# Building Classification Models with scikit-learn

UNDERSTANDING CLASSIFICATION AS A MACHINE LEARNING PROBLEM

**Janani Ravi**
CO-FOUNDER, LOONYCORN

www.loonycorn.com

# Overview

Logistic regression for classification
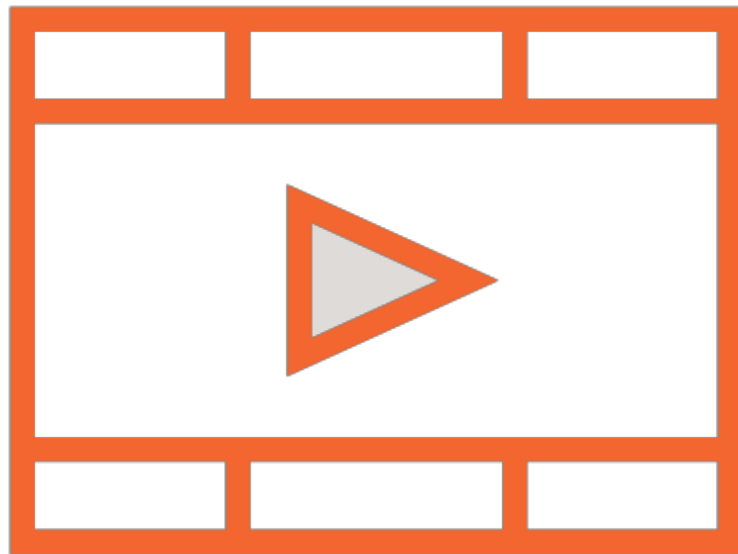
Evaluating classification models

Accuracy, precision, and recall

ROC curves

Binary, multi-label, and multi-class classification

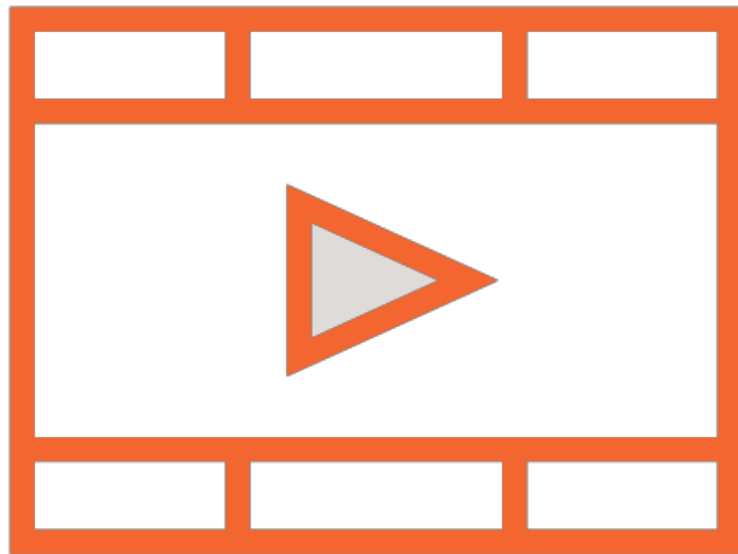# Prerequisites and Course Outline

# Prerequisites

**Basic Python programming**

**Basic understanding of the ML workflow**

**High school math**

# Prerequisite Courses

**Building Your First scikit-learn Solution**

# Course Outline

Understanding the classification problem
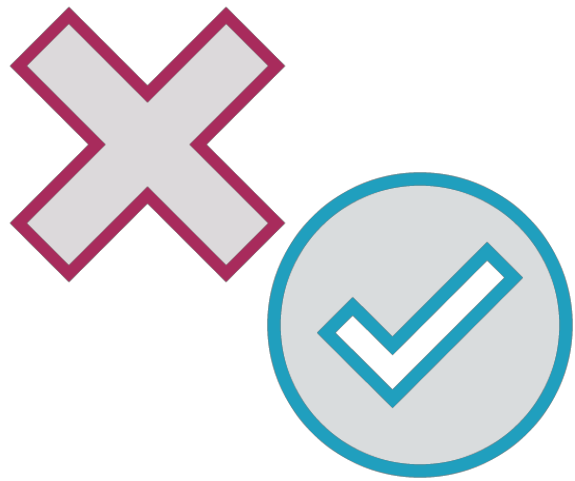
Building a simple ML classifier

Choosing and implementing classification technique
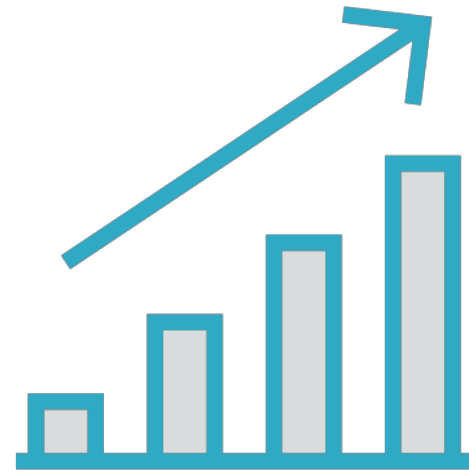
Hyperparameter tuning for classification

Classifying images
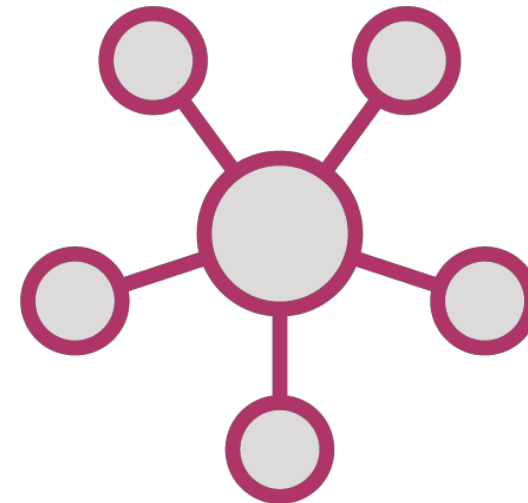
# Classification and Classifiers
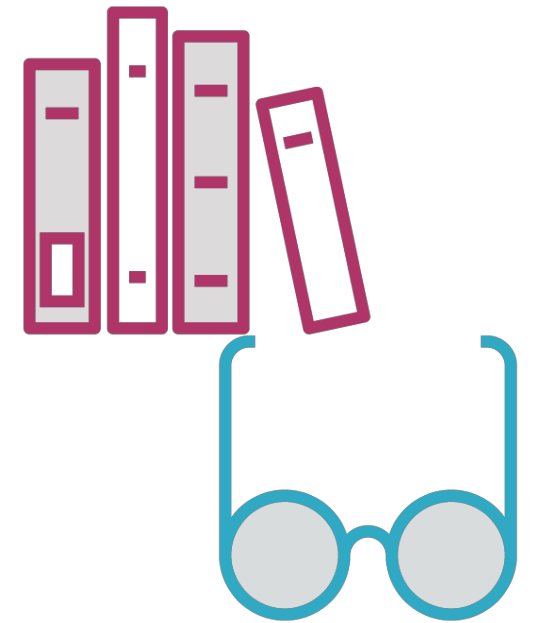
# Types of Machine Learning Problems



