

Lecture 2

Tian Han

Outline

- Polynomial Regression
- Binary Classification - Logistic regression

Warm-up: Linear Regression

Linear Regression (Task)

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$

Output: a vector $\mathbf{w} \in \mathbb{R}^d$ and scalar $b \in \mathbb{R}$ such that $\mathbf{x}_i^T \mathbf{w} + b \approx y_i$.

Tasks

Linear
Regression

assume y_i is a linear
function of \mathbf{x}_i .



Least Squares Regression (Method)

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$

1. Add one dimension to $\mathbf{x} \in \mathbb{R}^d$: $\bar{\mathbf{x}}_j = [\mathbf{x}_j; 1] \in \mathbb{R}^{d+1}$.
2. Solve least squares regression: $\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \|\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}\|_2^2$.

Tasks

Linear
Regression

Methods

Least Squares Regression

Least Squares Regression (Method)

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$

1. Add one dimension to $\mathbf{x} \in \mathbb{R}^d$: $\bar{\mathbf{x}}_j = [\mathbf{x}_j; 1] \in \mathbb{R}^{d+1}$.
2. Solve least squares regression: $\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \|\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}\|_2^2$.

Tasks

Linear
Regression

Methods

Least Squares Regression

Algorithms

Analytical Solution

Gradient Descent

Polynomial Regression

The Regression Task

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Output: a function $f: \mathbb{R}^d \mapsto \mathbb{R}$ such that $f(\mathbf{x}) \approx y$.

Question: f is unknown! So how to learn f ?

The Regression Task

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Output: a function $f: \mathbb{R}^d \mapsto \mathbb{R}$ such that $f(\mathbf{x}) \approx y$.

Question: f is unknown! So how to learn f ?

Answer: polynomial approximation; f is a polynomial function.

Taylor expansion: $f(x) = f(a) + f'(a)(a - x) + \frac{f''(a)}{2!}(a - x)^2 + \dots$

Polynomial Regression: 1D Example

Input: scalars $x_1, \dots, x_n \in \mathbb{R}$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Output: a function $f: \mathbb{R} \mapsto \mathbb{R}$ such that $f(x) \approx y$.

One-dimensional example: $f(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_p x^p$.

Polynomial Regression: 1D Example

Input: scalars $x_1, \dots, x_n \in \mathbb{R}$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Output: a function $f: \mathbb{R} \mapsto \mathbb{R}$ such that $f(x) \approx y$.

One-dimensional example: $f(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_p x^p$.

Polynomial regression:

1. Define a feature map $\Phi(x) = [1, x, x^2, x^3, \dots, x^p]$.
2. For $j = 1$ to n , do the mapping $x_j \mapsto \Phi(x_j)$.
 - Let $\Phi = [\Phi(x_1); \dots, \Phi(x_n)]^T \in \mathbb{R}^{n \times (p+1)}$
3. Solve the least squares regression $\min_{\mathbf{w} \in \mathbb{R}^{p+1}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2$.

Polynomial Regression: 2D Example

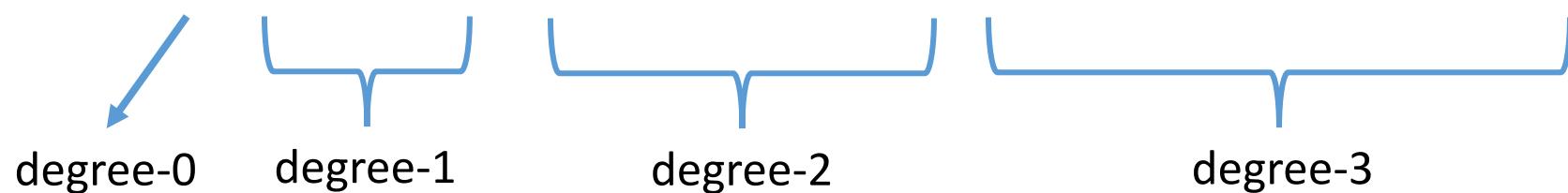
Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^2$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Output: a function $f: \mathbb{R}^2 \mapsto \mathbb{R}$ such that $f(\mathbf{x}_i) \approx y_i$.

Two-dimensional example: how to do feature mapping?

Polynomial features:

$$\Phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2].$$



Polynomial Regression

```
import numpy  
X = numpy.arange(6).reshape(3, 2)  
print('X = ')  
print(X)
```

```
X =  
[[0 1]  
 [2 3]  
 [4 5]]
```

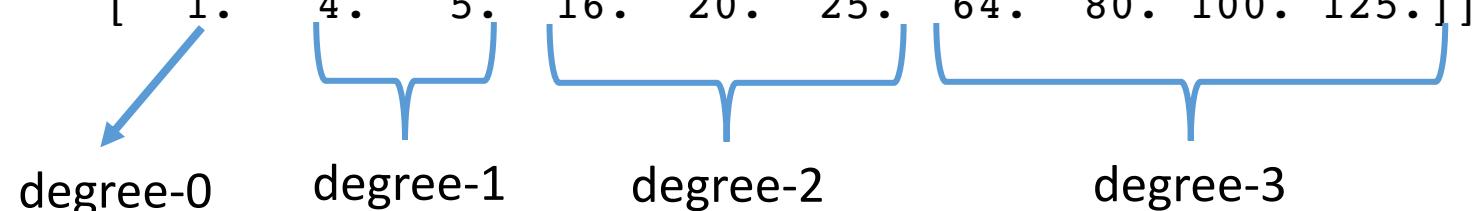
Polynomial Regression

```
import numpy  
X = numpy.arange(6).reshape(3, 2)  
print('X = ')  
print(X)
```

```
X =  
[[0 1]  
 [2 3]  
 [4 5]]
```

```
from sklearn.preprocessing import PolynomialFeatures  
poly = PolynomialFeatures(degree=3)  
Phi = poly.fit_transform(X)  
print('Phi = ')  
print(Phi)
```

```
Phi =  
[[ 1.   0.   1.   0.   0.   1.   0.   0.   0.   1.]  
 [ 1.   2.   3.   4.   6.   9.   8.  12.  18.  27.]  
 [ 1.   4.   5.  16.  20.  25.  64.  80. 100. 125.]]
```



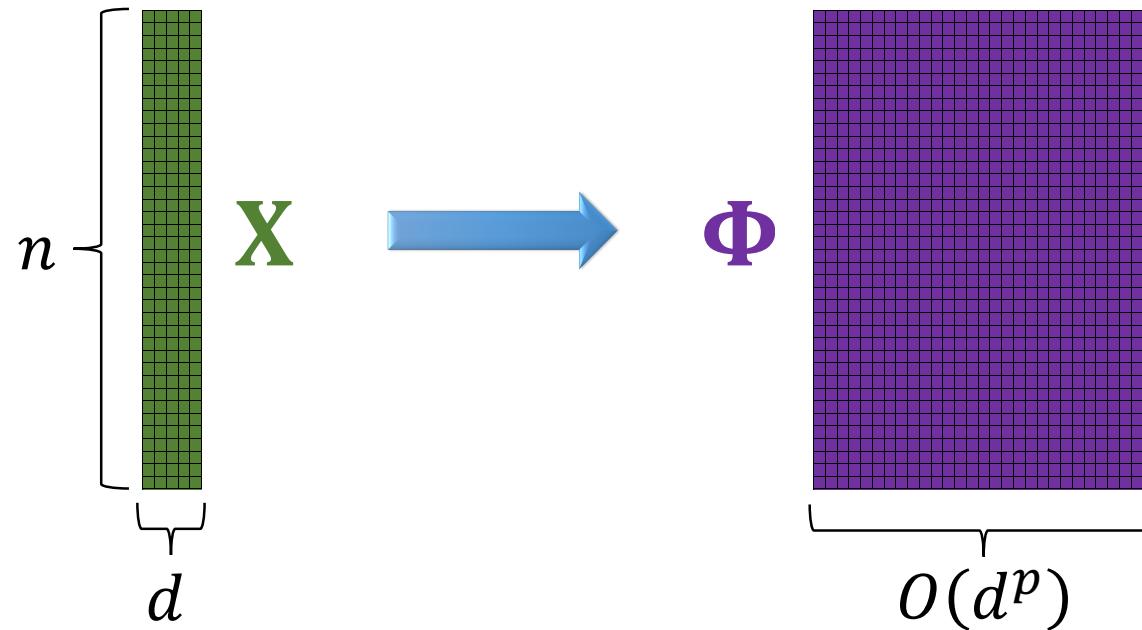
Polynomial Regression

- \mathbf{x} : d -dimensional
- $\Phi(\mathbf{x})$: degree- p polynomial
- The dimension of $\Phi(\mathbf{x})$ is $O(d^p)$

Polynomial Regression

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Output: a function $f: \mathbb{R}^d \mapsto \mathbb{R}$ such that $f(\mathbf{x}_i) \approx y_i$.



Training, Test, and Overfitting

Polynomial Regression: Training

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Feature map: $\phi(\mathbf{x})$. Its dimension is $O(d^p)$.

Least squares: $\min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2$.

Polynomial Regression: Training

Input: vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$.

Feature map: $\Phi(\mathbf{x})$. Its dimension is $O(d^p)$.

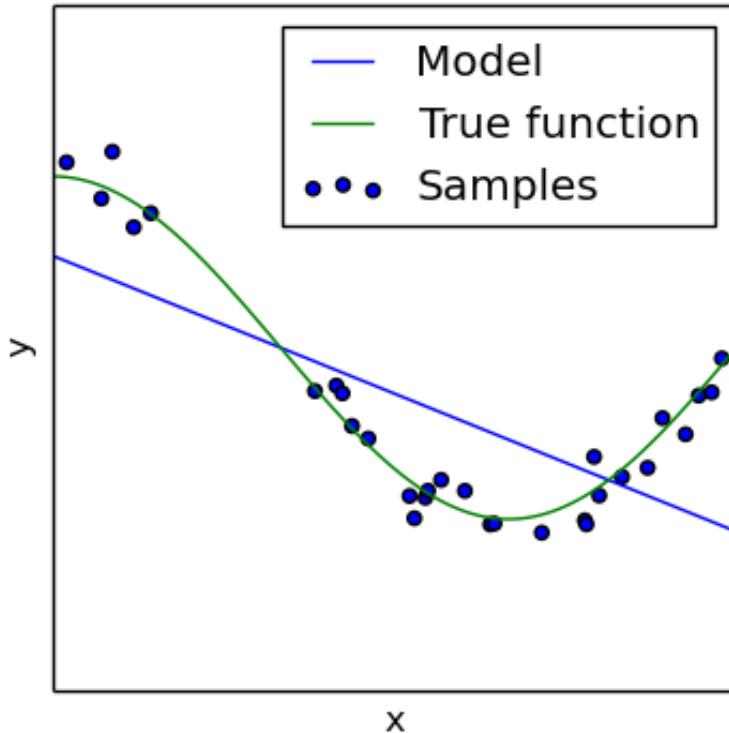
Least squares: $\min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2$.

Question: what will happen as p grows?

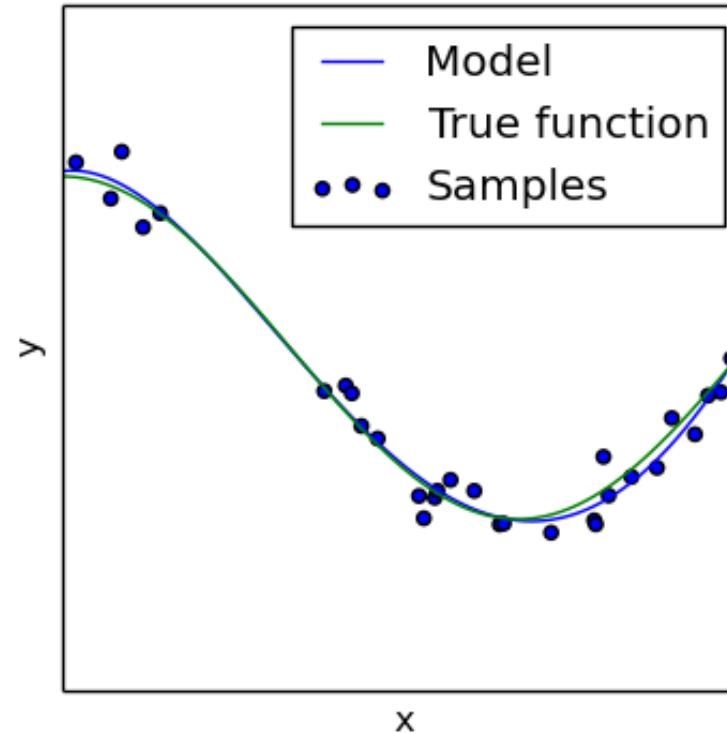
1. For sufficiently large p , the dimension of the feature $\Phi(\mathbf{x})$ exceeds n .
2. Then you can find \mathbf{w} such that $\Phi \mathbf{w} = \mathbf{y}$. (Zero training error!)

Training and Testing

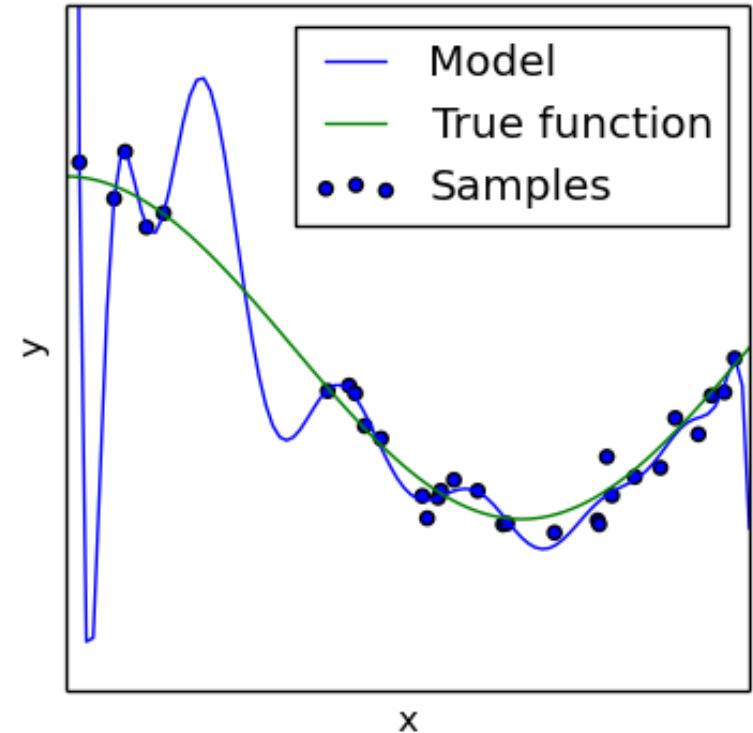
Degree 1



Degree 4



Degree 15



Underfitting

Overfitting

Training and Testing

Train:

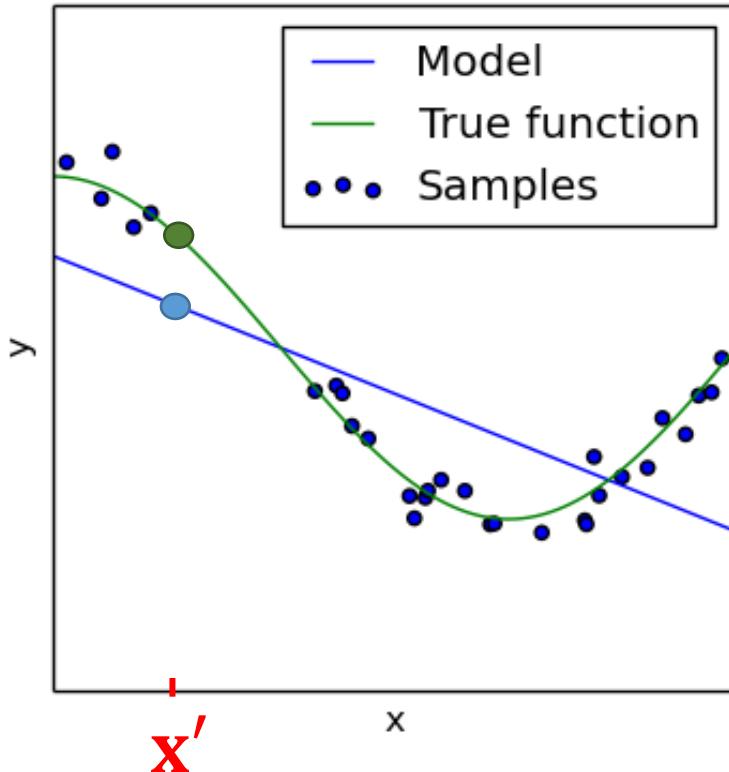
- Input:** vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \mathbb{R}$.
- Output:** a function $f: \mathbb{R}^d \mapsto \mathbb{R}$ such that $f(\mathbf{x}_i) \approx y_i$.

