

# Error Analysis for Skewed Data

*Handling Skewed Data*

Advice for Applying Machine Learning:

## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

```
function y = predictCancer(x)
    y = 0; %ignore x!
return
```

Windows'u Etkinleştir  
Windows'u etkinleştirmek için Ayarlar'a gidin.

## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
return
```

0.50% error

## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

skewed classes.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
return
```

0.5% error



## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

→ skewed classes.

```
function y = predictCancer(x)
```

```
→ y = 0; %ignore x!
```

```
return
```

0.5% error

99.2% accy (0.8% error)

99.5% accuracy (0.5% error)

Is there an  
improvement or  
not?

Windows'u Etkinleştirin  
Windows'u etkinleştirmek için Ayarlar'a gidin.

## Precision/Recall

$y = 1$  in presence of rare class that we want to detect

Actual class

Predicted 1 class

	1	0
1	•	
0		

Windows'u Etkinleştir  
Windows'u etkinleştirmek için Ayarlar'a gidin.



# Precision/Recall

$y = 1$  in presence of rare class that we want to detect

Actual class

1 ←      0

Predicted 1 class

→ 0

True positive	False positive
False negative	True negative

Windows'u Etkinleştir  
Windows'u etkinleştirmek için Ayarlar'a gidin.

# Precision/Recall

$y = 1$  in presence of rare class that we want to detect

		Actual class	
		1	0
Predicted class	1	True positive	False positive
	0	False negative	True negative

## → Precision

(Of all patients where we predicted  $y = 1$ , what fraction actually has cancer?)

## Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

# Precision/Recall

$y = 1$  in presence of rare class that we want to detect

		Actual class	
		1	0
Predicted class	1	True positive	False positive
	0	False negative	True negative

## → Precision

(Of all patients where we predicted  $y = 1$ , what fraction actually has cancer?)

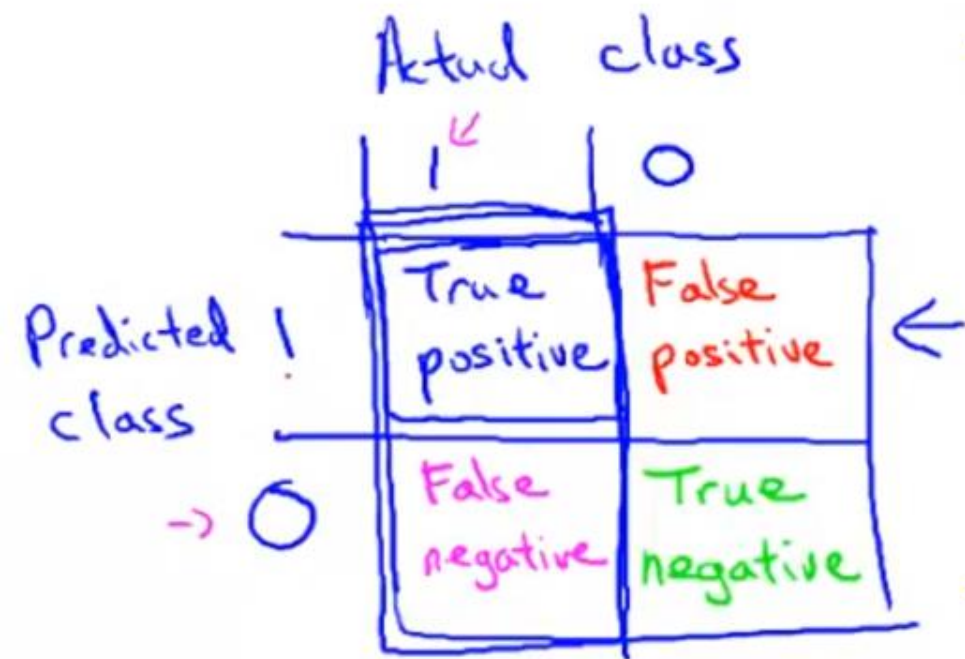
$$\frac{\text{True positives}}{\text{\# predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

## Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

# Precision/Recall

$y = 1$  in presence of rare class that we want to detect



## → Precision

(Of all patients where we predicted  $y = 1$ , what fraction actually has cancer?)

$$\frac{\text{True positives}}{\text{\#predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

## → Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

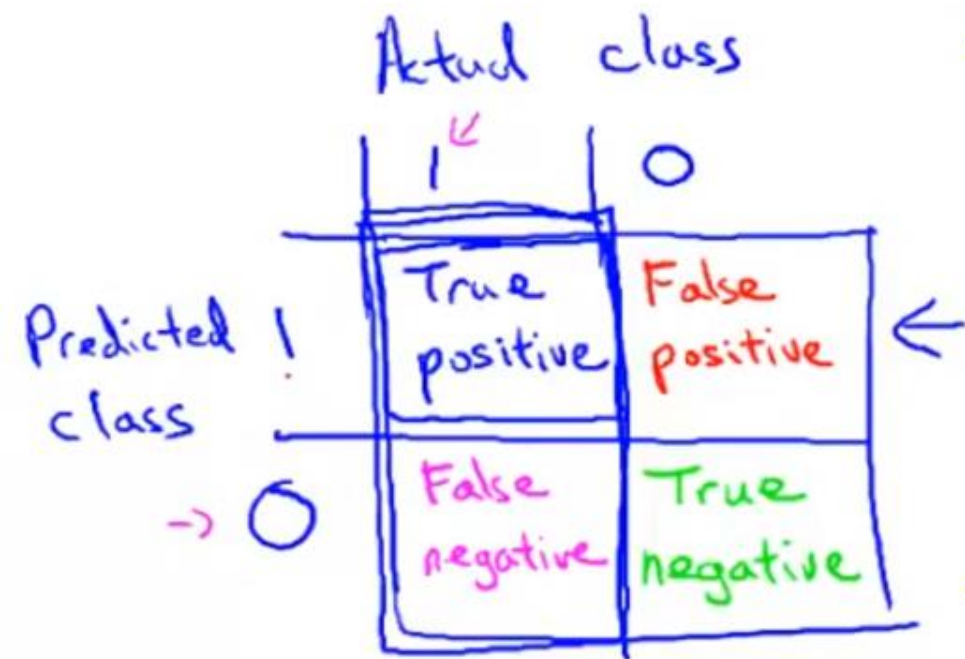
$$\frac{\text{True positives}}{\text{\#actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

Windows'u etkinleştirin  
Windows'u etkinleştirmek için Ayarlar'a gidin.



# Precision/Recall

$y = 1$  in presence of rare class that we want to detect



## → Precision

(Of all patients where we predicted  $y = 1$ , what fraction actually has cancer?)

$$\frac{\text{True positives}}{\text{\#predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

## → Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True positives}}{\text{\#actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

$y = 0$   
recall = 0

Windows'u etkinleştirin  
Windows'u etkinleştirmek için Ayarlar'a gidin.

# Example

- Your algorithm's performance on the test set is given to the right. What is the algorithm's precision and recall?

		Actual class	
		1	0
Predicted class	1	80	20
	0	80	820

# Precision/Recall

- Usually we use the convention that  $y$  is equal to 1, in the presence of the more rare class.
- So if we are trying to detect rare conditions such as cancer, precision and recall are defined setting  $y$  equals 1, rather than  $y$  equals 0.
- And by using precision and recall, we find, what happens is that even if we have very skewed classes, it's not possible for an algorithm to you know, "cheat" and predict  $y$  equals 1 all the time, or predict  $y$  equals 0 all the time, and get high precision and recall.
- And in particular, if a classifier is getting high precision and high recall, then we are actually confident that the algorithm has to be doing well, even if we have very skewed classes.