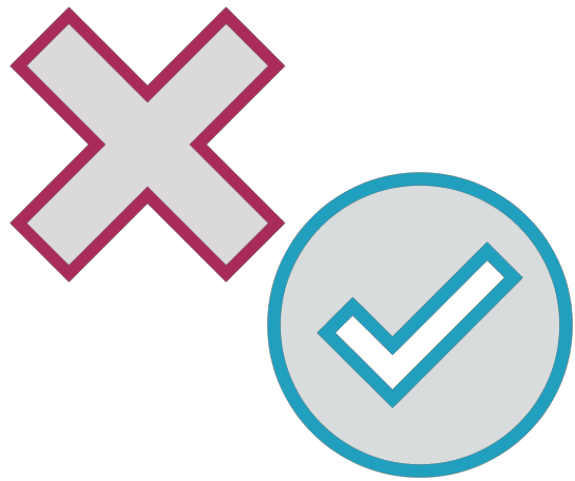**Classification**

**Regression**

**Clustering**

**Dimensionality reduction**

# Types of Machine Learning Problems



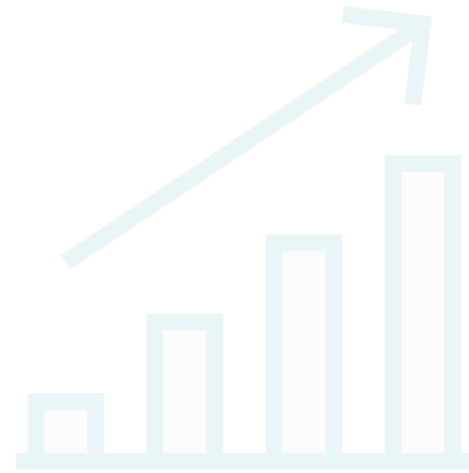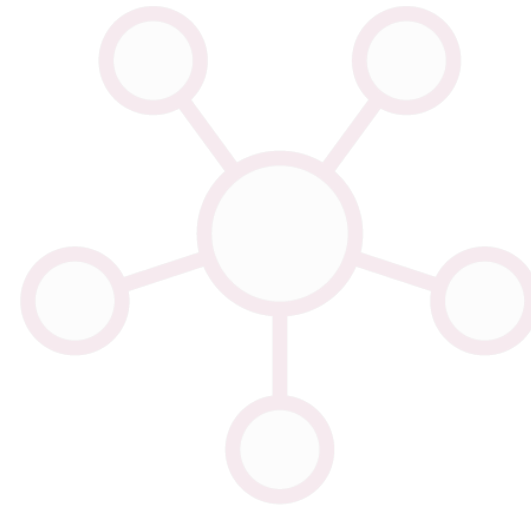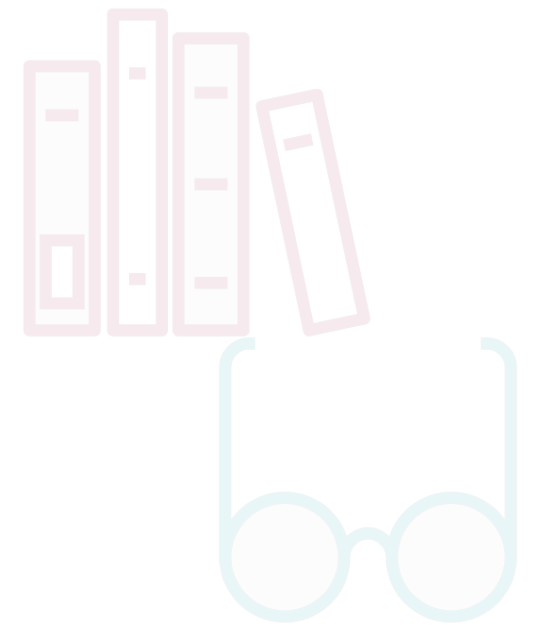**Classification**   Regression   Clustering   Dimensionality reduction