Test:

- Input:** a *never-seen-before* feature vectors $\mathbf{x}' \in \mathbb{R}^d$.
- output:** predict its label by $f(\mathbf{x}')$.

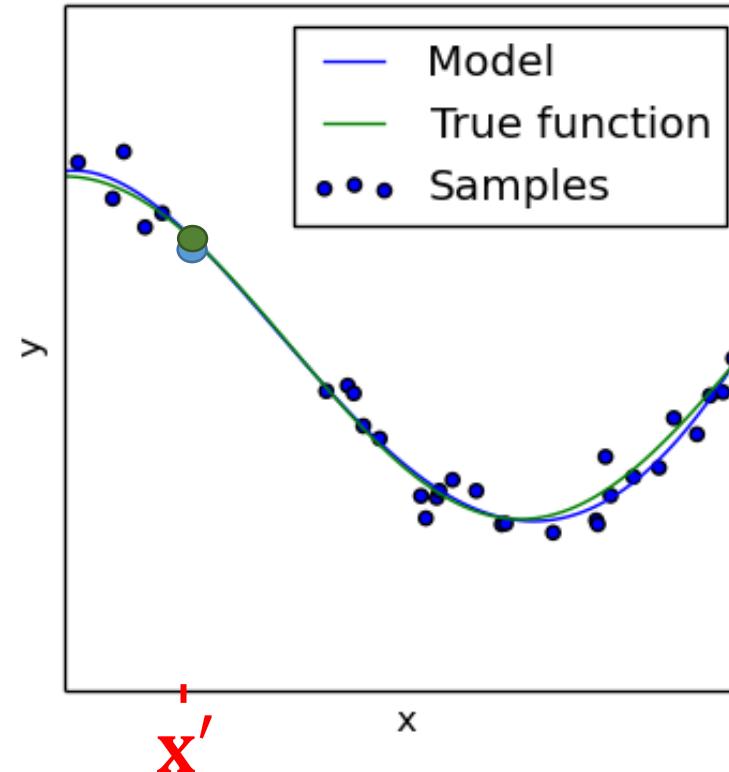
Training and Testing

Degree 1



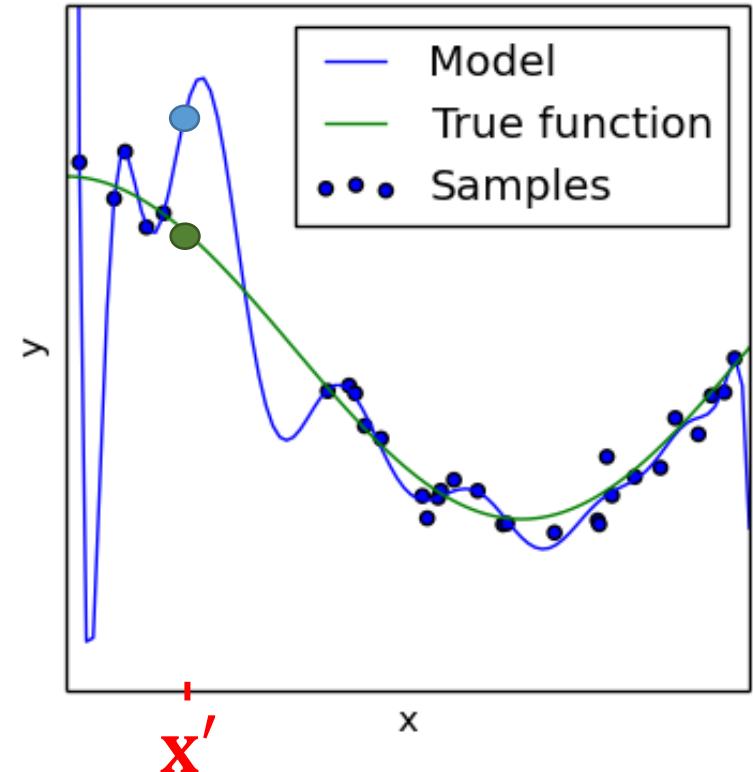
BAD

Degree 4



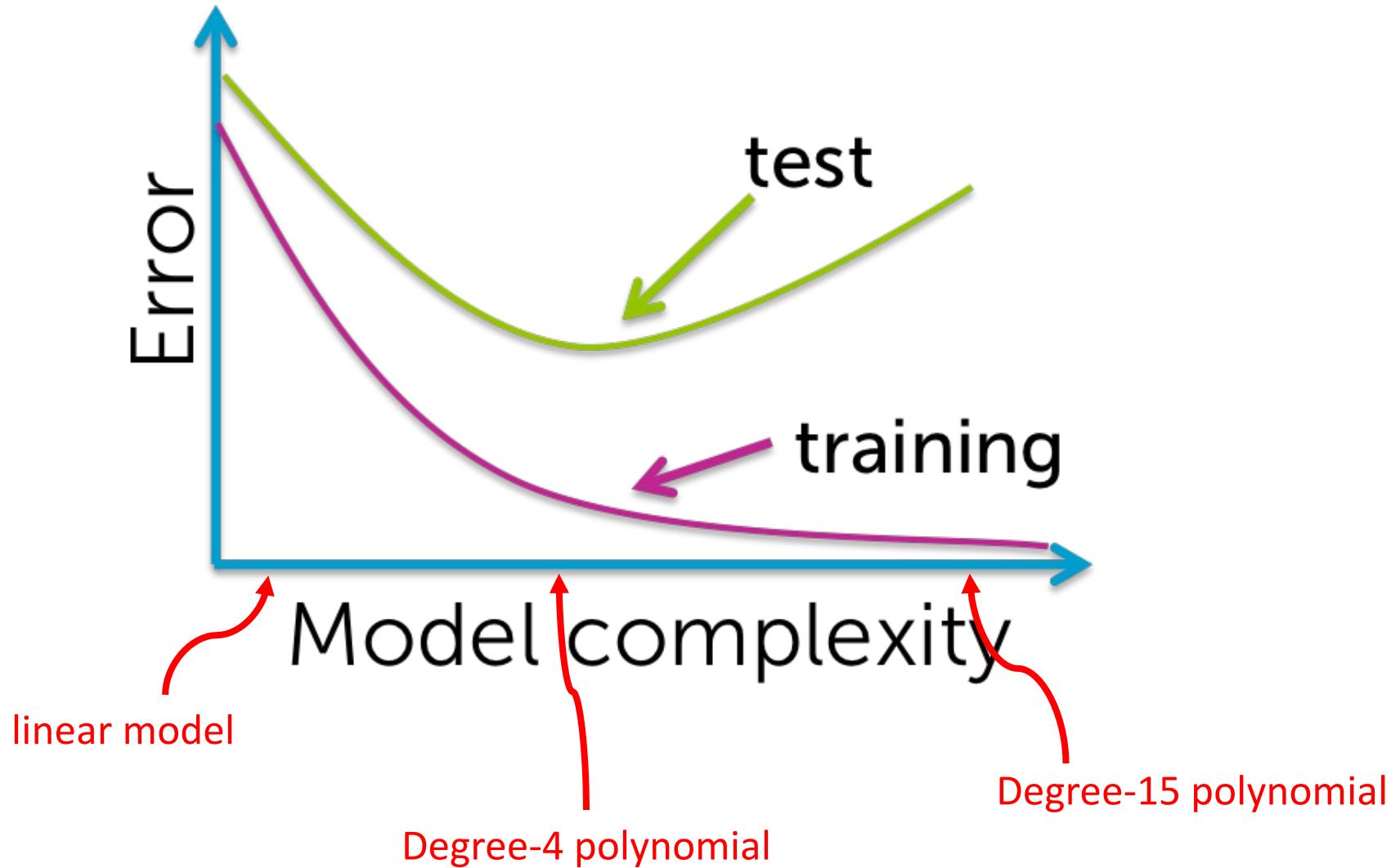
GOOD

Degree 15



BAD

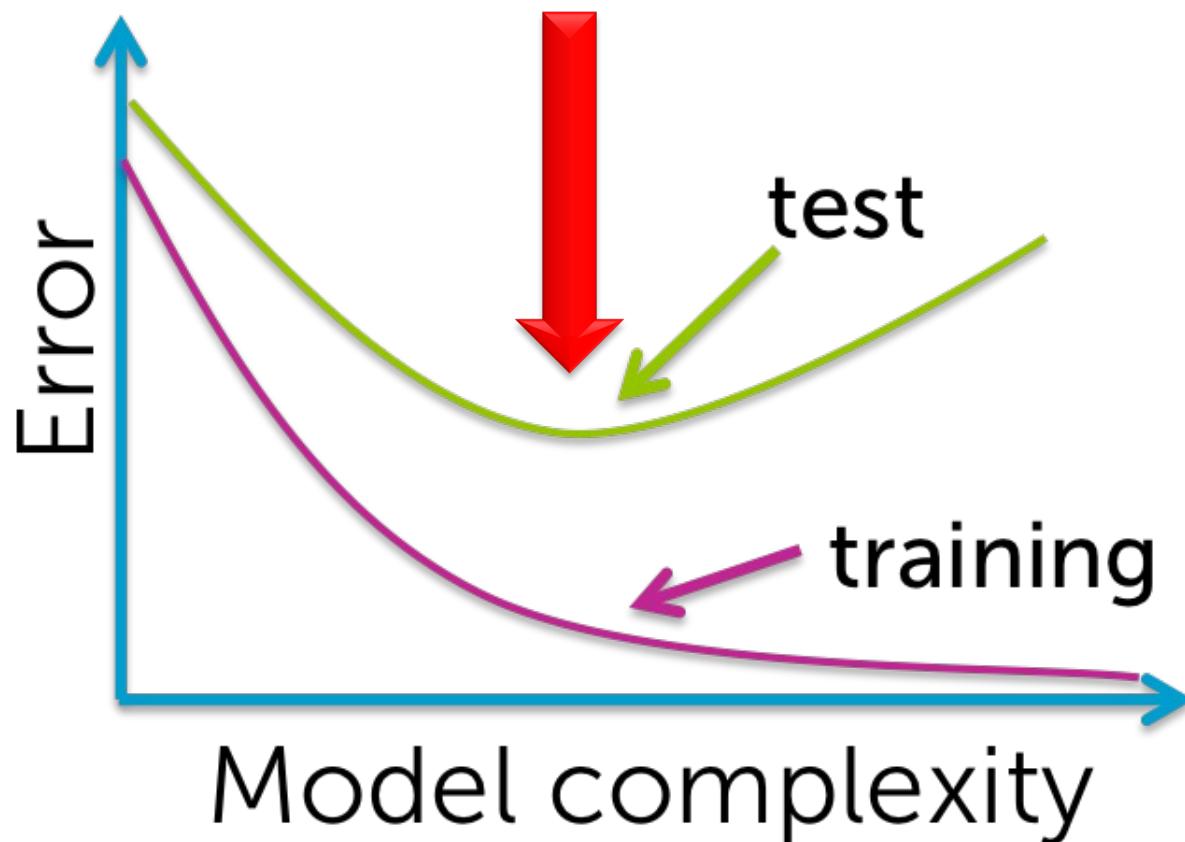
Training and Testing



Hyper-Parameter Tuning

Question: for the polynomial regression model, how to determine the degree p ?

Answer: the degree p leads to the smallest test error.



Hyper-Parameter Tuning

Training Set	Test Set
Train a degree-1 polynomial regression	Test MSE = 23.2
Train a degree-2 polynomial regression	Test MSE = 19.0
Train a degree-3 polynomial regression	Test MSE = 16.7
Train a degree-4 polynomial regression	Test MSE = 12.2
Train a degree-5 polynomial regression	Test MSE = 14.8
Train a degree-6 polynomial regression	Test MSE = 25.1
Train a degree-7 polynomial regression	Test MSE = 39.4
Train a degree-8 polynomial regression	Test MSE = 53.0

Hyper-Parameter Tuning

Training Set	Test Set
Train a degree-1 polynomial regression	Test MSE = 23.2
Train a degree-2 polynomial regression	Test MSE = 19.0
Train a degree-3 polynomial regression	Test MSE = 16.7
Train a degree-4 polynomial regression	Test MSE = 12.2
Train a degree-5 polynomial regression	Test MSE = 14.8
Train a degree-6 polynomial regression	Test MSE = 25.1
Train a degree-7 polynomial regression	Test MSE = 39.4
Train a degree-8 polynomial regression	Test MSE = 53.0

Hyper-Parameter Tuning



Cross-Validation (Naïve Approach) for Hyper-Parameter Tuning

Cross-Validation (Naïve Approach)

labels



features

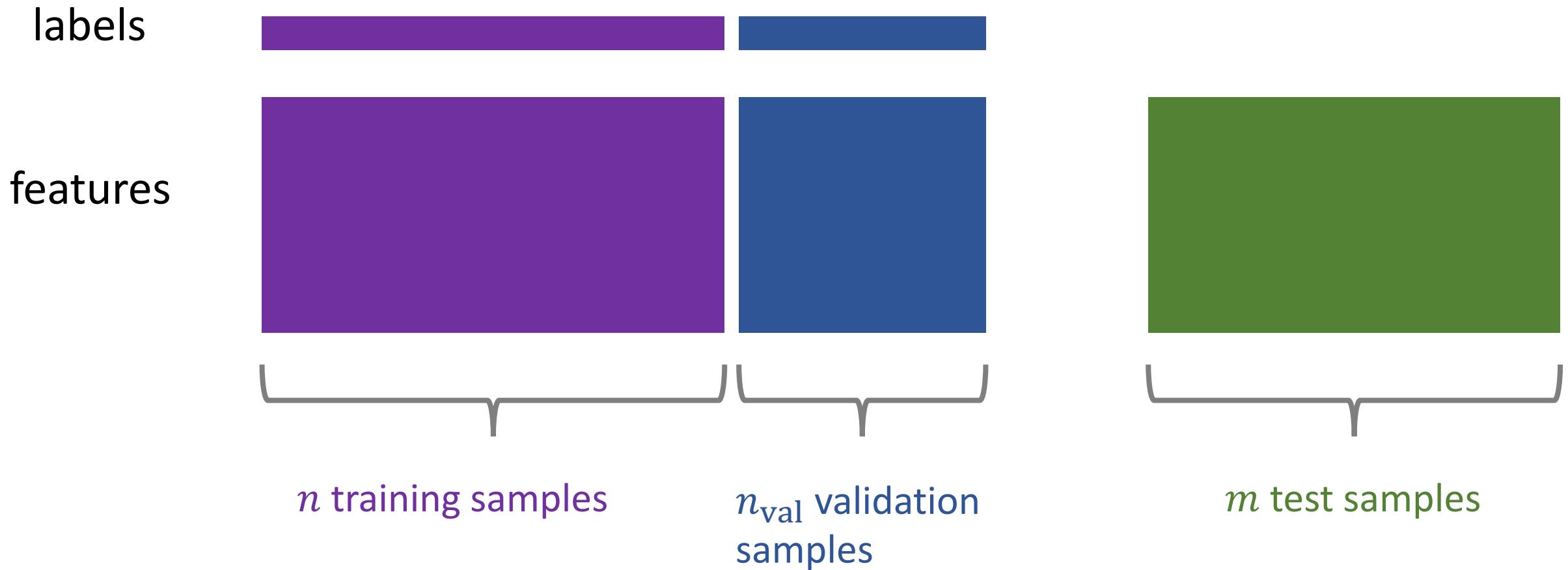


n training samples



m test samples

Cross-Validation (Naïve Approach)



Cross-Validation (Naïve Approach)

Training Set	Test MSE
Train a degree-1 polynomial regression	Test MSE = 23.2
Train a degree-2 polynomial regression	Test MSE = 19.0
Train a degree-3 polynomial regression	Test MSE = 16.7
Train a degree-4 polynomial regression	Test MSE = 12.2
Train a degree-5 polynomial regression	Test MSE = 14.8
Train a degree-6 polynomial regression	Test MSE = 25.1
Train a degree-7 polynomial regression	Test MSE = 39.4
Train a degree-8 polynomial regression	Test MSE = 53.0

