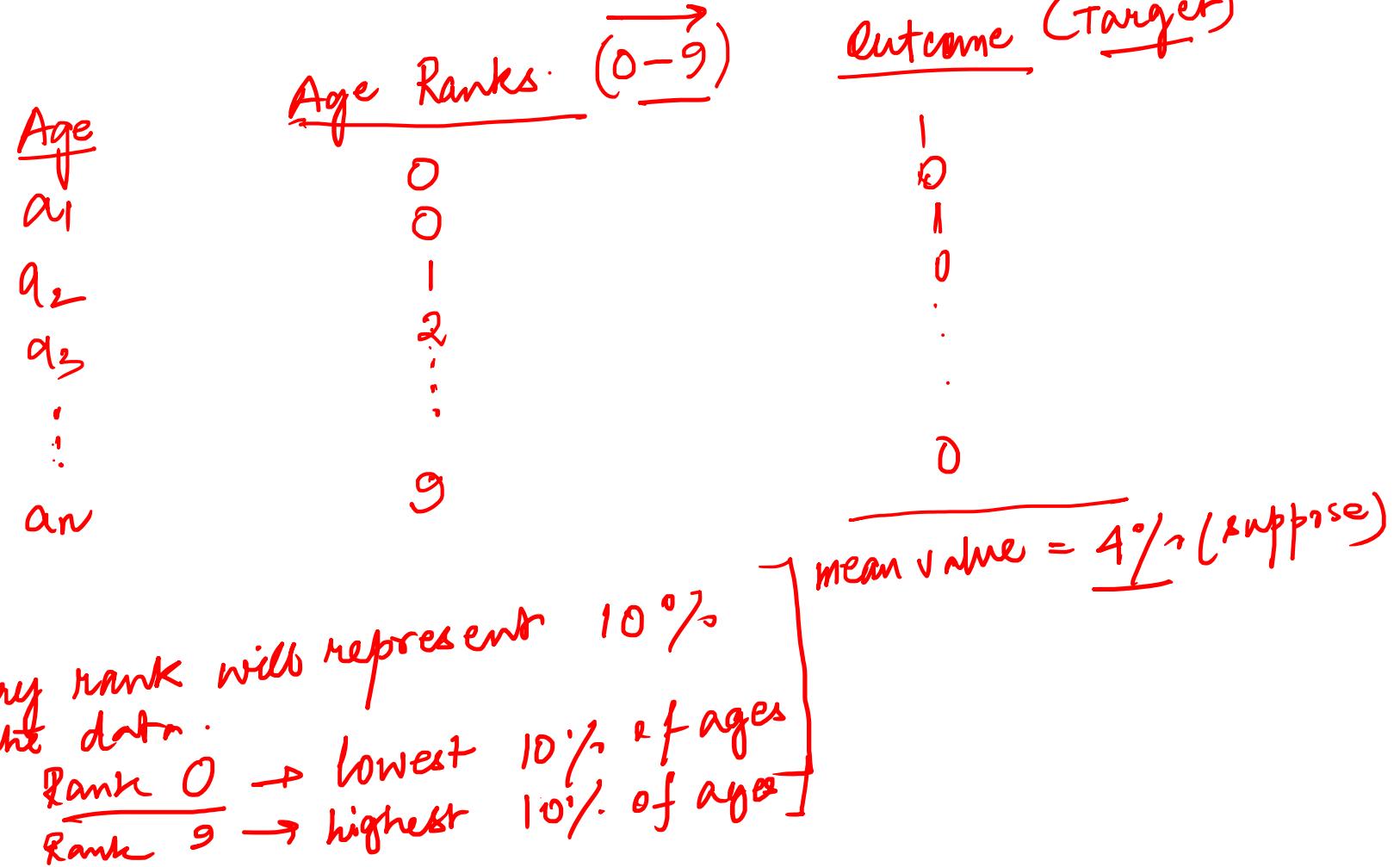


distribution of feature  $x$

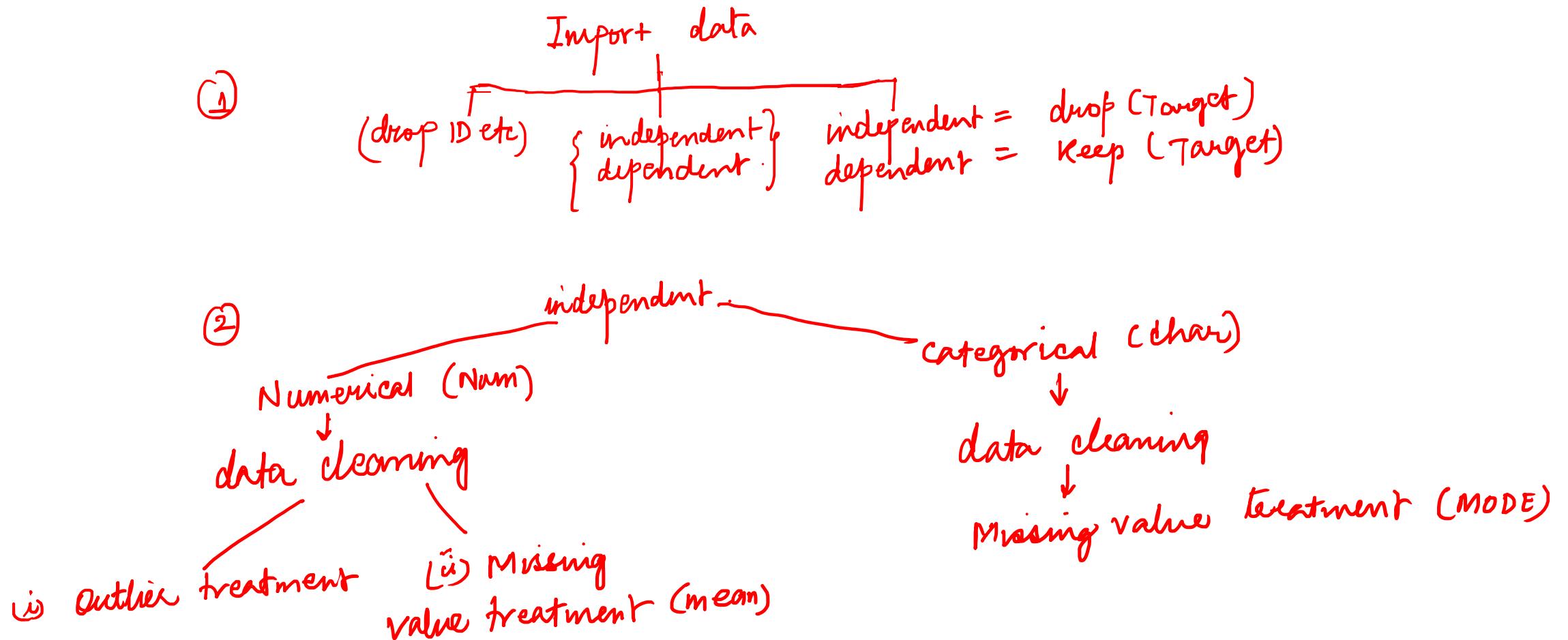
- As variance of a feature reduces the predictive power also reduces) \* \*
- Treat the extreme values in such a way that variance reduction is kept at a minimum
- if  $x < p1$  then  $x = p1$  else  $x \rightarrow$  FLOORING  
if  $x > p99$  then  $x = p99$  else  $x \rightarrow$  CAPPING

## Bivariate Analysis

white Box method for selecting features.



# Python Implementation Flowchart



## (Feature selection)

### ③ Num (Numerical df)

- i) Remove all those features that have a 0 variance
- ii) we will use the class called K Bins Discretizer
- iii) develop Bi-variate graphs → {remove non predictive features}
- iv) use (Select K Best) → select the final set of features from (num)

### ④ char (categorical df)

- i) Develop Bi-variate graphs → {remove non predictive features}
- ii) select K Best

## (Applying Algorithms)

statistical  
 { Logistic }  