# Whales: Fish or Mammals?



**Mammals**

Members of the infraorder *Cetacea*

**Fish**

Look like fish, swim like fish, move with fish

# Whales: Fish or Mammals?



**ML-based Classifier**

# ML-based Classifier

## Training

Feed in a large corpus of data classified correctly

## Prediction

Use it to classify new instances which it has not seen before

# Training the ML-based Classifier

**Corpus**

**ML-based Classifier**

Classification

**Improves model parameters**

**Feedback - loss function or cost function**

# ML-based Binary Classifier

Breathes like a mammal

Gives birth like a mammal



**ML-based Classifier**

Mammal

**Corpus**

# "Traditional" ML-based Binary Classifier

Moves like a fish,
Looks like a fish



**ML-based Classifier**

Fish

**Corpus**

# ML-based Binary Classifier



**Corpus**

**Classification Algorithm**

**ML-based Classifier**

# ML-based Binary Classifier



**Corpus**

**Naive Bayes, Support Vector Machines, Decision Trees**

**ML-based Classifier**

# ML-based Binary Classifier

Breathes like a mammal

Gives birth like a
mammal

**ML-based Classifier**

Mammal

**Corpus**

# ML-based Binary Classifier

**Breathes like a mammal**

**Gives birth like a mammal**

Input: Feature Vector

ML-based Classifier

Mammal

Corpus

# ML-based Binary Classifier

Breathes like a mammal
Gives birth like a
mammal

**Mammal**

ML-based Classifier

Predicted Label

Corpus

# ML-based Binary Classifier

Breathes like a mammal
Gives birth like a
mammal

ML-based Classifier

Mammal

Predicted Label
=
Actual Label

Corpus

# ML-based Binary Classifier

**Moves like a fish,**

**Looks like a fish**

ML-based Classifier

Fish

Input: Feature Vector

Corpus

# ML-based Binary Classifier

Moves like a fish,
Looks like a fish

ML-based Classifier

Fish

Predicted Label
≠
Actual Label

Corpus

# Logistic Regression: Intuition

# Two Approaches to Deadlines

**Start 5 minutes before deadline**

Good luck with that

**Start 1 year before deadline**

Maybe overkill

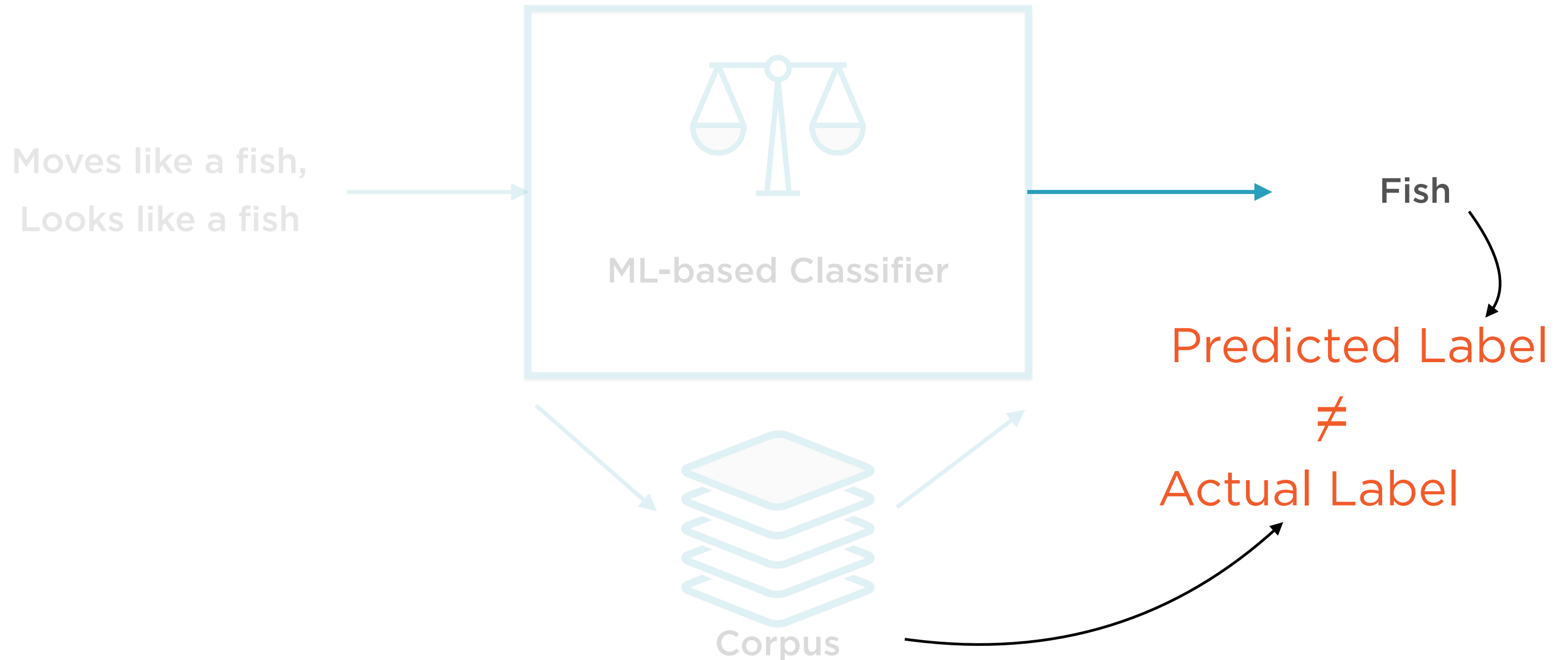## Neither approach is optimal

# Starting a Year in Advance

**Probability of meeting the deadline**

100%

**Probability of getting other important work done**

0%

# Starting Five Minutes in Advance

**Probability of meeting the deadline**

0%

..............................................................................................................................

**Probability of getting other important work done**

100%

# The Goldilocks Solution

## Work fast
Start very late and hope for the best

## Work smart
Start as late as possible to be sure to make it

## Work hard
Start very early and do little else

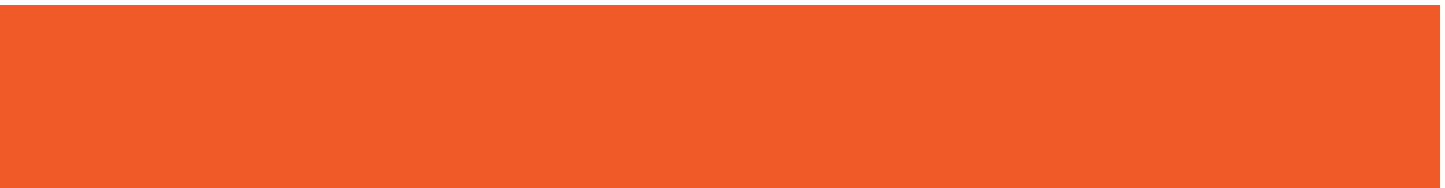**As usual, the middle path is best**

# Working Smart

**Probability of meeting the deadline**
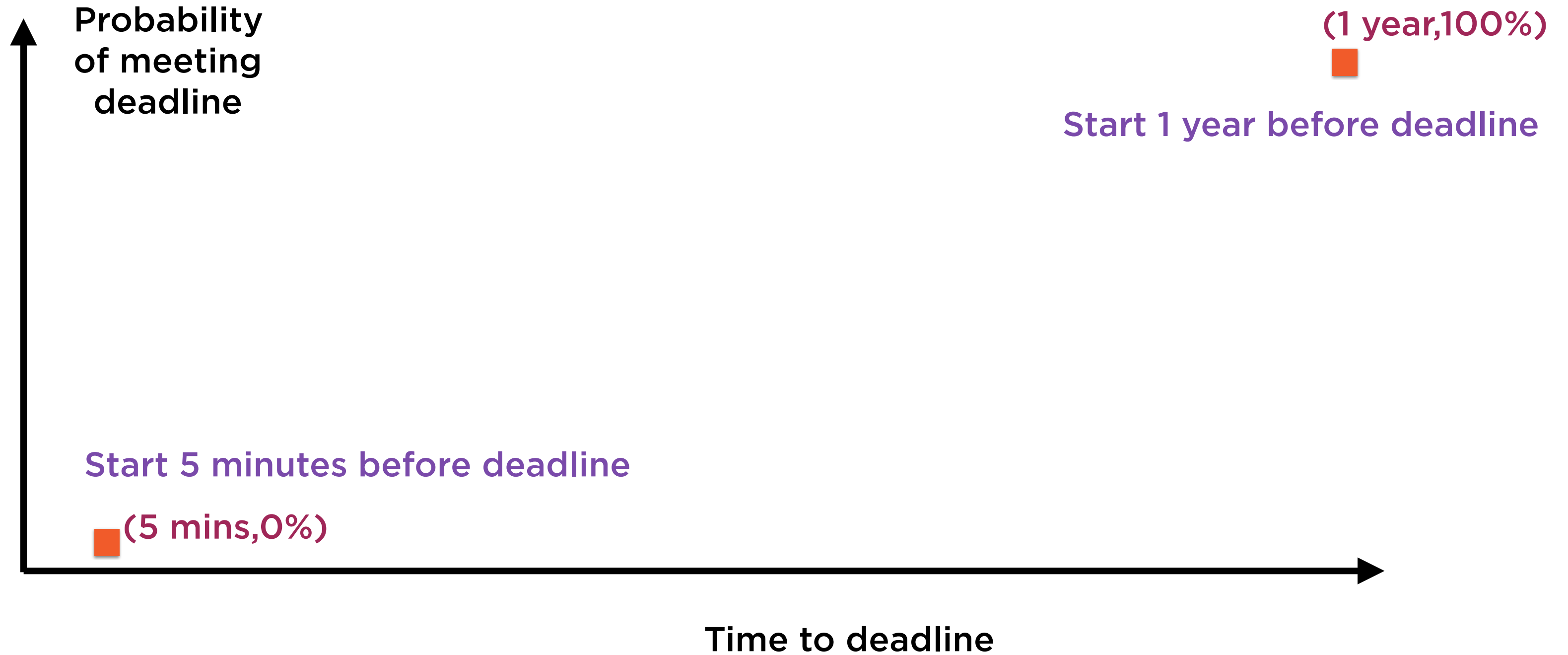
95%

**Probability of getting other important work done**

95%

# Working Hard, Fast, Smart

Probability of meeting deadline

**(1 year,100%)**

Start 1 year before deadline

Start 5 minutes before deadline

**(5 mins,0%)**

Time to deadline

# Working Hard, Fast, Smart



Probability of meeting deadline

(1 year,100%)

Work hard

(?,95%)

Work smart

Work fast
(5 mins,0%)