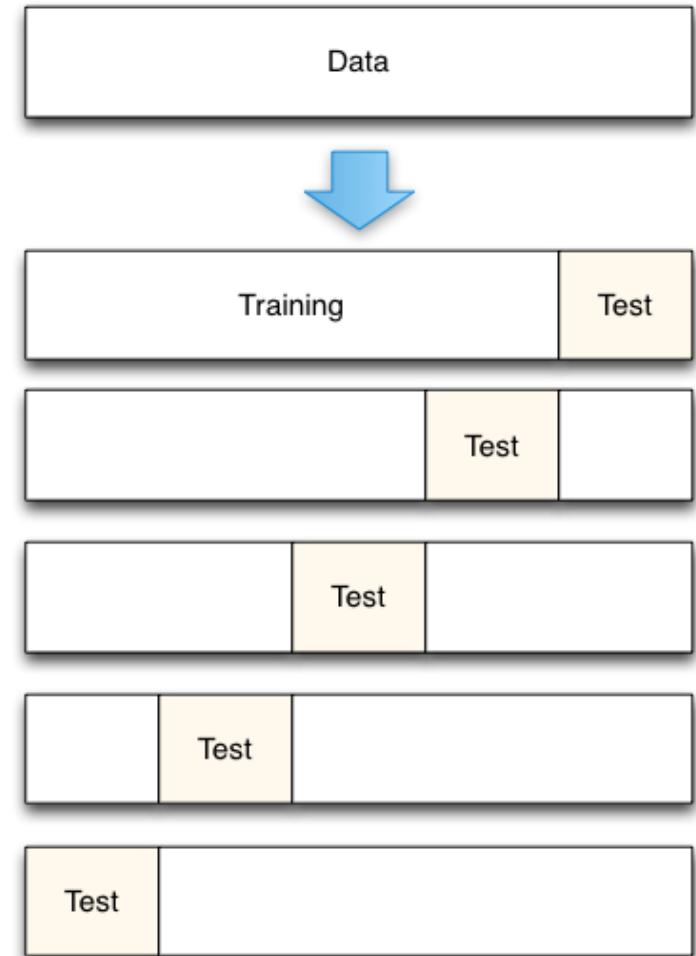
Cross-Validation (Naïve Approach)

Training Set	Validation Set Test Set
Train a degree-1 polynomial regression	Valid. MSE = 23.1
Train a degree-2 polynomial regression	Valid. MSE = 19.2
Train a degree-3 polynomial regression	Valid. MSE = 16.3
Train a degree-4 polynomial regression	Valid. MSE = 12.5
Train a degree-5 polynomial regression	Valid. MSE = 14.4
Train a degree-6 polynomial regression	Valid. MSE = 25.0
Train a degree-7 polynomial regression	Valid. MSE = 39.1
Train a degree-8 polynomial regression	Valid. MSE = 53.5

k-Fold Cross-Validation

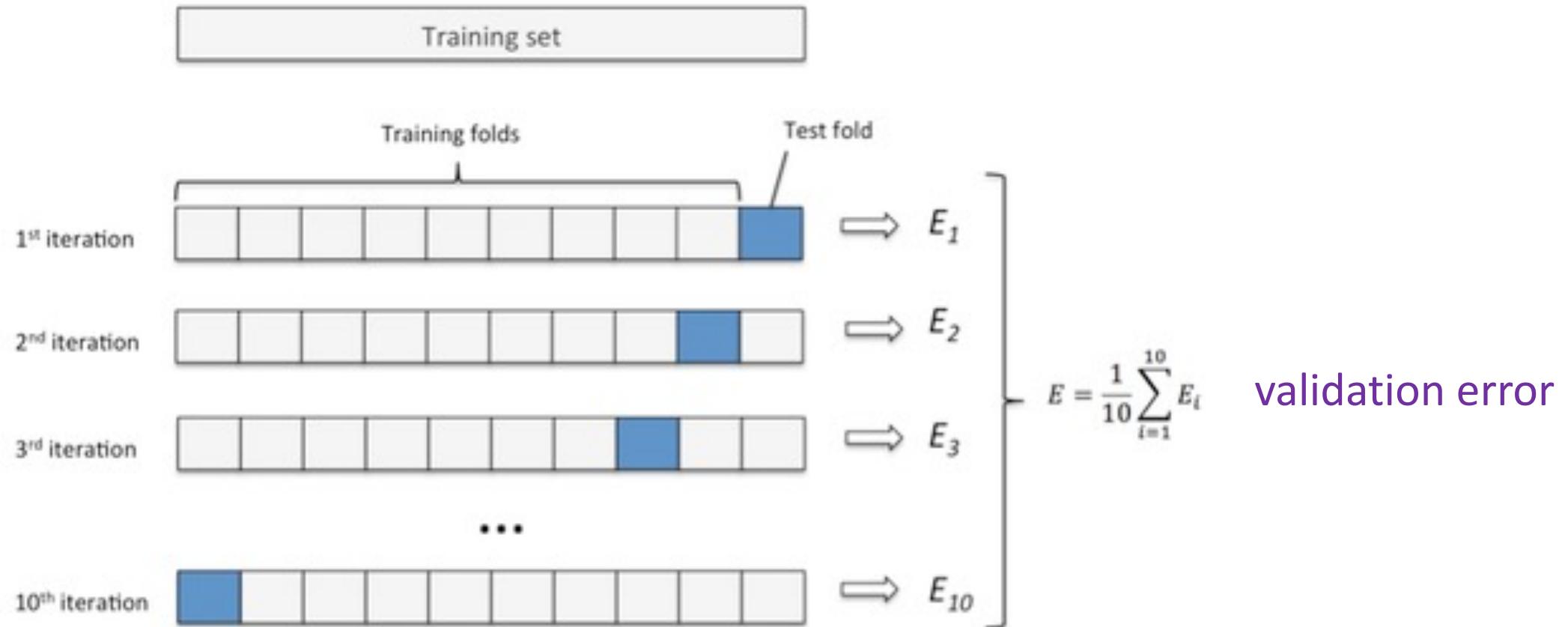
k -Fold Cross-Validation

1. Propose a grid of hyper-parameters.
 - E.g. $p \in \{1, 2, 3, 4, 5, 6\}$.
2. Randomly partition the training samples to k parts.
 - $k - 1$ parts for training.
 - One part for test.
3. Compute the averaged test errors of the k repeats.
 - The average is called the validation error.
4. Choose the hyper-parameter p that leads to the smallest validation error.



Example: 5-fold cross-validation

Example: 10-Fold Cross-Validation



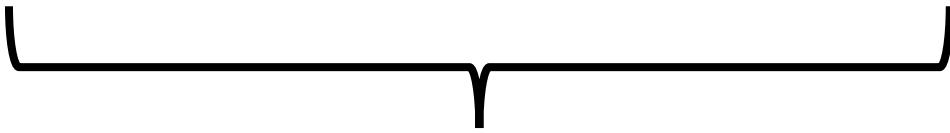
Example: 10-Fold Cross-Validation

hyper-parameter	validation error
p=1	23.19
p=2	21.00
p=3	18.54
p=4	24.36
p=5	27.96
p=6	33.10

Real-World Machine Learning Competition

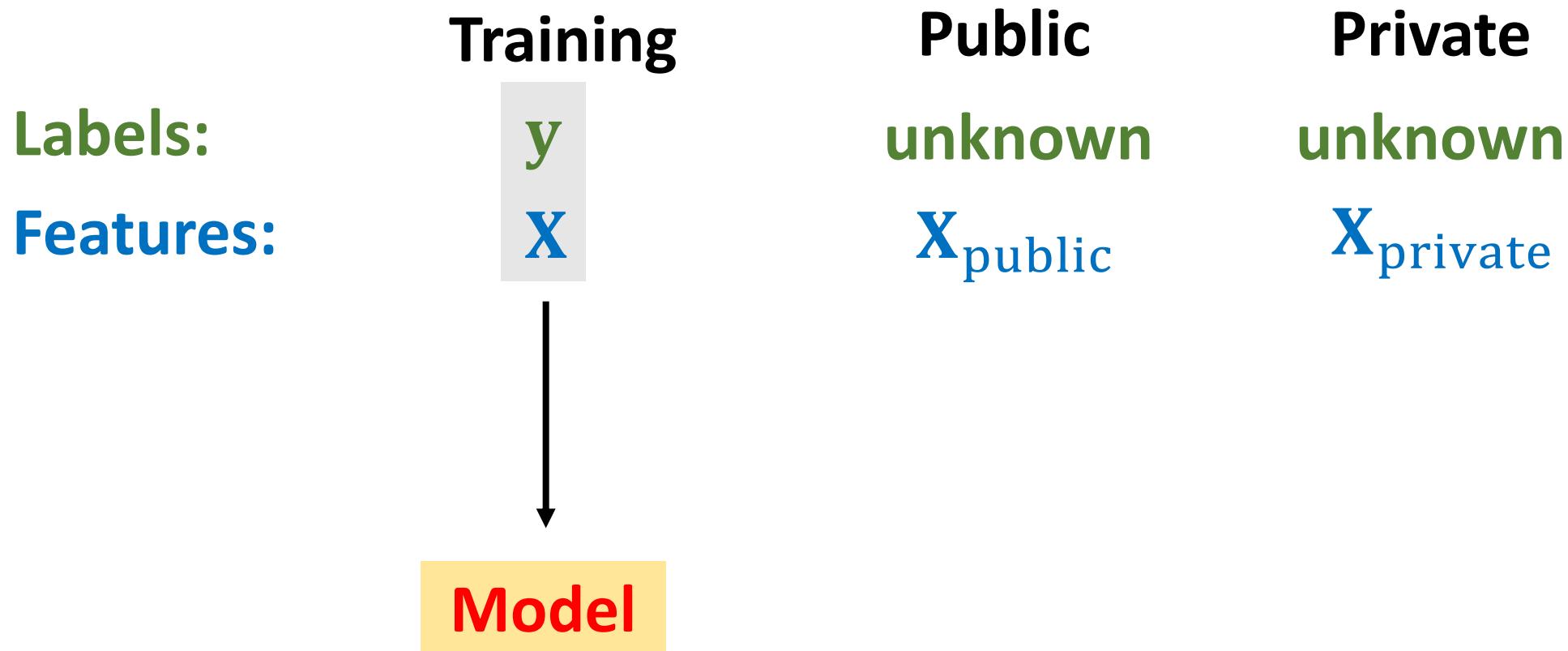
The Available Data

	Training	Public	Private
Labels:	y	unknown	unknown
Features:	X	X _{public}	X _{private}

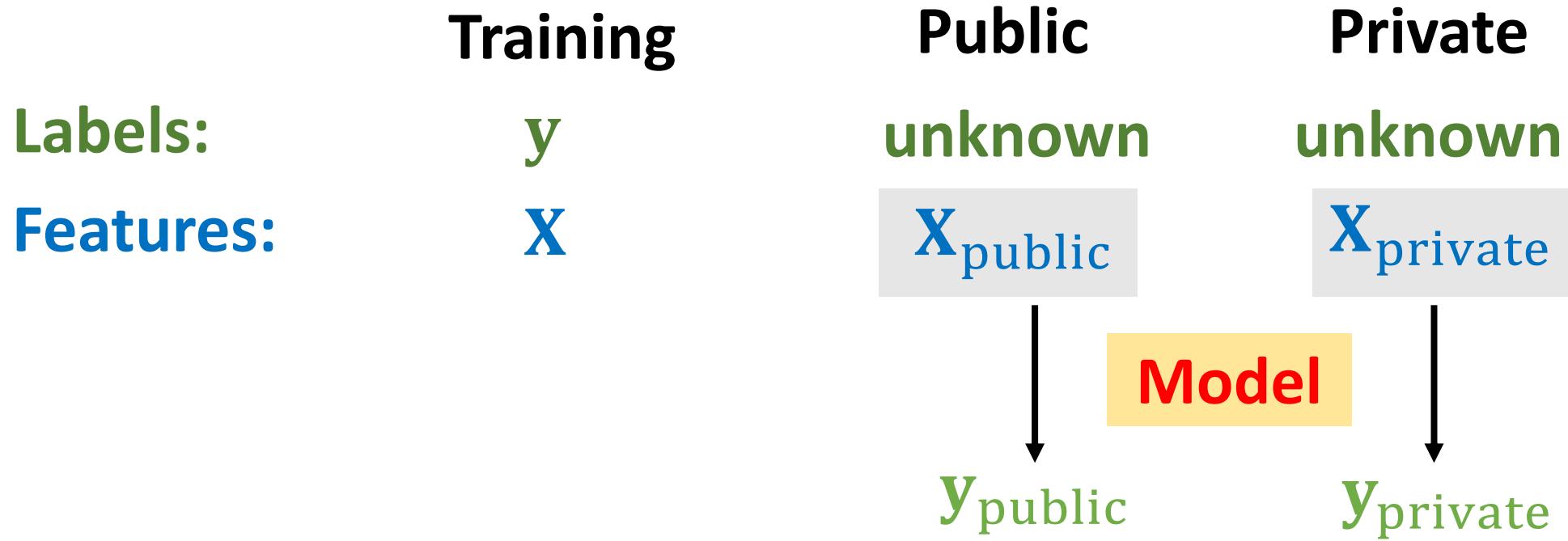

Test Data

The public and private are mixed;
Participants cannot distinguish them.