 $\log\left(\frac{y}{1-y}\right) = f(x)$

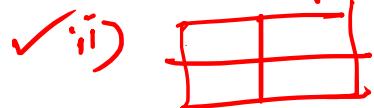
ML (tree Based)  
 { Decision Tree  
 Random forest  
 GBM }



(predict 50% cut off)

### Evaluation

✓ i) Accuracy, precision etc



✓ ii) confusion matrix

[Probability of  $Y=1$ ]

Reality → modelled output

iii) LORENZ CURVE

$x_1$	$x_2$	$x_3$	$x_4$	$y$		
1	2	3	4	1	=	
2	3	4	5	2	=	
3	4	5	6	3	=	
4	5	6	7	4	=	
5	6	7	8	5	=	
6	7	8	9	6	=	
7	8	9	10	7	=	
8	9	10	11	8	=	
9	10	11	12	9	=	
10	11	12	13	10	=	

10 Ranks / Deciles  
 in ascending order.

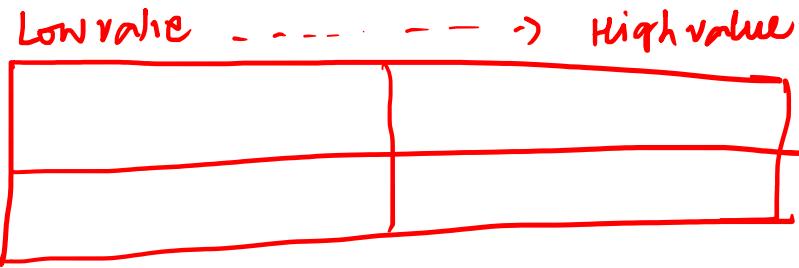
$$\left\{ \frac{\text{Total defaulter}}{\text{Total customer}} = x\% \right\} \uparrow x\%$$

Top 3

(Based on prob of  $Y=1$ )

Bottom 7

"



- {
- Ⓐ we will identify the audience for targeting
  - Ⓑ we will identify the treatments to apply (what to offer)
  - Ⓒ we will chalk out a prioritization from Low → High priority for targeting
- create recommendations based on the model outcome — ★★★