Time to deadline

Working Hard, Fast, Smart

# Working Hard, Fast, Smart

Probability
of meeting
deadline

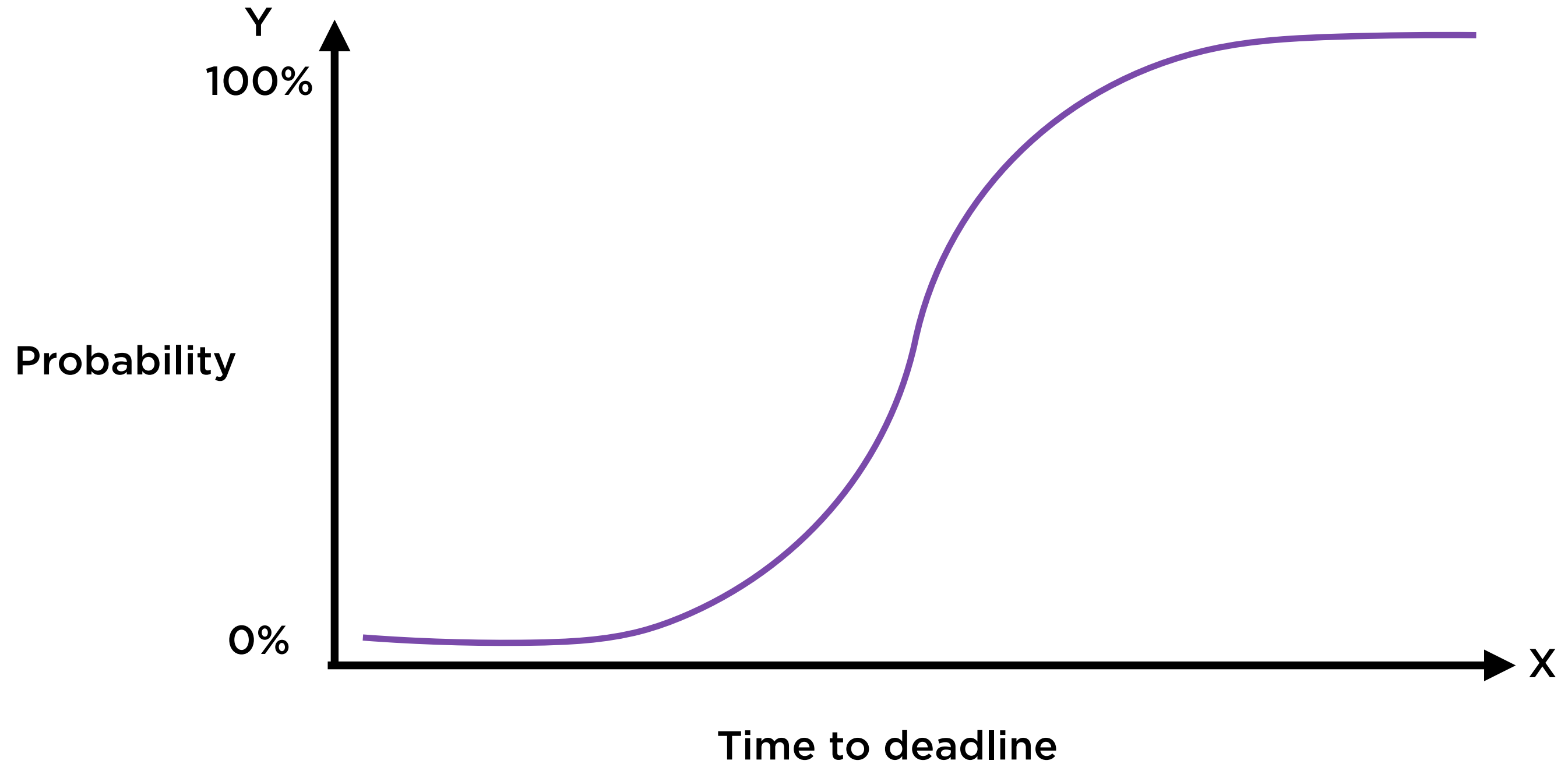Work hard

Work smart

Work fast

Time to deadline

Logistic Regression helps find how probabilities are changed by actions

# Working Smart with Logistic Regression

Working Smart with Logistic Regression

Start too late, and you'll definitely miss

# Working Smart with Logistic Regression

Y
100%

Probability

<50%    >50%

0%

X

**Time to deadline**

**Working smart is knowing when to start**

# Logistic Regression S-curves



y: hit or miss? (0 or 1?)

x: start time before deadline

p(y) : probability of y = 1

# Logistic Regression S-curves

$$p(y_i) = \frac{1}{1 + e^{-(A+Bx_i)}}$$
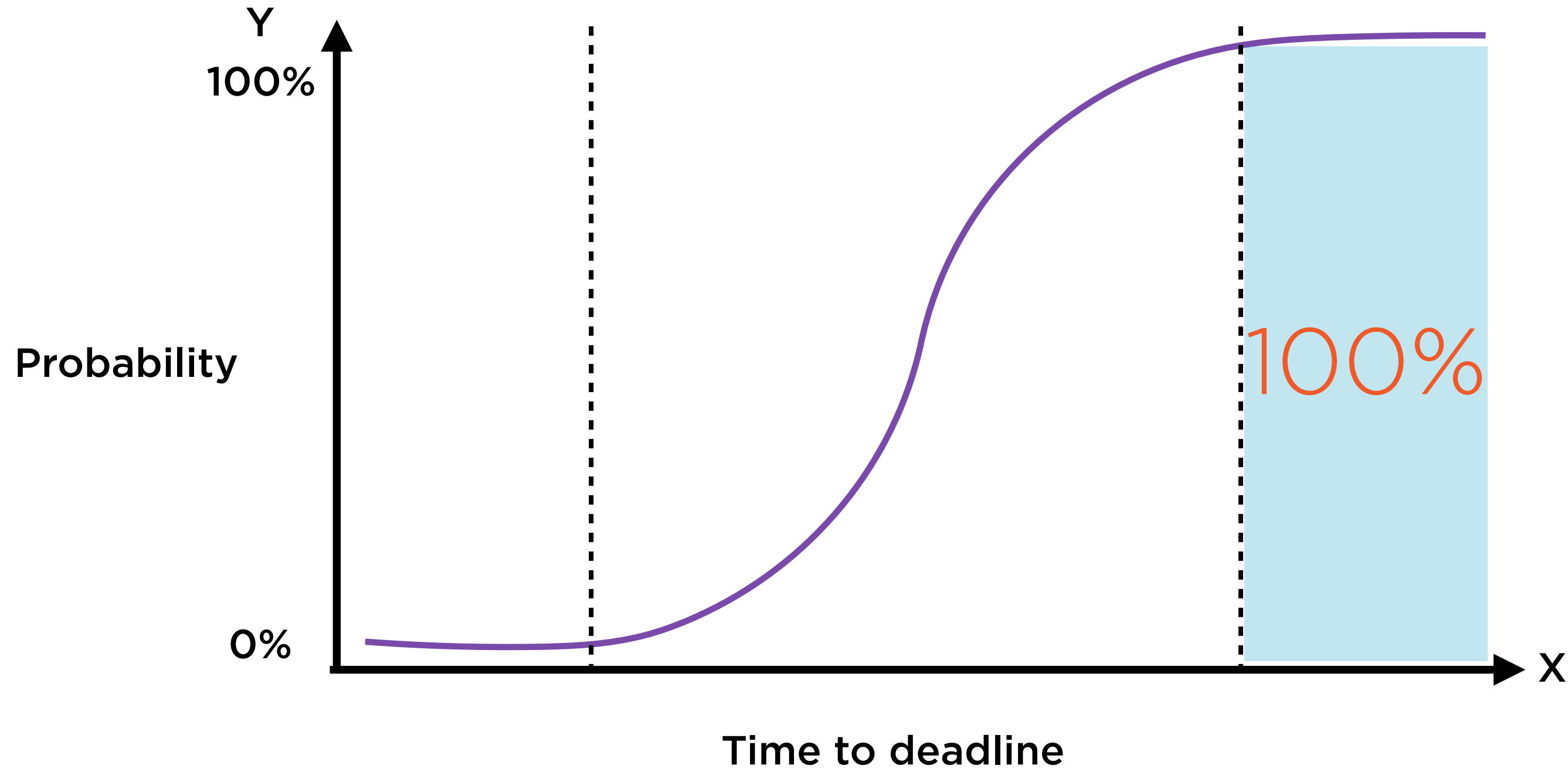
**Logistic regression involves finding the "best fit" such curve**