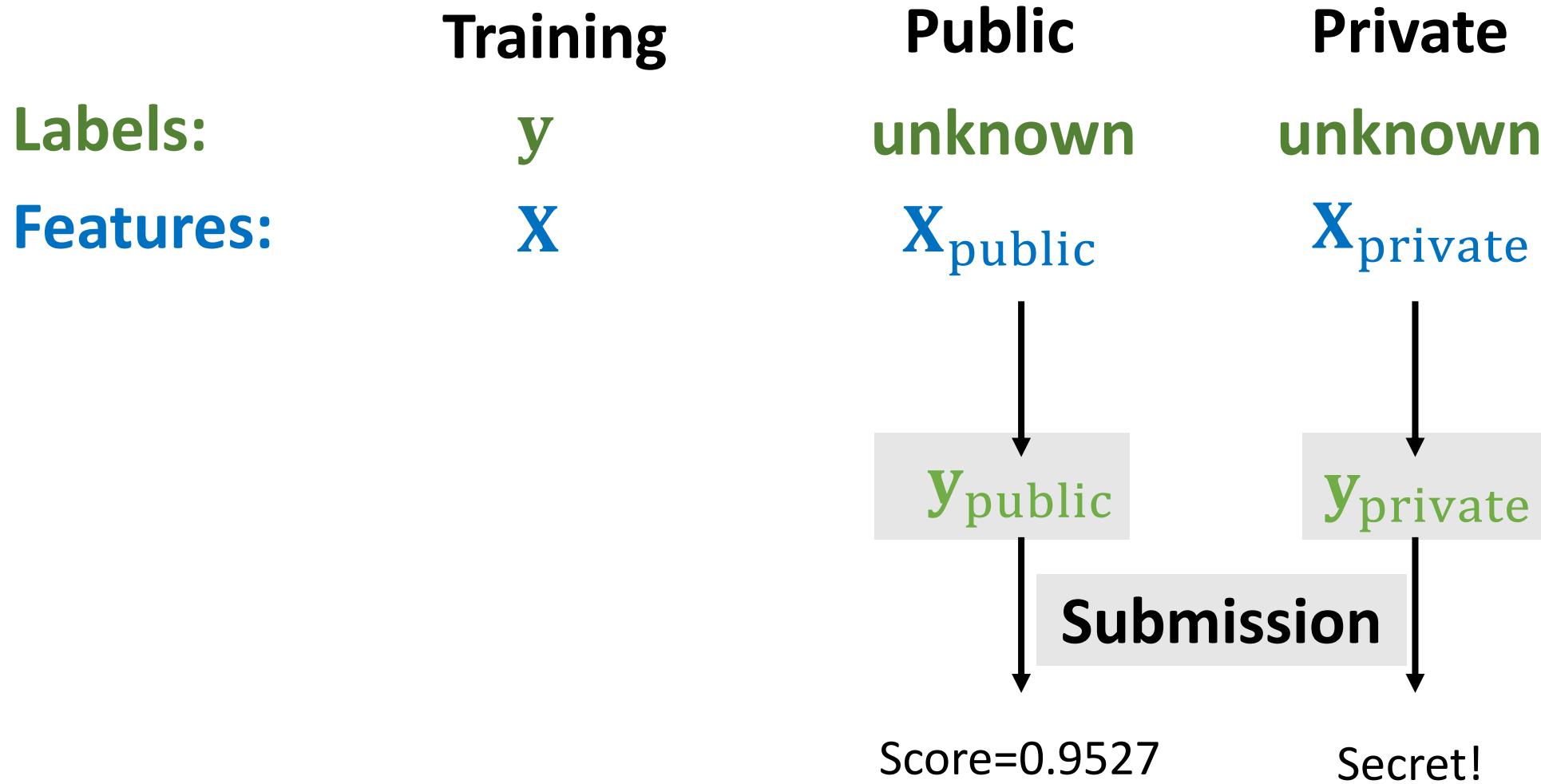
Train A Model



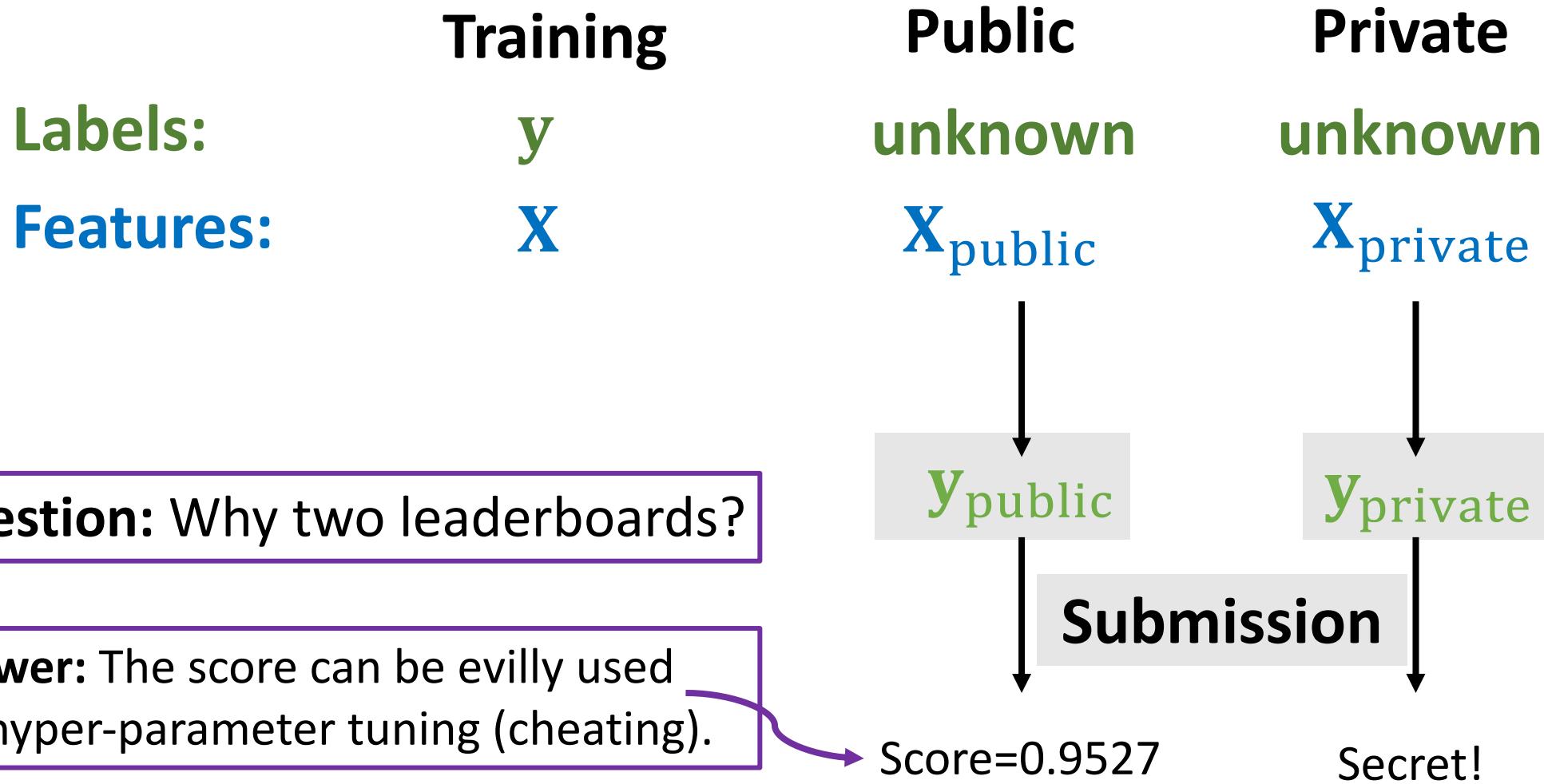
Prediction



Submission to Leaderboard



Submission to Leaderboard



Summary

- Polynomial regression for non-linear problems.
- Polynomial regression has a **hyper-parameter p** .
- Underfitting (very small p) and overfitting (very big p) .
- Tune the **hyper-parameters** using cross-validation.
- Make your model parameters and **hyper-parameters** independent of the test set!!!

Binary Classification

(Logistic Regression)

Vector and Matrix Derivatives

Derivative of Scalar w.r.t. Scalar

Examples:

- $y = x^2; \frac{dy}{dx} = 2x.$

- $y = e^x; \frac{dy}{dx} = e^x.$

Derivative of Vector w.r.t. Scalar

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a scalar $x \in \mathbb{R}$:

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{bmatrix}$$

- Example:

$$\mathbf{y} = \begin{bmatrix} 3x^2 \\ x + 1 \\ \log x \\ e^x \end{bmatrix}, \quad \frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} 6x \\ 1 \\ 1/x \\ e^x \end{bmatrix}$$

Derivative of Scalar w.r.t. Vector

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{bmatrix}$$

- Example 1:

$$y = \|\mathbf{x}\|_2^2 = \sum_{i=1}^m x_i^2, \quad \frac{\partial y}{\partial \mathbf{x}} = 2\mathbf{x}.$$

Derivative of Scalar w.r.t. Vector

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{bmatrix}$$

- Example 2:

$$y = \mathbf{x}^T \mathbf{z} = \sum_{i=1}^m x_i z_i, \quad \frac{\partial y}{\partial \mathbf{x}} = \mathbf{z}.$$

Derivative of Scalar w.r.t. Vector

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{bmatrix}$$

- Example 3:

$$y = \sum_{i=1}^m \log(1 + e^{-x_i}),$$

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \log(1+e^{-x_1})}{\partial x_1} \\ \vdots \\ \frac{\partial \log(1+e^{-x_m})}{\partial x_m} \end{bmatrix} = \begin{bmatrix} -\frac{1}{1+e^{x_1}} \\ \vdots \\ -\frac{1}{1+e^{x_m}} \end{bmatrix}$$

Derivative of Vector w.r.t. Vector

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

m × n matrix

- Example 1:

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \underbrace{\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}}_{m \times m}$$

The (i, j) -th entry is $\frac{\partial y_j}{\partial x_i}$

Derivative of Vector w.r.t. Vector

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \quad m \times n \text{ matrix}$$

- Example 2:

$$\mathbf{y} = \begin{bmatrix} a_1 x_1^2 \\ a_2 x_2^2 \\ \vdots \\ a_m x_m^2 \end{bmatrix} \in \mathbb{R}^m, \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \underbrace{\begin{bmatrix} 2a_1 x_1 & 0 & \dots & 0 \\ 0 & 2a_2 x_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 2a_m x_m \end{bmatrix}}_{m \times m}$$

Derivative of Vector w.r.t. Vector

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

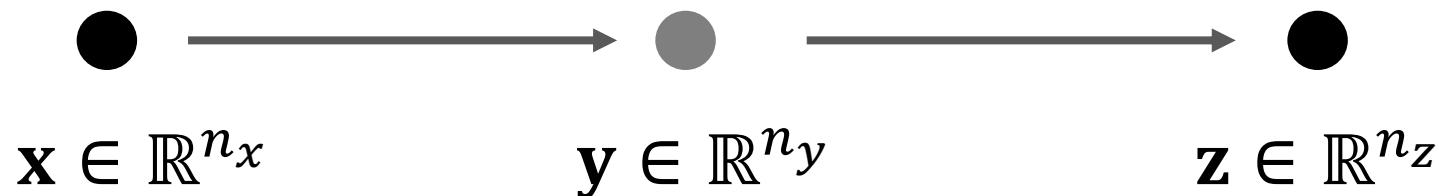
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \quad m \times n \text{ matrix}$$

- Example 3:

$$\mathbf{A} \in \mathbb{R}^{n \times m}, \quad \mathbf{y} = \mathbf{Ax} \in \mathbb{R}^n, \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}^T \in \mathbb{R}^{m \times n}$$

Chain Rule

- Let $\mathbf{z} \in \mathbb{R}^{n_z}$ be a function of $\mathbf{y} \in \mathbb{R}^{n_y}$ and \mathbf{y} be a function of $\mathbf{x} \in \mathbb{R}^{n_x}$.



$$\underbrace{\frac{d\mathbf{z}}{d\mathbf{x}}}_{n_x \times n_z} = \underbrace{\frac{d\mathbf{y}}{d\mathbf{x}}}_{n_x \times n_y} \underbrace{\frac{d\mathbf{z}}{d\mathbf{y}}}_{n_y \times n_z}$$

Derivative of Scalar w.r.t. Matrix

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a matrix $\mathbf{Z} \in \mathbb{R}^{p \times q}$:
 1. Vectorization: $\mathbf{x} = \text{vec}(\mathbf{Z}) \in \mathbb{R}^{pq \times 1}$.
 2. Compute $\frac{\partial y}{\partial \mathbf{x}} \in \mathbb{R}^{pq \times 1}$.
 3. Reshape the resulting $pq \times 1$ vector to $p \times q$ matrix.

Derivative of Vector w.r.t. Matrix

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a matrix $\mathbf{Z} \in \mathbb{R}^{p \times q}$:
 1. Vectorization: $\mathbf{x} = \text{vec}(\mathbf{Z}) \in \mathbb{R}^{pq \times 1}$.
 2. Compute $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{pq \times n}$.
 3. Reshape the resulting $pq \times n$ matrix to $p \times q \times n$ tensor.

Binary Classification

Tasks

Methods

Algorithms

Binary Classification

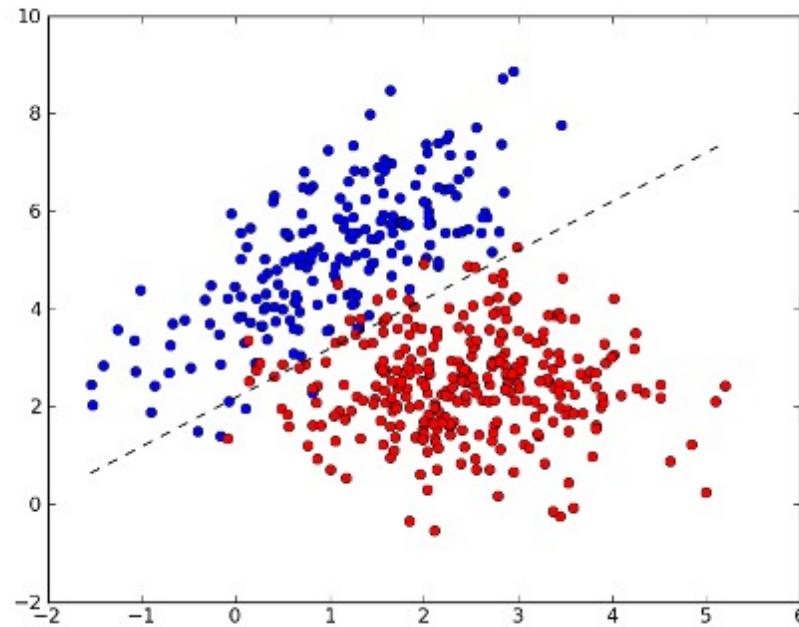
Input: feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \{-1, +1\}$.

Output: a function $f: \mathbb{R}^d \mapsto \{-1, +1\}$.

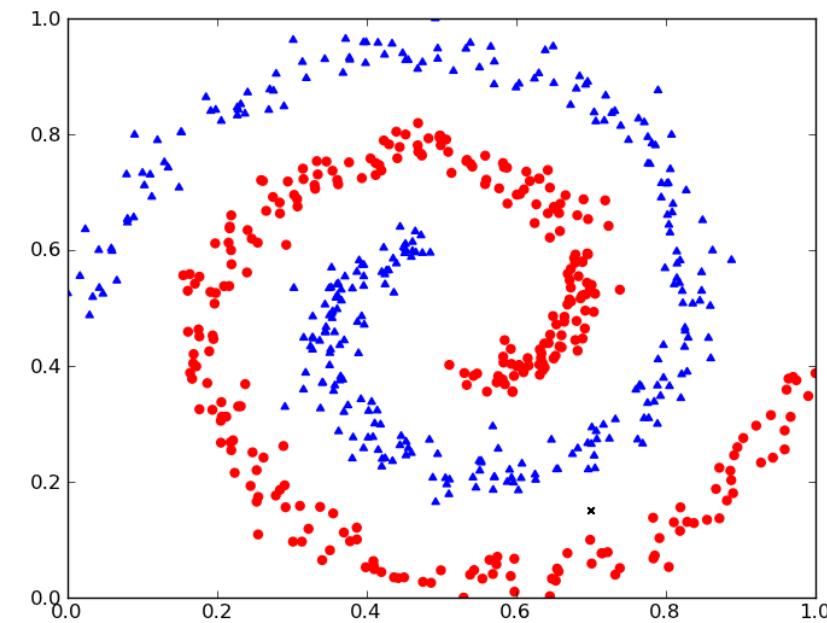
Binary Classification

Input: feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \dots, y_n \in \{-1, +1\}$.

Output: a function $f: \mathbb{R}^d \mapsto \{-1, +1\}$.



Linear Classification



Nonlinear Classification

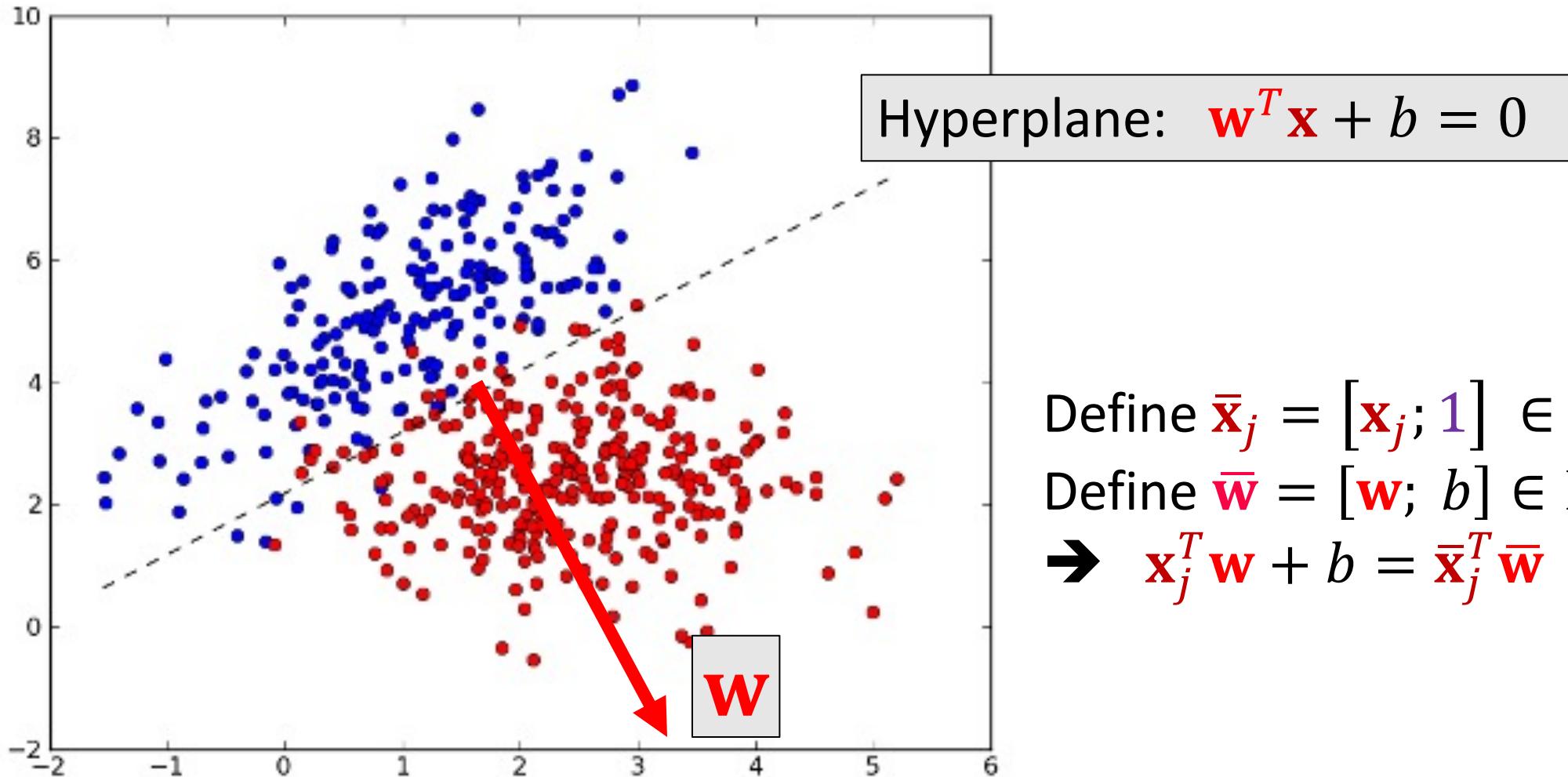
Logistic Regression (Linear Classifier)

