

## Last class (Aug 17)

- 1) Overview of AT&T Churn Prediction
- 2) Accuracy Metric
- 3) Hyper-parameter Tuning
- 4) Logit / Log Odds
- 5) Impact of Outliers
- 6) Multi-class classification - OVR vs multi-nomial

## Today's class

- 1) Spam vs Non-Spam: Business Case
- 2) Issue with Accuracy
- 3) Confusion Matrix & Code
- 4) Precision & Code
- 5) Recall & Code
- 6) F1 score & Code

$$850 + 350 = 1200$$

← class 0:  $\frac{850}{1200} \times 100\% = 70.83$   
non-spam

← class 1:  $\frac{350}{1200} \times 100\% = 29.17$   
spam

$$850 \rightarrow 1100 \quad \underline{1200}$$

$$\text{class 0: } \frac{1100}{1200} \times 100\% = 91.67\%$$

Dumb Model  $\rightarrow$  91.67%  
 ref  $\uparrow$

ML Model  $\rightarrow$  92% Acc

dumb model $\leftarrow$	Ref	ML Model
$\downarrow$ majority class	50%	93% ✓✓ very good
	71%	93%
	92%	93% $\rightarrow$ almost useless model

class 0 :  $\rightarrow$  50%  
 class 1 :  $\rightarrow$  50%

class 0, 1, 2  $\rightarrow$  33.33%

Acc  $\rightarrow$  class 0      class 1  
 $\searrow$   
 X  
 doesn't give  
 Separate performance

...

		Predicted	
		0	1
Ground Truth	0	TN ①	FP ②
	1	FN ③	TP ④

2x2  
matrix  
for  
2 classes

y	$\hat{y}$	Name
→ 0	0	TN ✓
→ 0	1	FP
→ 1	0	FN
→ 1	1	TP ✓

$$\begin{aligned}
 \text{Acc} &= \frac{\# \text{ correct pred}}{\# \text{ total pred}} \\
 &= \frac{TP + TN}{TP + TN + FP + FN}
 \end{aligned}$$

		Pred		
		class 0	class 1	class 2
Ground Truth	class 0			
	class 1			
	class 2			

FP → Type 1 error

FN → Type 2 error

Spam vs Non-Spam Email

↓

pos class

class 1

↓

neg class

class 0

[ Scenario I: inbox receives a promotional email  
↳ FN

Scenario II: Your Google often letter

Maximize Rec  $\Rightarrow$  minimize FP  $\hookrightarrow$  goes to Spam folder  $\boxed{\text{FP}}$   $\downarrow$  Dangerous

$\downarrow$   
 Non-Cancer

Chen 1  
↓  
Cancers

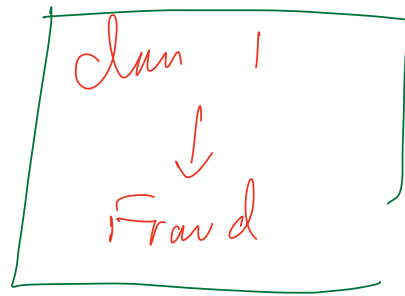
Scenario 1: ML model says a healthy patient has cancer  $\rightarrow$  FP

Scenario 2: ML model says a Cancerous patient  
Recall is maximized is healthy  $\rightarrow$   $\downarrow$  FN very dangerous

FN is minimized

Class 0

↓  
Non-Fraud  
(Legit)



→ FP

Scenario 1: ML model says a legit transaction is Fraud

Scenario 2: ML model says a fraud transaction is legit  
→ Customer is facing inconvenience

Bank/Person is suffering loss

FN

TP: Fraud labelled as Fraud

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

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

$$FP \uparrow, \text{Prec} \downarrow$$
$$\text{Prec} \propto \frac{1}{FP}$$

$$FN \uparrow, \text{Rec} \downarrow$$
$$\text{Recall} \propto \frac{1}{FN}$$

Precision  $\times$  FN

Recall  $\times$  FP

$$HM(v_1, v_2) = \frac{2 \times v_1 \times v_2}{v_1 + v_2}$$

			diff
M1	0.30	0.80	→ 0.5
M2	0.20	0.90	→ 0.7
M3	0.70	0.40	→ 0.3
M4	0.29	0.30	→ 0.01

Revenue → 10:35

$$H.M(v_1, v_2) = \frac{[GM(v_1, v_2)]^2}{AM(v_1, v_2)}$$

$$= \frac{(\sqrt{v_1 v_2})^2}{v_1 + v_2} \times \frac{2}{v_1 + v_2}$$

Dataset: 400 email samples

→ 40 email → spam (class 1)

→ 360 email → non spam (class 0)

Ideal Model

TP, TN, FP, FN

$$\text{Precision} = \frac{TP}{TP + \textcircled{FP}} = 1 \quad \begin{array}{c} \downarrow \\ 0 \end{array} \quad \begin{array}{c} \downarrow \\ 0 \end{array}$$

$$\text{Recall} = \frac{TP}{TP + \textcircled{FN}} = 1 \quad \begin{array}{c} \downarrow \\ 0 \end{array}$$