- A is the intercept

- B is the regression coefficient

*(e is the constant 2.71828)*

# Logistic Regression S-curves

S-curves are widely studied, well understood

Logistic regression uses S-curve to estimate probabilities

$$p(y) = \frac{1}{1 + e^{-(A+Bx)}}$$

# Linear Regression



Regression Line:
y = A + Bx

Y

X

**Finding the best fit line through these points**

# Logistic Regression



Regression Curve

$$p(y) = \frac{1}{1 + e^{-(A+Bx)}}$$

**Finding the best fit S-curve through these points**

# Logistic Regression

**Regression Equation:**

$$p(y_i) = \frac{1}{1 + e^{-(A+Bx_i)}}$$

**Solve for A and B that "best fit" the data**

# Cross Entropy: Loss Function

# Linear Regression Cost Function



Regression Line:
**y = A + Bx**

**The mean square error measures how far
the line is from the actual points**

# Logistic Regression Cost Function



$$p(y) = \frac{1}{1 + e^{-(A+Bx)}}$$

**Cross entropy measures how well the
estimated probabilities match actual labels**

Intuition: Low Cross Entropy

$Y_{actual}$

$Y_{predicted}$

# Intuition: Low Cross Entropy



**Y**actual

**Y**predicted

**The labels of the two series are in-synch**

Intuition: High Cross Entropy

$Y_{actual}$

$Y_{predicted}$

# Intuition: High Cross Entropy



Y<sub>actual</sub>

Y<sub>predicted</sub>

**The labels of the two series are out-of-synch**

# Accuracy, Precision, Recall

# Accuracy

Compare predicted and actual labels

More matches = higher accuracy

High accuracy is good, but...

An algorithm might have high accuracy but still be a poor machine learning model

Its predictions are **useless**

# All-is-well Binary Classifier

Medical reports →

**Always classify as "normal"**

→ No Cancer

Here, accuracy for rare cancer may be 99.9999%, but...

# Accuracy

Some labels maybe much more **common/rare** than others

Such a dataset is said to be **skewed**

Accuracy is a poor evaluation metric here

# Confusion Matrix

**Predicted Labels**

**Actual Label**

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | **10 instances** | **4 instances** |
| **No Cancer** | **5 instances** | **1000 instances** |

# Confusion Matrix

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | 10 | 4 |
| **No Cancer** | 5 | 1000 |

Actual Label

# True Positive

Predicted Labels

Actual Label

|  | **Cancer** | **No Cancer** |
|---|---|---|
| **Cancer** | 10 | 4 |
| **No Cancer** | 5 | 1000 |

Actual Label = Predicted Label

# True Positive

Predicted Labels

Actual Label

| | Cancer | No Cancer |
|---|---|---|
| **Cancer** | **10** TP | 4 |
| **No Cancer** | 5 | **1000** |

Actual Label  =  Predicted Label

# False Positive

Predicted Labels

Actual Label

| | Cancer | No Cancer |
|---|---|---|
| **Cancer** | 10 | 4 |
| **No Cancer** | 5 | 1000 |

Actual Label ≠ Predicted Label

# False Positive

Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | 10 | 4 |
| **No Cancer** | **5** FP | 1000 |

Actual Label  ≠  Predicted Label

# True Negative

# True Negative

Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | 10 | 4 |
| No Cancer | 5 | 1000 **TN** |

Actual Label = Predicted Label

# False Negative



Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | 10 | **4** |
| **No Cancer** | 5 | 1000 |

Actual Label ≠ Predicted Label

# False Negative

Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | 10 | 4 **FN** |
| **No Cancer** | 5 | 1000 |

Actual Label  ≠  Predicted Label

# Confusion Matrix



**Predicted Labels**

**Actual Label**

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | **TP** 10 | **FN** 4 |
| No Cancer | **FP** 5 | **TN** 1000 |

# Accuracy

Predicted Labels

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | **TP** 10 | **FN** 4 |
| **No Cancer** | **FP** 5 | **TN** 1000 |

Actual Label

# Accuracy

Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | TP 10 | FN 4 |
| No Cancer | FP 5 | TN 1000 |

Actual Label = Predicted Label

# Accuracy

Predicted Labels

Actual Label

| | Cancer | No Cancer |
|---|---|---|
| Cancer | **TP** 10 | **FN** 4 |
| No Cancer | **FP** 5 | **TN** 1000 |

$$\text{Accuracy} = \frac{TP + TN}{\text{Num Instances}} = \frac{1010}{1019} = 99.12\%$$

# Accuracy

Accuracy = 99.12%

Classifier gets it right 99.12% of the time

But...

# Accuracy

**Predicted Labels**

**Actual Label**

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | TP 10 | FN 4 |
| **No Cancer** | FP 5 | TN 1000 |

People on chemotherapy, radiation when not required

# Accuracy

Predicted Labels

Actual Label

| | Cancer | No Cancer |
|---|---|---|
| Cancer | TP 10 | FN 4 |
| No Cancer | FP 5 | TN 1000 |

Cancer not detected, no treatment prescribed

Accuracy is not a good metric to evaluate whether this model performs well

# Precision



**Predicted Labels**

**Actual Label**

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | **TP** 10 | **FN** 4 |
| No Cancer | **FP** 5 | **TN** 1000 |

# Precision

Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | TP 10 | FN 4 |
| No Cancer | FP 5 | TN 1000 |

Precision = Accuracy when classifier flags cancer

# Precision

## Predicted Labels

## Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | 10 **TP** | 4 **FN** |
| **No Cancer** | 5 **FP** | 1000 **TN** |

$$\text{Precision} \quad = \quad \frac{\text{TP}}{\text{TP + FP}} \quad = \quad \frac{10}{15} \quad = \quad 66.67\%$$

# Precision

**Precision = 66.67%**

**1 in 3 cancer diagnoses is incorrect**

# Recall

## Predicted Labels

|  | Cancer | No Cancer |
|---|---|---|
| **Cancer** | **TP** 10 | **FN** 4 |
| **No Cancer** | **FP** 5 | **TN** 1000 |

**Actual Label**

# Recall

Predicted Labels

Actual Label

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | **TP** 10 | **FN** 4 |
| No Cancer | **FP** 5 | **TN** 1000 |