Tasks

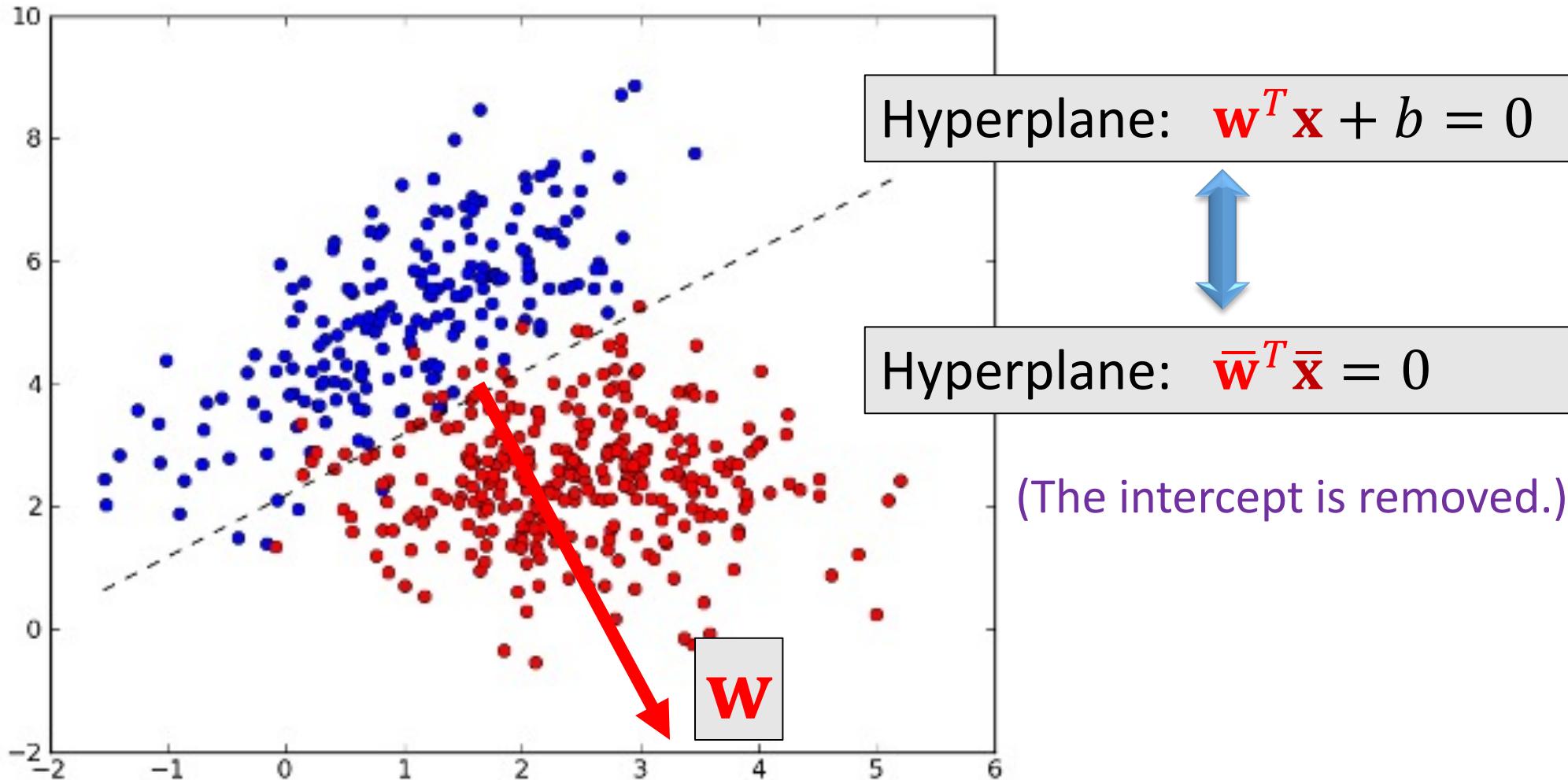
Methods

Algorithms

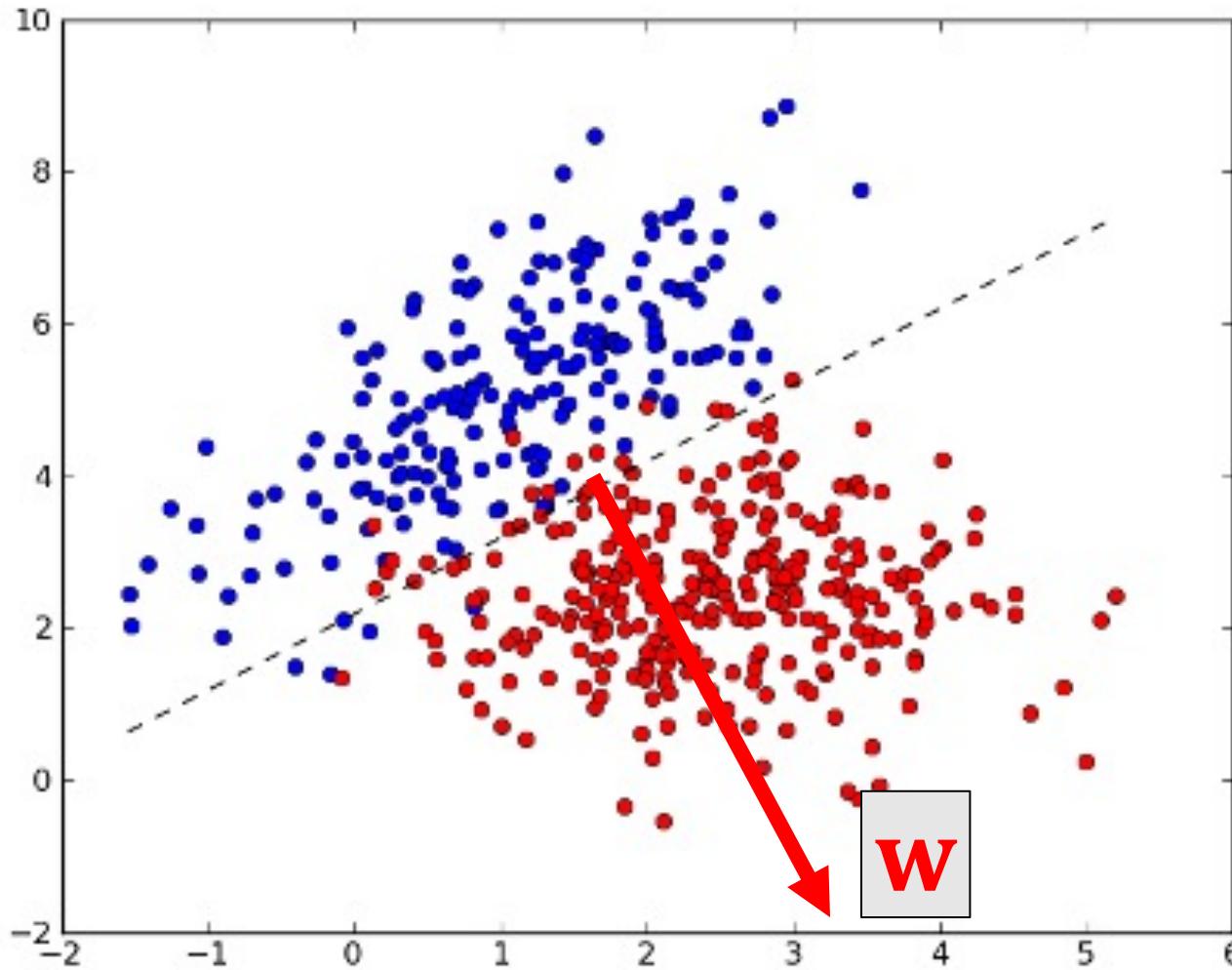
Linear Classifier



Linear Classifier

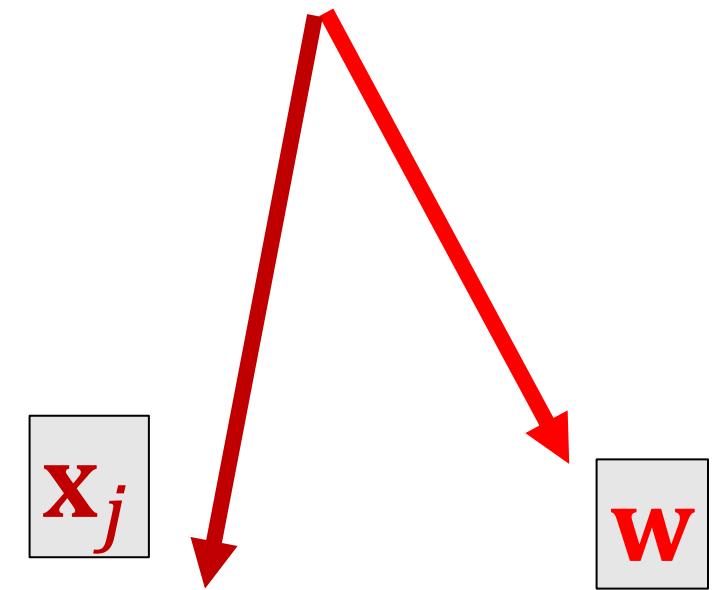


Linear Classifier

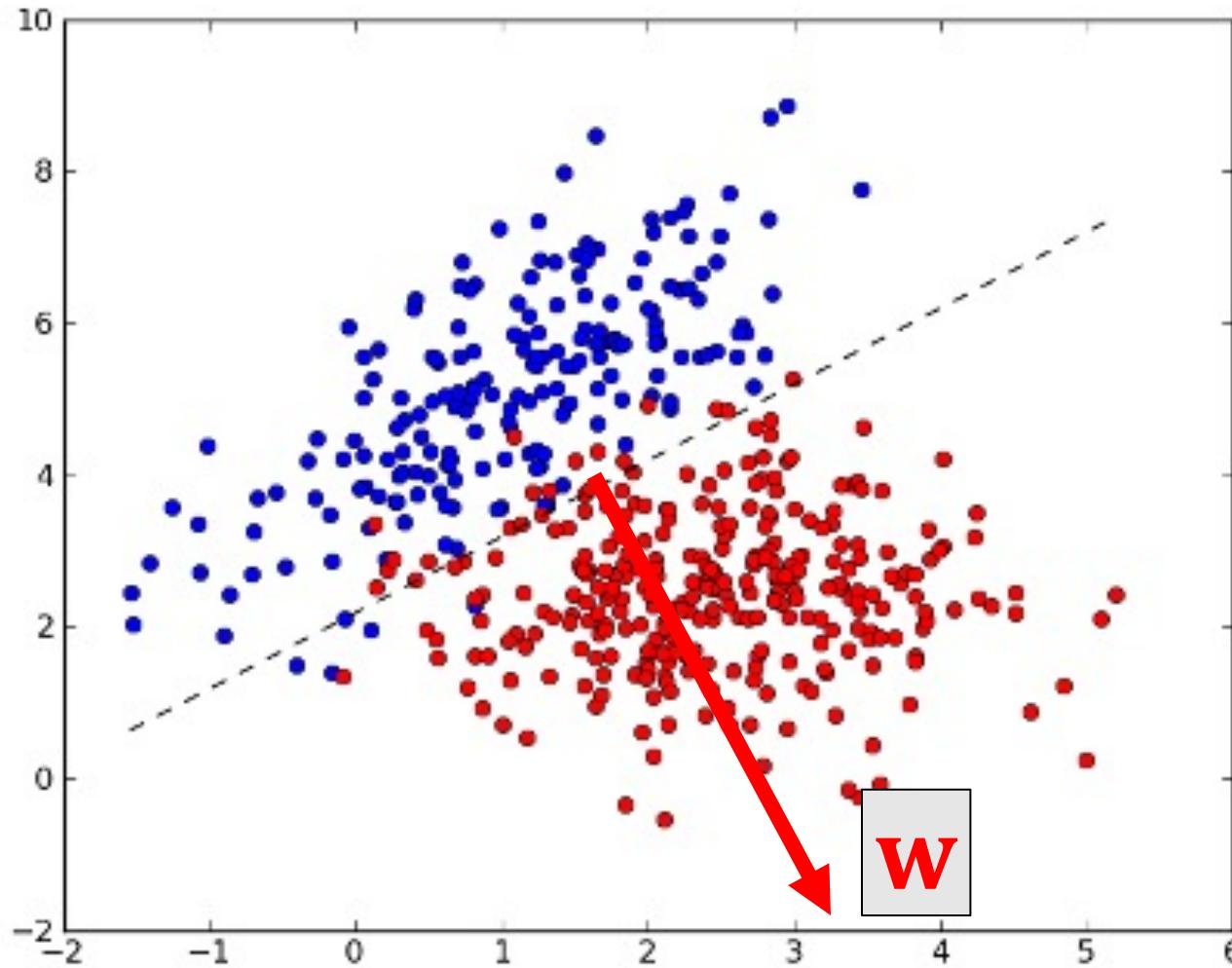


Learn a vector \mathbf{w} such that

- If $y_j = +1$, then $\mathbf{w}^T \mathbf{x}_j > 0$.



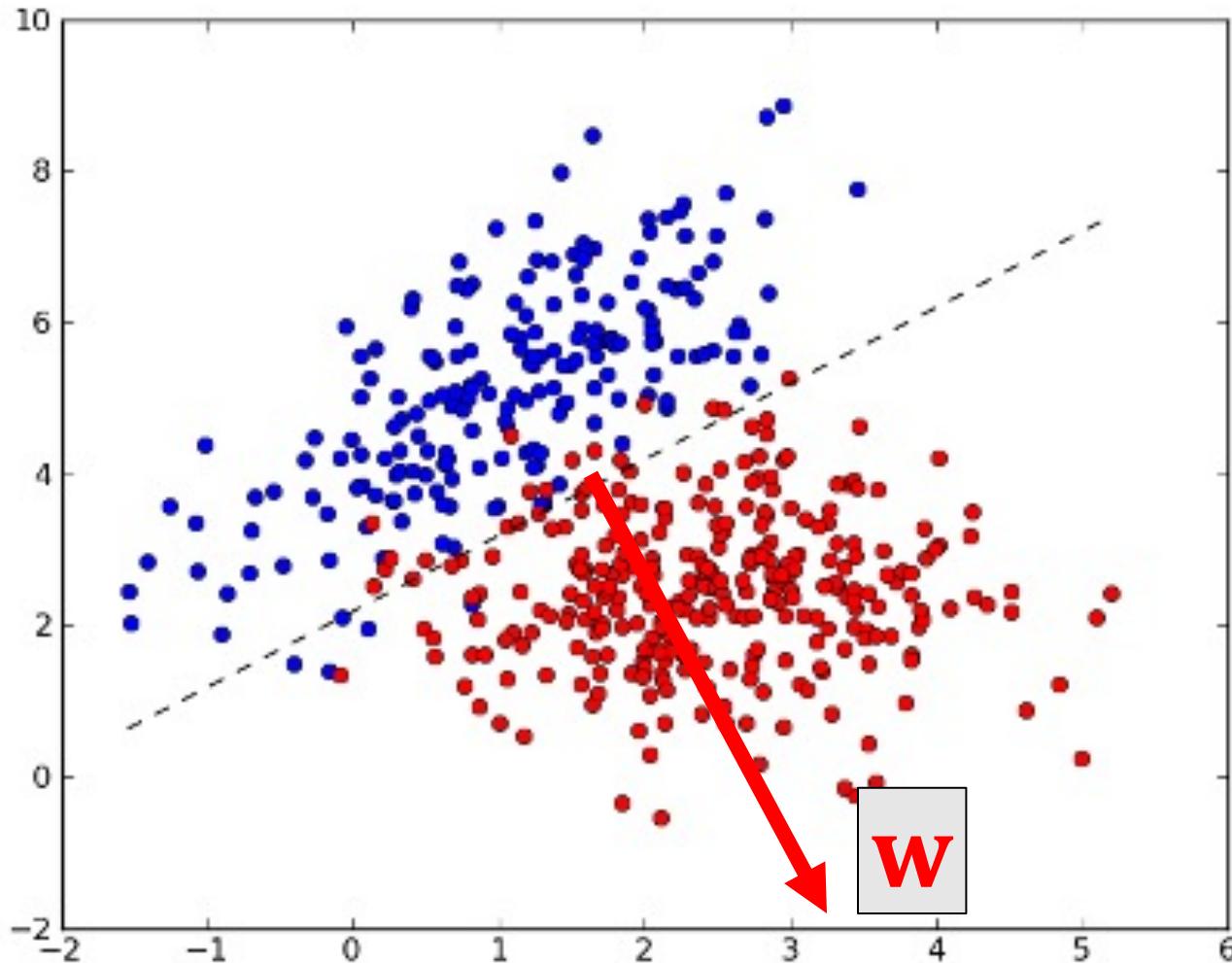
Linear Classifier



Learn a vector \mathbf{w} such that

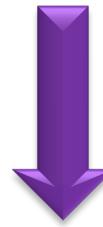
- If $y_j = +1$, then $\mathbf{w}^T \mathbf{x}_j > 0$.
- If $y_j = -1$, then $\mathbf{w}^T \mathbf{x}_j < 0$.