Dumb Model

Predictions = Majority

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{0}{0 + 0} = \frac{0}{0} \rightarrow \text{undefined}$$

= class 0  
= -ve class

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{0}{0 + 40} = 0$$

		Prediction	
		0	1
Ground Truth	0	360 TN	0 FP
	1	40 FN	0 TP

TN:  $y = 0$   
 $\hat{y} = 0$   
 FN:  $y = 1$ ,  
 $\hat{y} = 0$

Dumb Model (-ve class majority)

Prec  $\rightarrow$  undefined

Recall  $\rightarrow$  0

Precision =

$$\frac{0}{0 + 0 + \boxed{10^{-6}}} = 0$$

$\epsilon \rightarrow 10^{-6}$

TN	00	FP	01
FN	10	TP	11

$$\frac{1,1}{1,1 + 0,1}$$

$\downarrow$   
Precision



$$\frac{TP}{TP + FN} \quad \text{Recall} = \frac{1,1}{1,0 + 1,1}$$

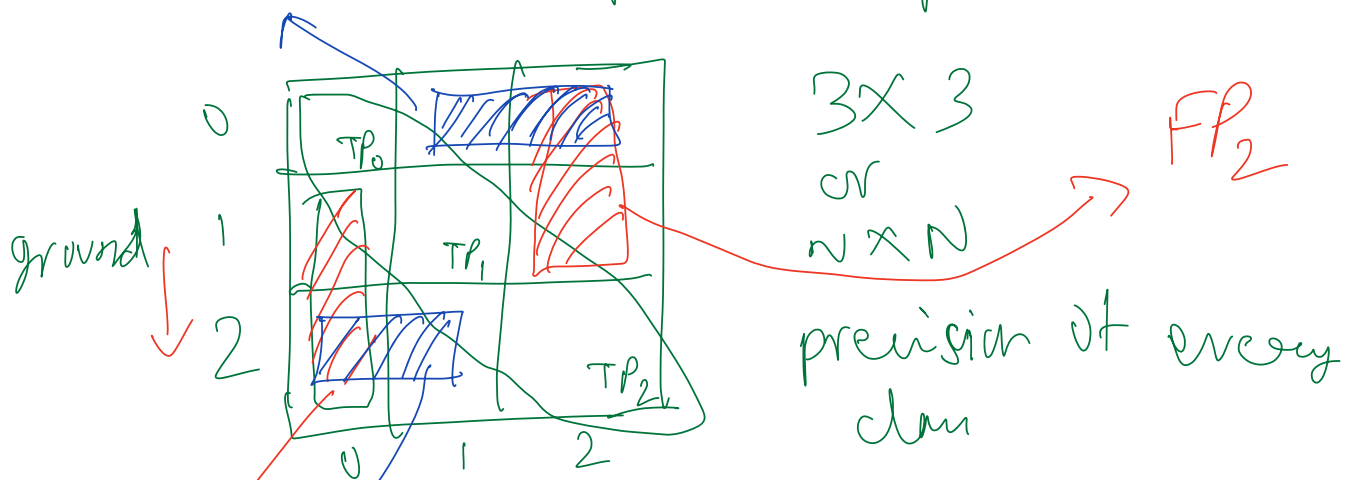
$$\epsilon = 10^{-6}, \quad \epsilon_1 = 10^{-4}$$

$$F_\beta = (1 + \beta^2) \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \text{ precision}) + \text{recall}}$$

$\beta = 2 \Rightarrow$  2x imp. to precision  
as compared to recall

$\beta = 0.5 \Rightarrow$  recall is 2x more

imp than precision



$\swarrow$   $\nwarrow$   
 $FN_2$   $\xrightarrow{\text{pred}}$   
 $FP_0$

Precision (class 0)

$$= \frac{TP_0}{TP_0 + FP_0}$$

Precision (class 2)

$$= \frac{TP_2}{TP_2 + FP_2}$$

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

$$\boxed{y = 1, \hat{y} = 0}$$

$\rightarrow$   
 $\text{pred} = -ve$   
 $gd = +ve$

$$\text{Recall (class 0)} = \frac{TP_0}{TP_0 + FN_0}$$

400 email  $\rightarrow$

GD

		0	1
0	0	TN	40 FP
1	0	FN	360 TP
		pred	

360 spam  $\rightarrow$  class 1

40 non-spam  $\rightarrow$  class 0

$\hookrightarrow$  pred of dumb model  
 $\hookrightarrow$  class 1

$$\text{Precision} = \frac{360}{360+40} = 0.9$$

$$\text{Recall} = \frac{360}{360+0} = 1$$

$$F1 = 0.947$$

↙ dumb

$$\boxed{ML} \rightarrow F1 \rightarrow \begin{matrix} 0.95 & 0.053 \\ & 0.947 \end{matrix}$$

$$ML \rightarrow F1 = 0.98 \quad 0.947$$

$$(0.33)$$