Recall = Accuracy when cancer actually present

# Recall

|  | Cancer | No Cancer |
|---|---|---|
| Cancer | **TP** 10 | **FN** 4 |
| No Cancer | **FP** 5 | **TN** 1000 |

Actual Label

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{14} = 71.42\%$$

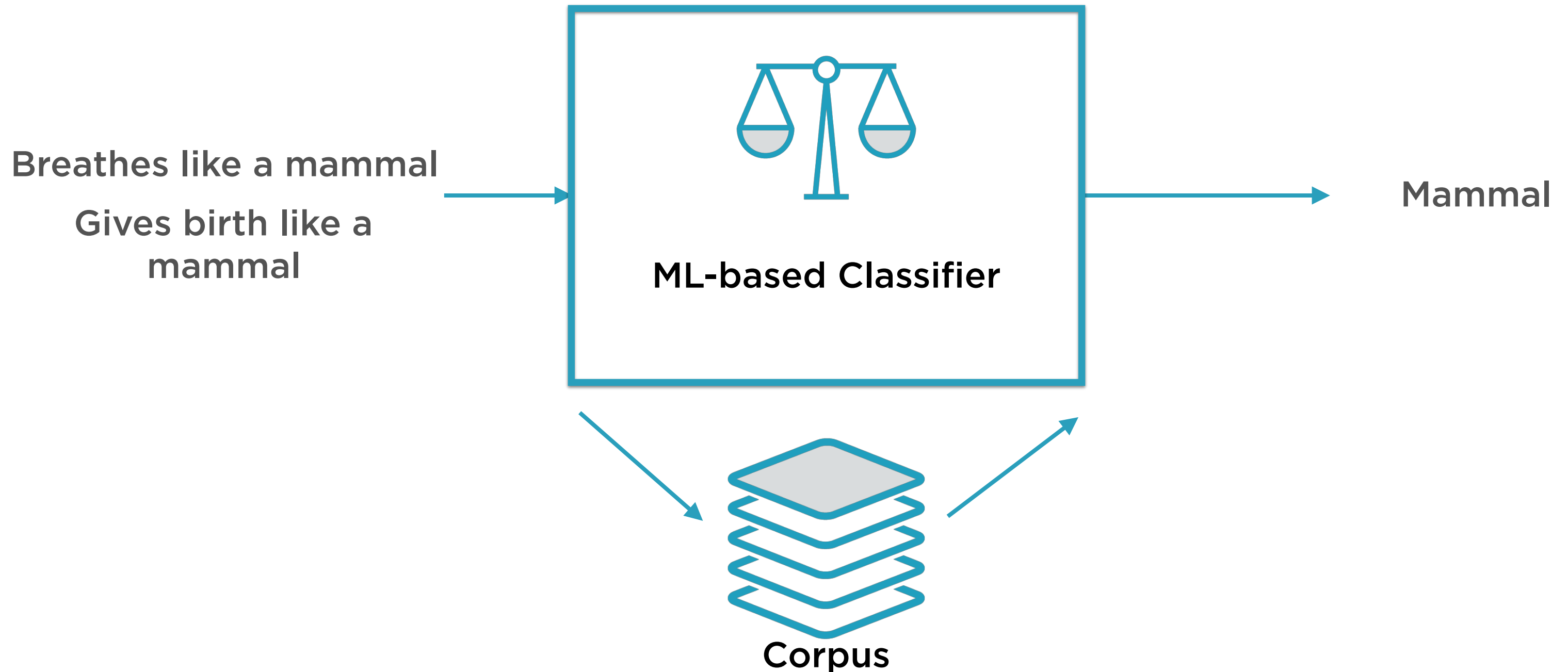# Recall

**Recall = 71.42%**
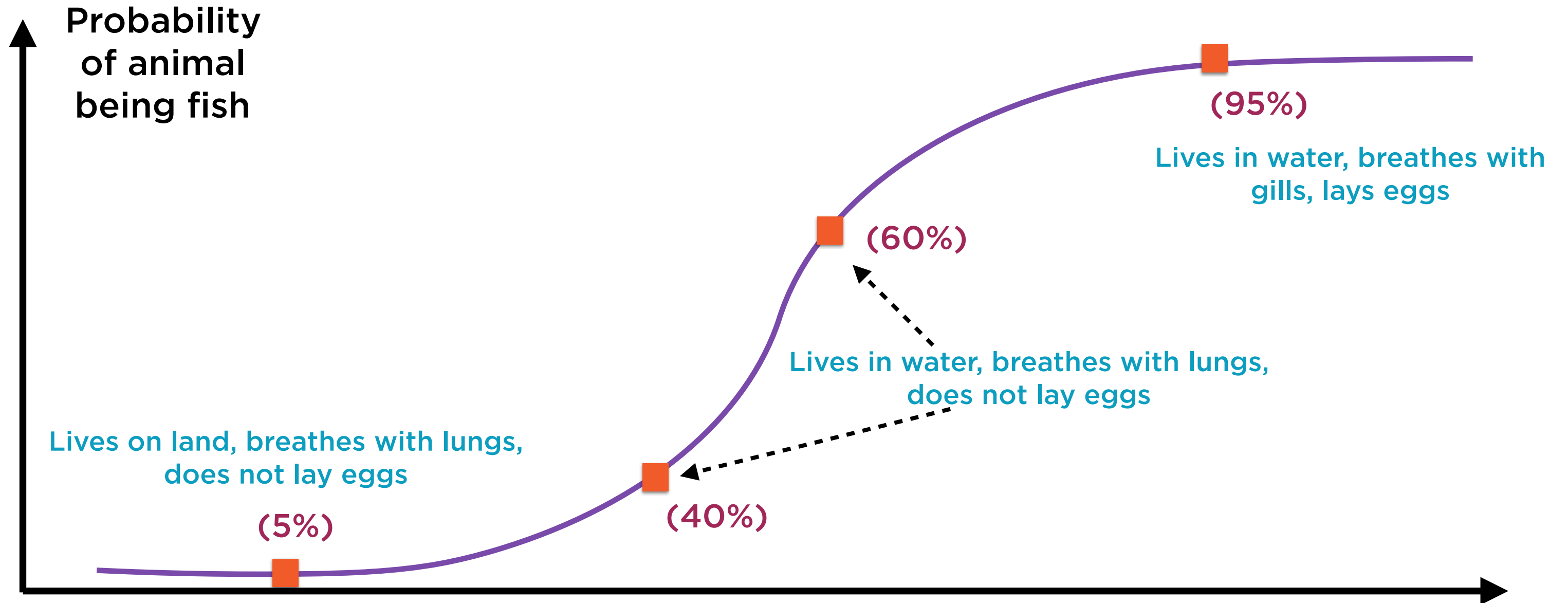
**2 in 7 cancer cases missed**

# Evaluating Classifiers

# ML-based Binary Classifier

Breathes like a mammal

Gives birth like a
mammal

**ML-based Classifier**

**Corpus**

Mammal

# ML-based Binary Classifier

Breathes like a mammal

Gives birth like a
mammal

**ML-based Classifier**

P(fish) = 0.45

**Corpus**

# Applying Logistic Regression



**Probability of animal being fish**

Lives on land, breathes with lungs, does not lay eggs

**(5%)**

**(40%)**

Lives in water, breathes with lungs, does not lay eggs

**(60%)**

Lives in water, breathes with gills, lays eggs

**(95%)**

**Whales: Fish or Mammals?**

# Choosing Decision Threshold

**Probability of animal being fish**

**P**<sub>**threshold**</sub>

(50%)