Linear Classifier



Learn a vector \mathbf{w} such that

- If $y_j = +1$, then $\mathbf{w}^T \mathbf{x}_j > 0$.
- If $y_j = -1$, then $\mathbf{w}^T \mathbf{x}_j < 0$.

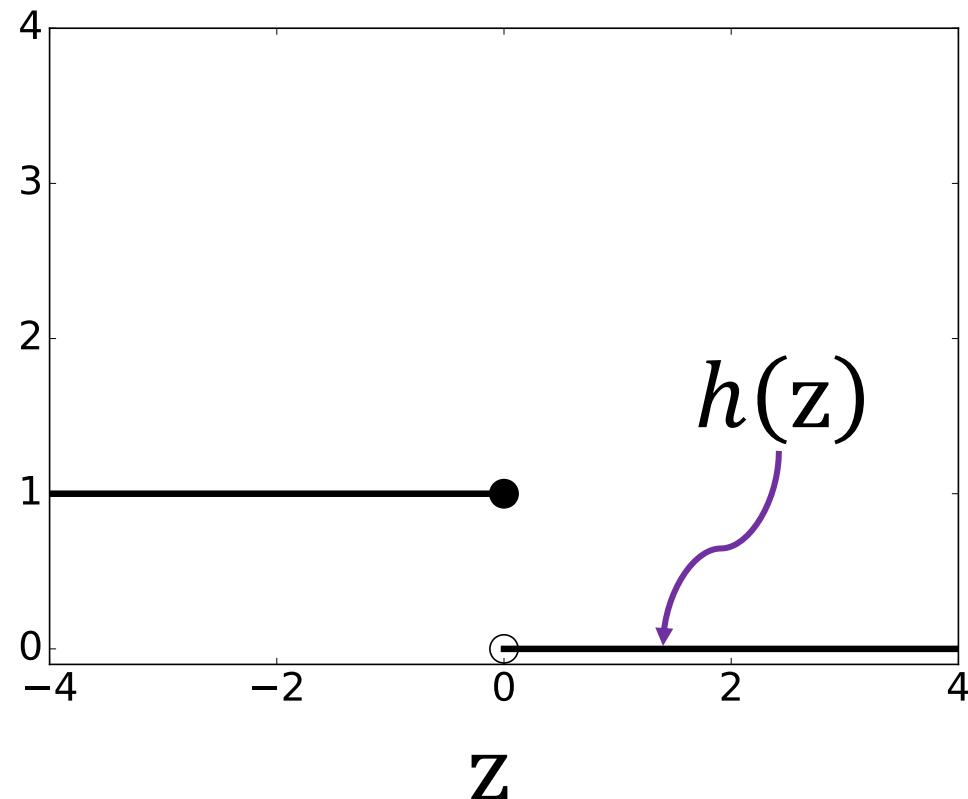


Key Idea:

Encourage $y_j \mathbf{w}^T \mathbf{x}_j$ to be positive

Directly Minimize the Classification Error?

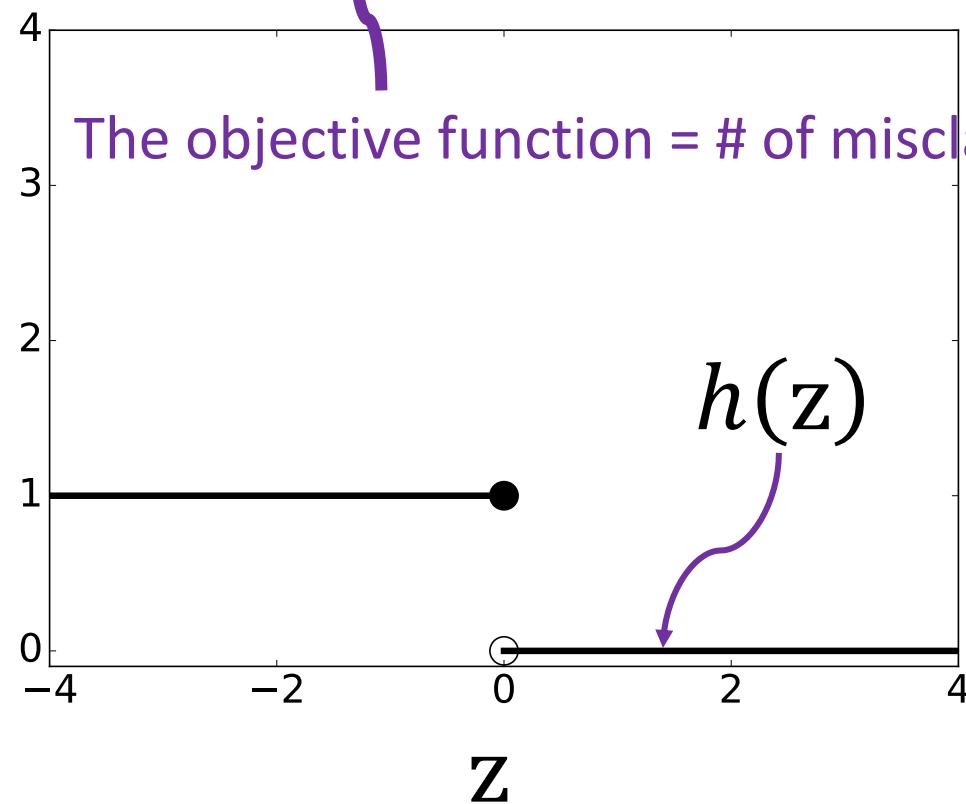
Minimize $\sum_j h(y_j \mathbf{w}^T \mathbf{x}_j)$, where $h(z) = \begin{cases} 1, & \text{if } z < 0; \\ 0, & \text{if } z \geq 0. \end{cases}$



Key Idea:
Encourage $y_j \mathbf{w}^T \mathbf{x}_j$ to be positive

Directly Minimize the Classification Error?

Minimize $\sum_j h(y_j \mathbf{w}^T \mathbf{x}_j)$, where $h(z) = \begin{cases} 1, & \text{if } z < 0; \\ 0, & \text{if } z \geq 0. \end{cases}$



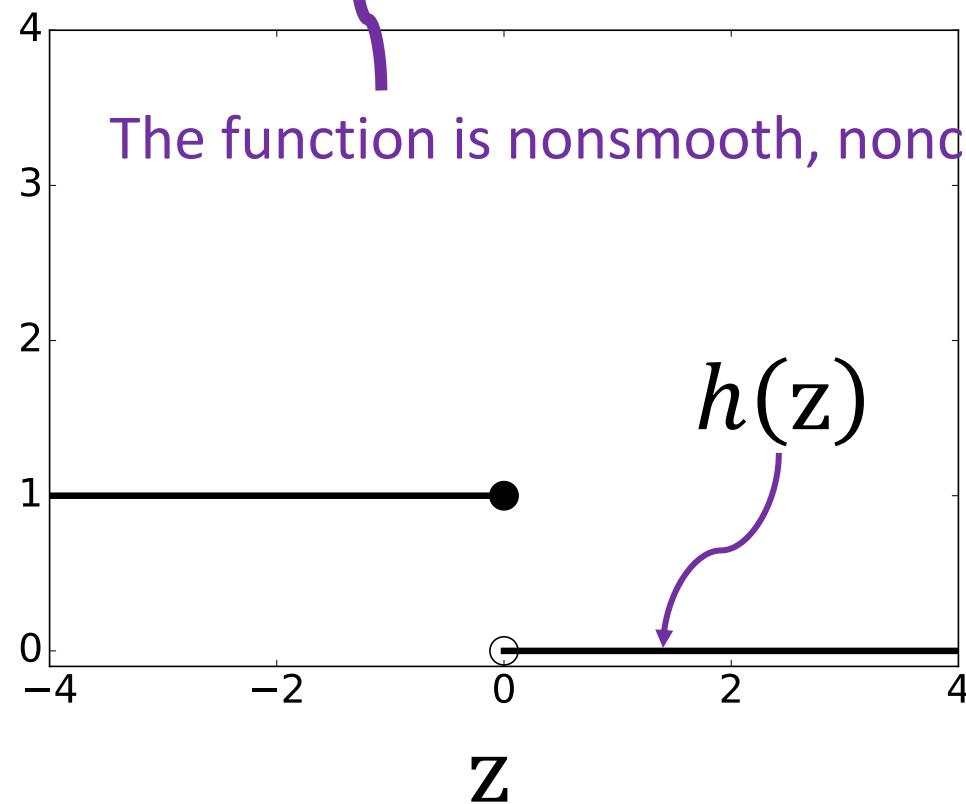
$h(z)$

The objective function = # of misclassified training samples

Key Idea:
Encourage $y_j \mathbf{w}^T \mathbf{x}_j$ to be positive

Directly Minimize the Classification Error?

Minimize $\sum_j h(y_j \mathbf{w}^T \mathbf{x}_j)$, where $h(z) = \begin{cases} 1, & \text{if } z < 0; \\ 0, & \text{if } z \geq 0. \end{cases}$

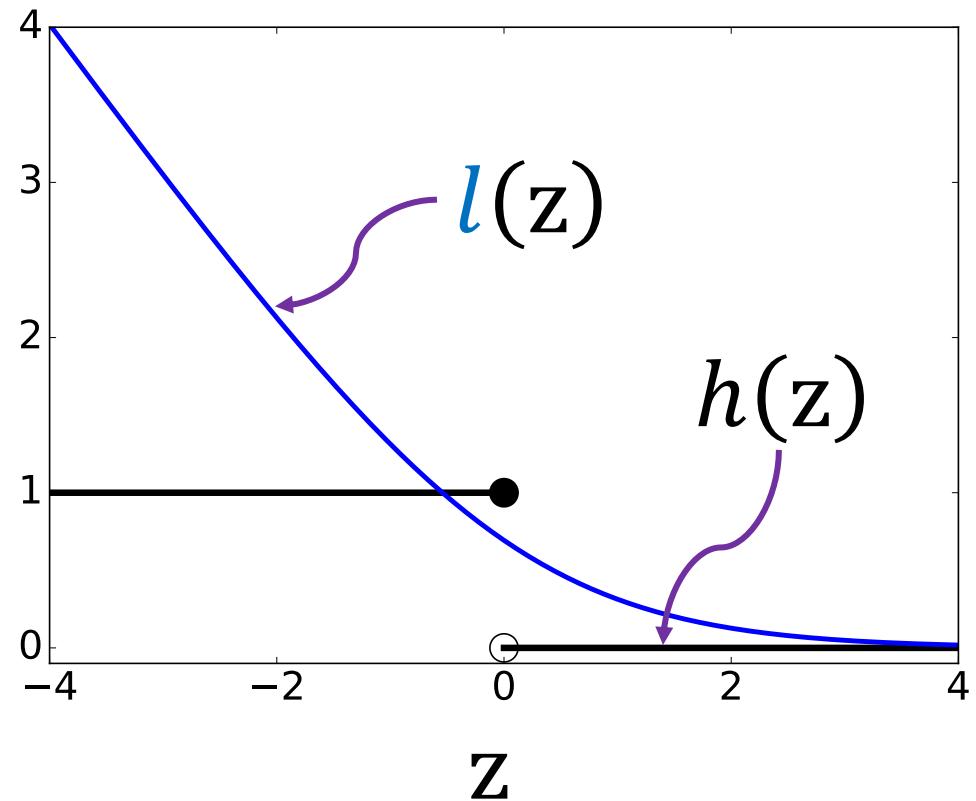


The function is nonsmooth, nonconvex, and hard to optimize.

Key Idea:
Encourage $y_j \mathbf{w}^T \mathbf{x}_j$ to be positive

Logistic Regression

Minimize $\sum_j l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.



Key Idea:
Encourage $y_j \mathbf{w}^T \mathbf{x}_j$ to be positive

Logistic Regression

Tasks

Methods

Algorithms

Logistic Regression

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

Tasks

Binary Classification

Multi-Class Classification

Methods

Logistic Regression

SVM

Neural Networks

Algorithms

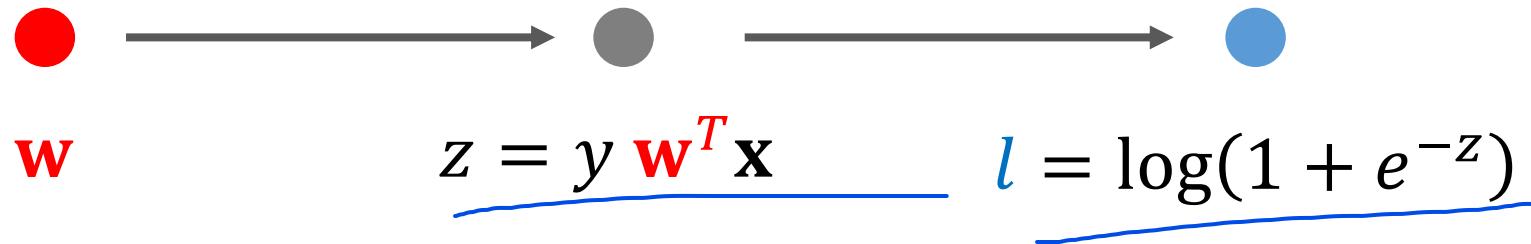
Gradient Descent (GD)

Accelerated GD

Stochastic GD

Gradient

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.



Gradient

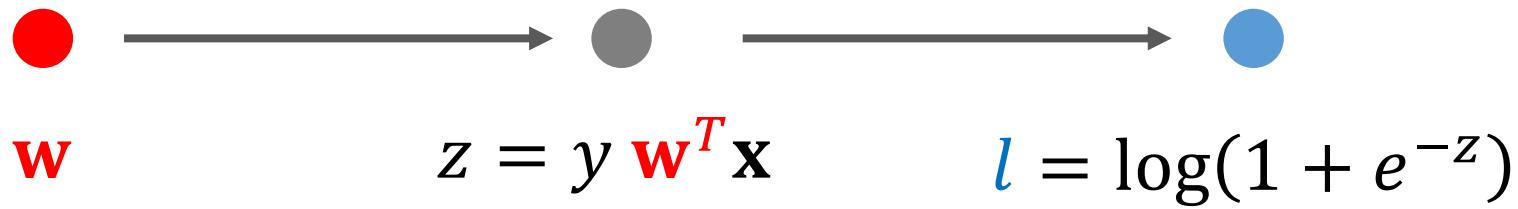
Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

The diagram illustrates the logistic regression process. It starts with a red circle labeled w , followed by a horizontal arrow pointing to a gray circle. This is followed by another horizontal arrow pointing to a blue circle. Below the first arrow is the label w . Below the second arrow is the equation $z = y \mathbf{w}^T \mathbf{x}$. Below the third arrow is the equation $l = \log(1 + e^{-z})$.

A blue convex curve, representing a function like $l(z)$, is shown. A tangent line is drawn at a point on the curve. A dot on the curve is connected by a line to the point of tangency, representing the gradient vector. The equation $\bullet \frac{\partial z}{\partial w} = y\mathbf{x}$ is written inside a blue box.

Gradient

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.


$$\begin{array}{ccc} \text{Red circle} & \xrightarrow{\hspace{2cm}} & \text{Grey circle} \\ \mathbf{w} & & z = y \mathbf{w}^T \mathbf{x} \\ & & \mathcal{l} = \log(1 + e^{-z}) \end{array}$$

- $\frac{\partial z}{\partial \mathbf{w}} = y \mathbf{x}, \quad \frac{\partial \mathcal{l}(z)}{\partial z} = \frac{-e^{-z}}{1+e^{-z}} = \boxed{\frac{1}{1+e^z}}$

Gradient

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.