(5%)

(20%)

(40%)

(60%)

(80%)

(95%)

# Choosing Decision Threshold



**Probability of animal being fish**

$P_{threshold}$

(5%)

(20%)

(40%)

(60%)

(80%)

(95%)

**If probability < $P_{threshold}$, it's a mammal**

# Applying Logistic Regression

Probability of animal being fish

$P_{threshold}$

(95%)

(80%)

(60%)

(40%)

(20%)

(5%)

If probability > $P_{threshold}$, it's a fish

# Recall vs. "Conservativeness"
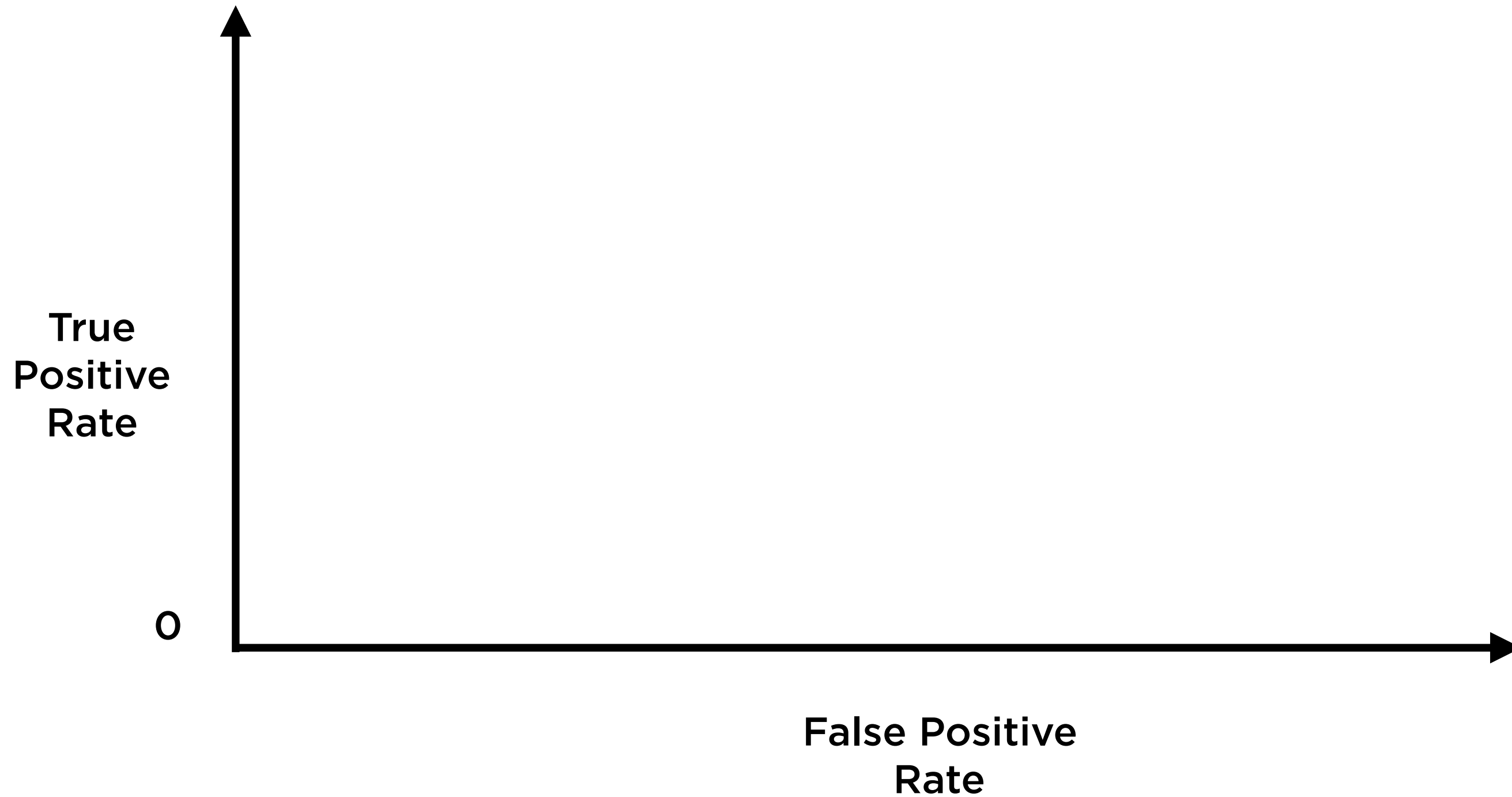


1.0

0

Recall

1.0

"Conservativeness" of Decision Threshold

# Precision-Recall Tradeoff

# Precision-Recall Tradeoff

Choosing P<sub>threshold</sub>

True Positive Rate

False Positive Rate

# Choosing P~threshold~
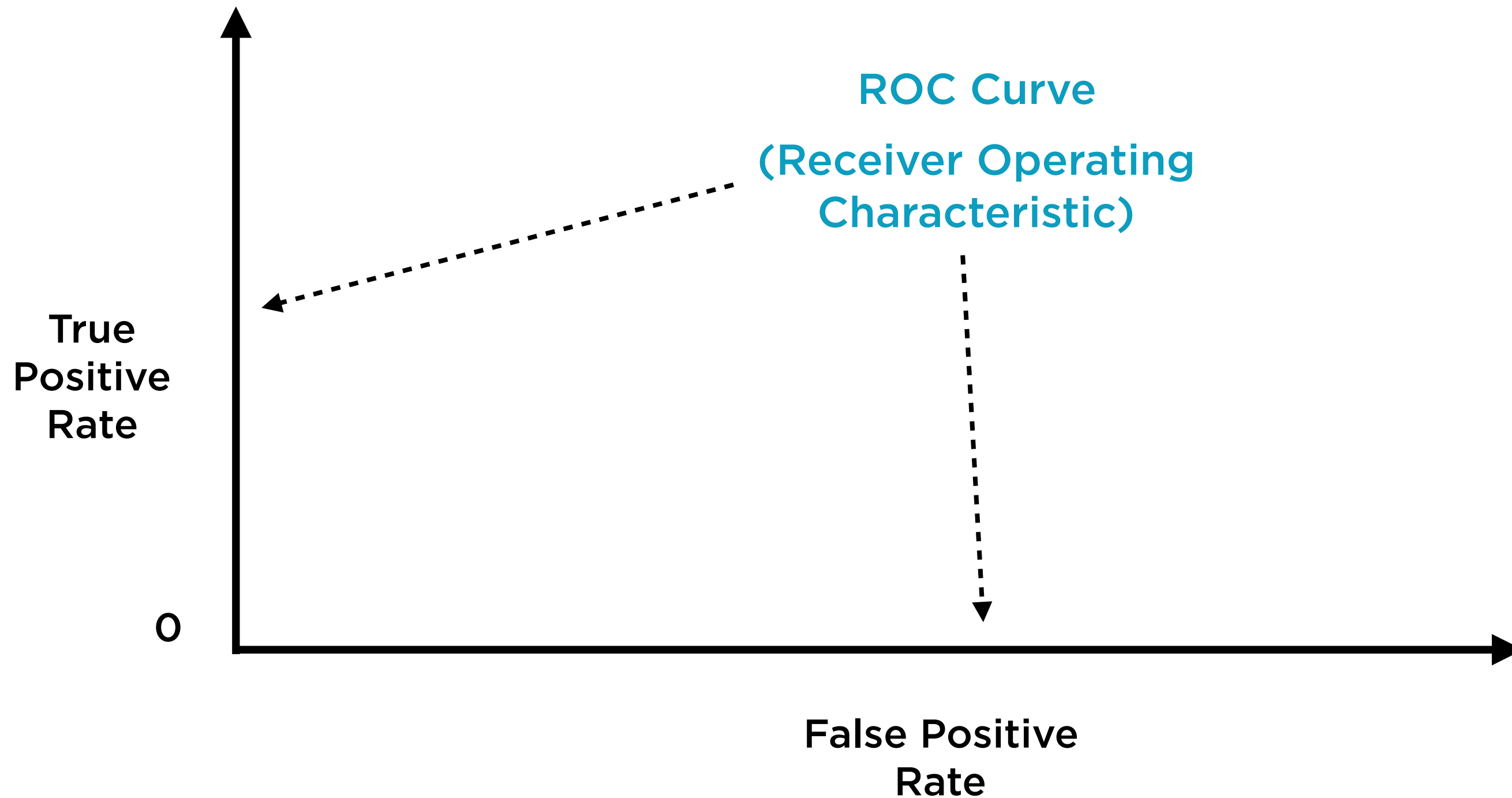
**True Positive Rate**

**Should be as high as possible**

0

False Positive Rate

# Choosing P~threshold~



**True Positive Rate** (y-axis)

0

**Should be as low as possible**

**False Positive Rate** (x-axis)

# Choosing P<sub>threshold</sub>

# Choosing P_threshold

True Positive Rate

1.0

0

False Positive Rate

Fit ROC curve from different values of P_threshold

# ROC Curve

True Positive Rate

1.0

0

**Pick top-left corner point as P$_{threshold}$**

**Why? Maximises True Positive Rate, minimizes False Positive Rate**

False Positive Rate

# ROC of Perfect Classifier



1.0

TP = 100%

FP = 0%

True
Positive
Rate

0

False Positive
Rate