- $\frac{\partial z}{\partial \mathbf{w}} = y\mathbf{x}$, $\frac{\partial l(z)}{\partial z} = \frac{-e^{-z}}{1+e^{-z}} = -\frac{1}{1+e^z}$.
- Chain rule: $\frac{\partial l}{\partial \mathbf{w}} = \frac{\partial z}{\partial \mathbf{w}} \cdot \frac{\partial l}{\partial z} = (y\mathbf{x}) \left(-\frac{1}{1+e^z} \right) = -\frac{y\mathbf{x}}{1+\exp(y\mathbf{w}^T \mathbf{x})}$.

Gradient

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

- We have shown: $\frac{\partial l(y \mathbf{w}^T \mathbf{x})}{\partial \mathbf{w}} = \frac{-y \mathbf{x}}{1 + \exp(y \mathbf{w}^T \mathbf{x})}$.

- Objective function: $f(\mathbf{w}) = \frac{1}{n} \sum_j l(y_j \mathbf{w}^T \mathbf{x}_j)$.

Gradient

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

- We have shown: $\frac{\partial l(y \mathbf{w}^T \mathbf{x})}{\partial \mathbf{w}} = \frac{-y \mathbf{x}}{1 + \exp(y \mathbf{w}^T \mathbf{x})}$.
- Objective function: $f(\mathbf{w}) = \frac{1}{n} \sum_j l(y_j \mathbf{w}^T \mathbf{x}_j)$.
- $\frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_j \frac{\partial l(y_j \mathbf{w}^T \mathbf{x}_j)}{\partial \mathbf{w}} = \frac{1}{n} \sum_j \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}^T \mathbf{x}_j)}$.

Gradient

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

- We have shown: $\frac{\partial l(y \mathbf{w}^T \mathbf{x})}{\partial \mathbf{w}} = \frac{-y \mathbf{x}}{1 + \exp(y \mathbf{w}^T \mathbf{x})}$.
- Objective function: $f(\mathbf{w}) = \frac{1}{n} \sum_j l(y_j \mathbf{w}^T \mathbf{x}_j)$.
- $\frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_j \frac{\partial l(y_j \mathbf{w}^T \mathbf{x}_j)}{\partial \mathbf{w}} = \frac{1}{n} \sum_j \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}^T \mathbf{x}_j)}$.

Gradient Descent (GD) Algorithm

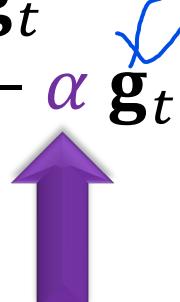
Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

t is epoch or iteration

GD repeat:

1. Compute gradient: \mathbf{g}_t
2. Update: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \mathbf{g}_t$



Tune the step size (learning rate) α

Algorithms

Gradient Descent (GD)

Accelerated GD

Stochastic GD

AGD Algorithm

Logistic regression: $\min_{\mathbf{w}} \frac{1}{n} \sum_{j=1}^n l(y_j \mathbf{w}^T \mathbf{x}_j)$, where $l(z) = \log(1 + e^{-z})$.

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

AGD repeat:

1. Compute gradient: \mathbf{g}_t
2. Update momentum: $\mathbf{v}_{t+1} = \beta \mathbf{v}_t + \mathbf{g}_t$
3. Update: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \mathbf{v}_{t+1}$

allows choosing how steep the GD should go

Algorithms

Gradient Descent (GD)

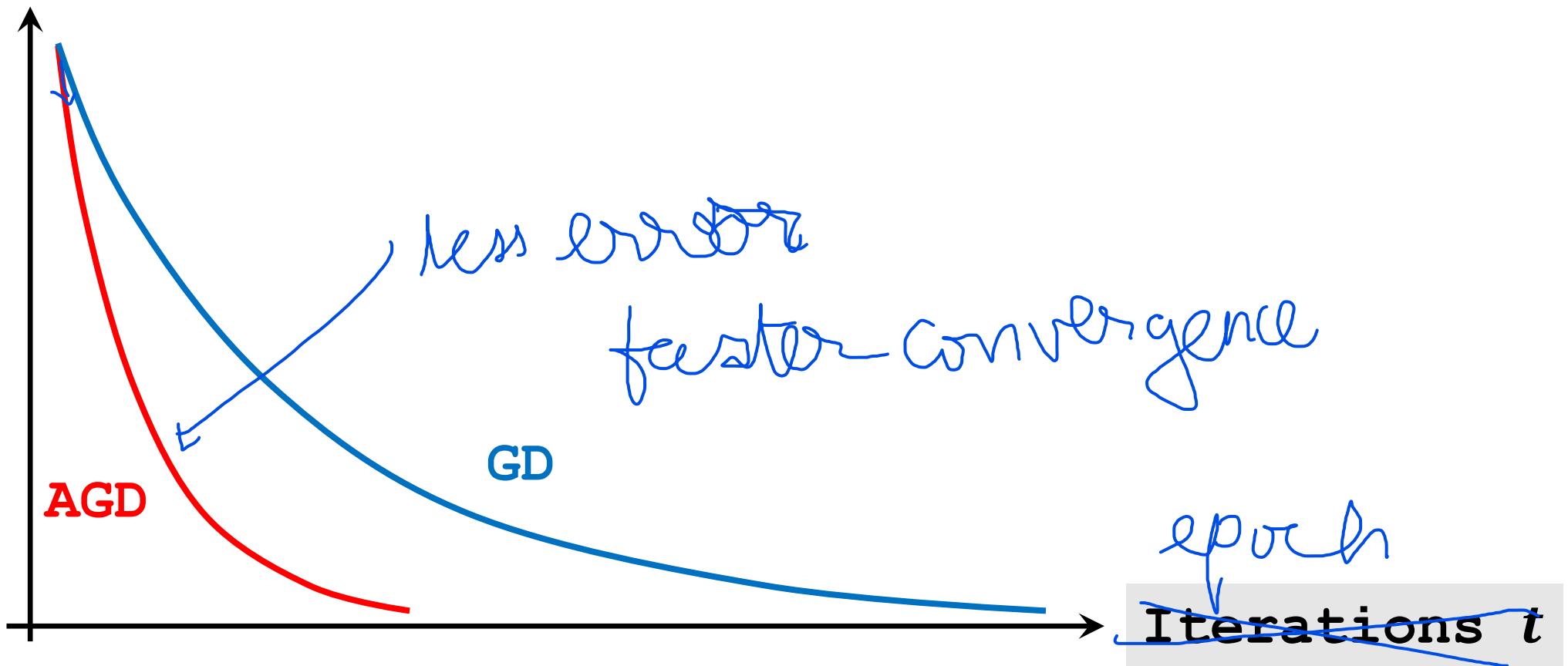
Accelerated GD

Tune α and β ($0 \leq \beta < 1$)

Stochastic GD

GD versus AGD

$$\| \mathbf{w}_t - \mathbf{w}^* \|_2 = l_2 \text{ norm error}$$



Time Complexity

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \tilde{\mathbf{g}}_{t,j}$, where $\tilde{\mathbf{g}}_{t,j} = \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

Per-iteration time complexity is $O(nd)$.

- $O(d)$ time for computing $\mathbf{w}_t^T \mathbf{x}_j$.
- $O(d)$ time for computing $\tilde{\mathbf{g}}_{t,j}$.
- $O(nd)$ time for computing all the $\tilde{\mathbf{g}}_{t,j}$.

Algorithms

Gradient Descent (GD)

Accelerated GD

Stochastic GD

SGD Algorithm

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \tilde{\mathbf{g}}_{t,j}$, where $\tilde{\mathbf{g}}_{t,j} = \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

The stochastic gradient is close to the full gradient:

$$\mathbf{g}_t = \mathbb{E}_j [\tilde{\mathbf{g}}_{t,j}],$$

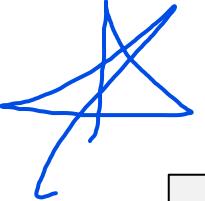
where j is randomly sampled from $\{1, \dots, n\}$.

Algorithms

Gradient Descent (GD)

Accelerated GD

Stochastic GD



SGD Algorithm

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \tilde{\mathbf{g}}_{t,j}$, where $\tilde{\mathbf{g}}_{t,j} = \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

SGD repeats

1. Randomly draw j from $\{1, 2, \dots, n\}$.
2. Compute the stochastic gradient $\tilde{\mathbf{g}}_{t,j}$.
3. Update: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \tilde{\mathbf{g}}_{t,j}$.

Algorithms

Per-iteration time complexity is $O(d)$.

Gradient Descent (GD)

Accelerated GD

Stochastic GD

Accelerated SGD Algorithm

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \tilde{\mathbf{g}}_{t,j}$, where $\tilde{\mathbf{g}}_{t,j} = \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

Accelerated SGD repeats

1. Randomly draw j from $\{1, 2, \dots, n\}$.
2. Compute the stochastic gradient $\tilde{\mathbf{g}}_{t,j}$.
3. Update momentum: $\mathbf{v}_{t+1} = \beta \mathbf{v}_t + \tilde{\mathbf{g}}_{t,j}$.
4. Update: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \mathbf{v}_{t+1}$.

Algorithms

Gradient Descent (GD)

Accelerated GD

Stochastic GD

SGD Algorithm

Gradient at \mathbf{w}_t : $\mathbf{g}_t = \frac{1}{n} \sum_{j=1}^n \tilde{\mathbf{g}}_{t,j}$, where $\tilde{\mathbf{g}}_{t,j} = \frac{-y_j \mathbf{x}_j}{1 + \exp(y_j \mathbf{w}_t^T \mathbf{x}_j)}$.

Output of SGD:

- Option 1: output the last iteration \mathbf{w}_{T+1}
- Option 2: output the average of \mathbf{w} produced by the last tens of iteration.

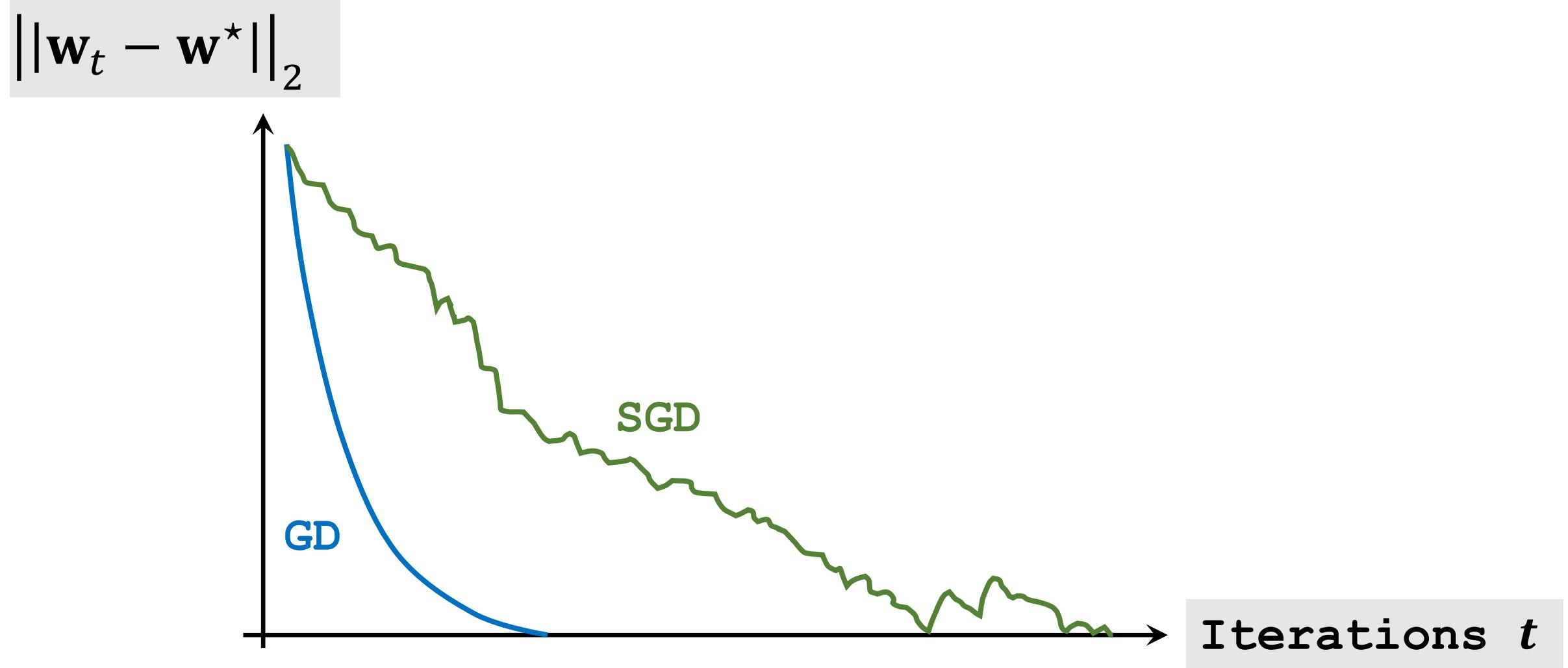
Algorithms

Gradient Descent (GD)

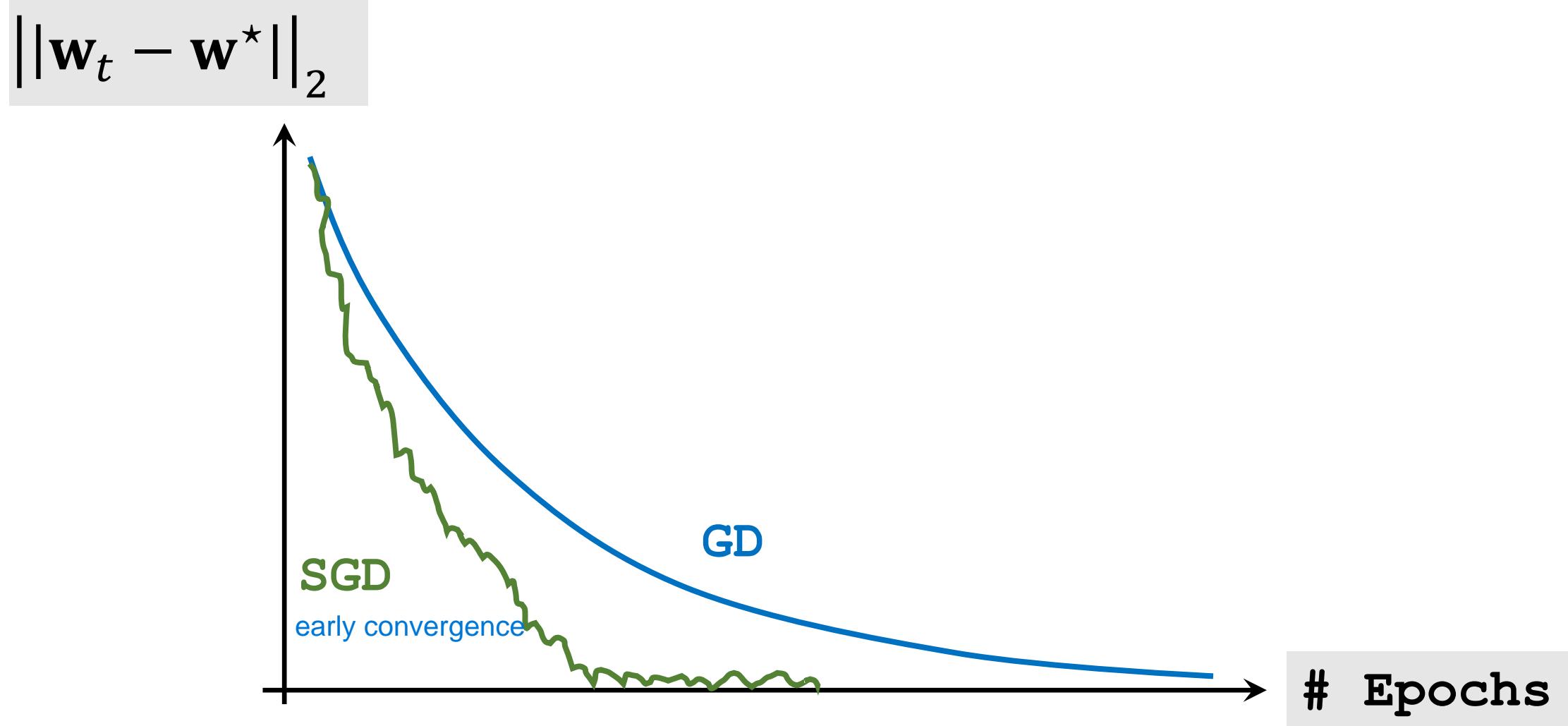
Accelerated GD

Stochastic GD

GD versus SGD



GD versus SGD



Training and Prediction

- Training:

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \sum_j l(y_j \mathbf{w}^T \mathbf{x}_j), \text{ where } l(z) = \log(1 + e^{-z}).$$

- For a test feature vector $\mathbf{x}' \in \mathbb{R}^d$, make prediction by

$$\operatorname{sign}(\mathbf{x}'^T \mathbf{w}^*).$$

Summary

- Logistic regression model for *linear binary* classification.