ROC of Random Classifier

# Types of Classification

# Types of Classification Tasks

**Binary**

**"Yes/No", "True/False", "Up/Down"**

Output is binary categorical variable

**Multilabel**

**("True", "Female"), ("False", "Female")**

Output is tuple of multiple binary variables (not disjoint)

**Multiclass**

**Digit classification**

Output variable takes 1 of N (>2) values

**Multioutput**

**("Sunday", "January")**

Multiclass + multilabel

# Multilabel



**Some algorithms are inherently multilabel**
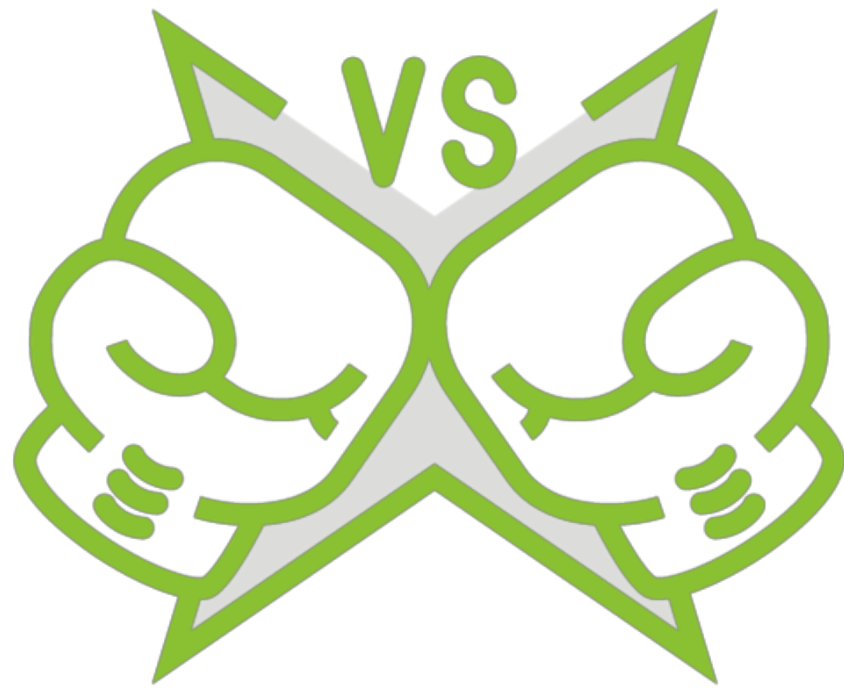
- Naive Bayes

# Multilabel



**Many classification algorithms are inherently binary**

- Logistic regression

- Support Vector Machines

**Inherently binary classifiers can be generalised for multilabel classification**
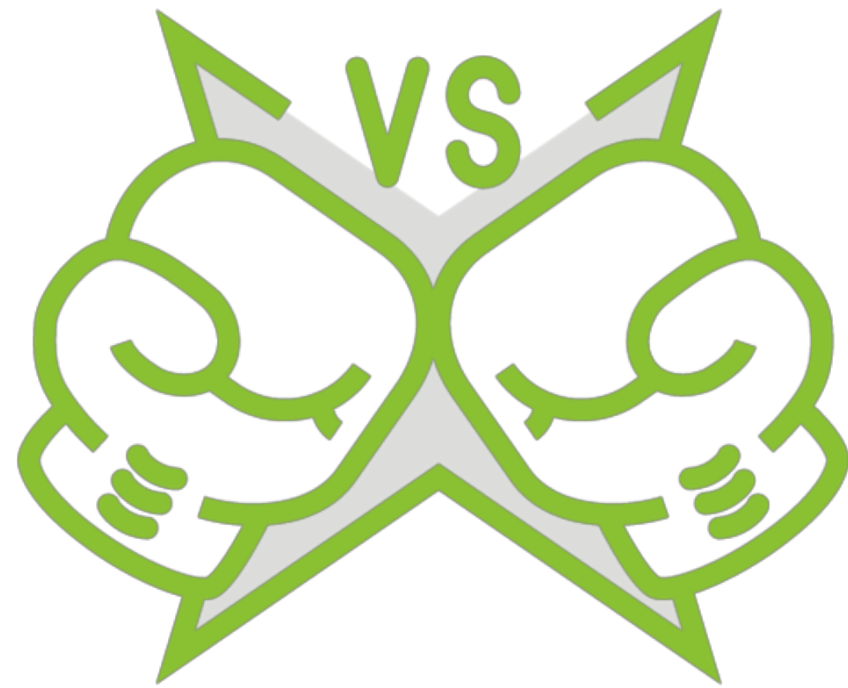
# One vs. All



**One-versus-all**

**Classifying digits 0-9**

**Train 10 binary classifiers**

- 0-detector, 1-detector...

- Predicted label = output of detector with highest score

# One vs. One

**One-versus-one**

**Train 45 binary classifiers**

- One detector for each pair of digits

- For N labels, need N(N-1)/2 classifiers

- Predicted label = output of digit that wins most duels

# Summary

Logistic regression for classification

Evaluating classification models

Accuracy, precision, and recall

ROC curves

Binary, multi-label, and multi-class classification