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_j l(y_j \mathbf{w}^T \mathbf{x}_j), \text{ where } l(z) = \log(1 + e^{-z}).$$

- Compute the gradient using vector derivatives and the chain rule.
- Gradient-based algorithms: GD, AGD, SGD, etc.
- Make prediction using $\operatorname{sign}(\mathbf{x}'^T \mathbf{w}^*)$.

Evaluate Binary Classification

Evaluate Binary Classification

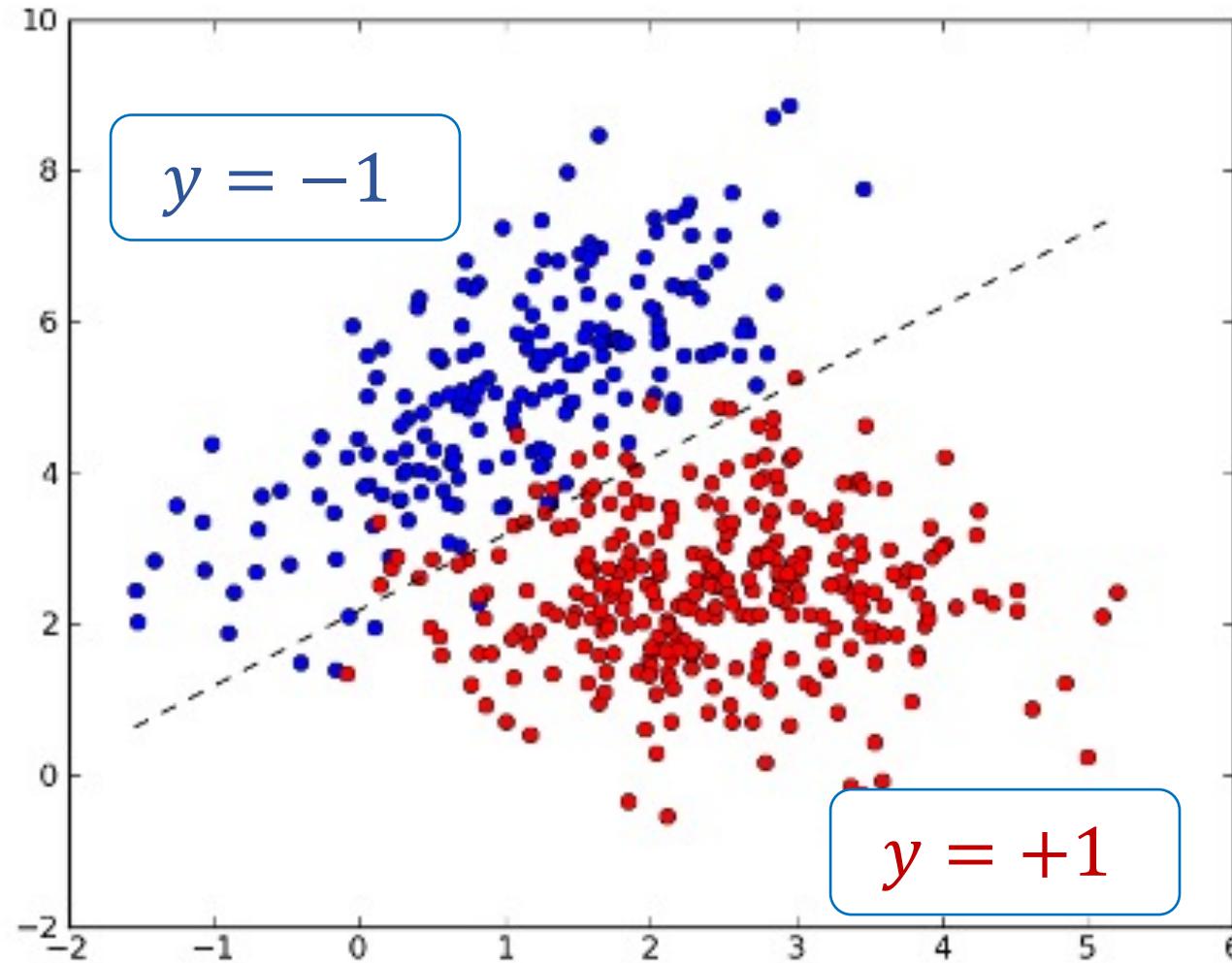
- Error Rate =
$$\frac{\text{\# Classification Errors}}{\text{\# Samples}}$$
- Accuracy = 1 - Error Rate

Evaluate Binary Classification

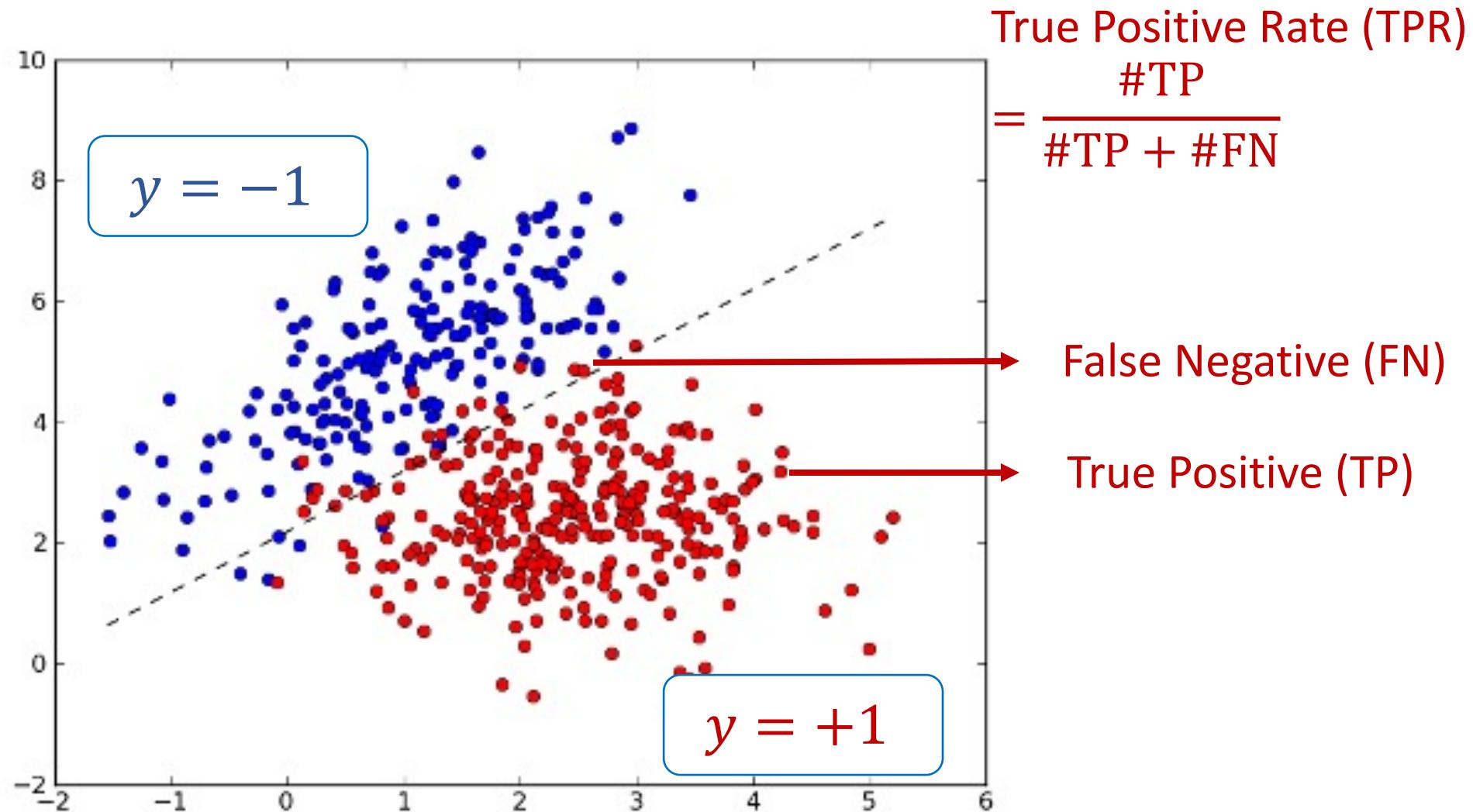
- Error Rate =
$$\frac{\text{\# Classification Errors}}{\text{\# Samples}}$$
- Accuracy = 1 - Error Rate

Error rate and Accuracy are not meaningful in class-imbalanced problems.

Evaluate Binary Classification



Evaluate Binary Classification



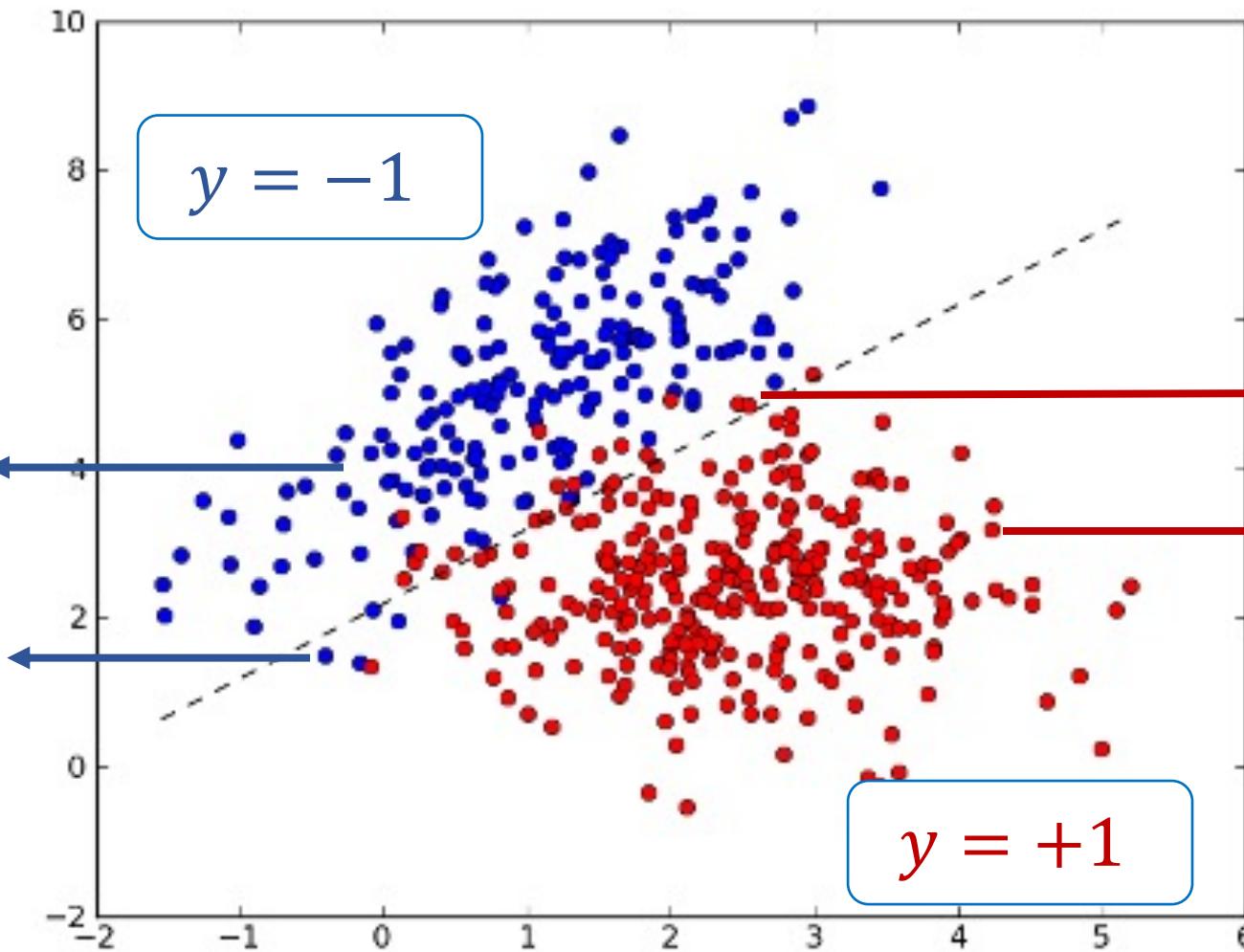
Evaluate Binary Classification

False Positive Rate (FPR)

$$= \frac{\#FP}{\#FP + \#TN}$$

True Negative (TN)

False Positive (FP)



True Positive Rate (TPR)

$$= \frac{\#TP}{\#TP + \#FN}$$

False Negative (FN)

True Positive (TP)

Evaluate Binary Classification

False Positive Rate (FPR)

$$= \frac{\#FP}{\#FP + \#TN}$$

True Negative (TN)

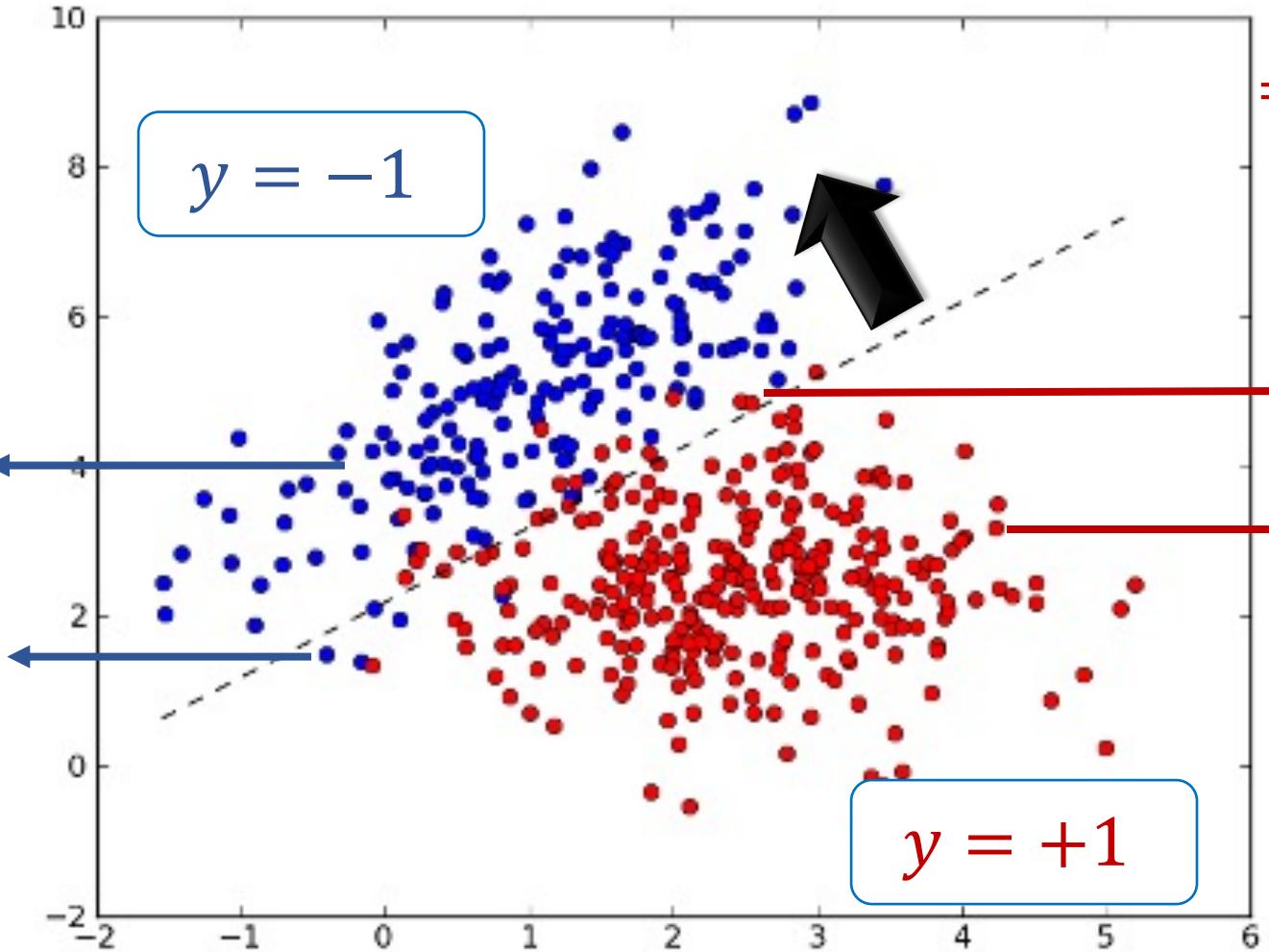
False Positive (FP)

True Positive Rate (TPR)

$$= \frac{\#TP}{\#TP + \#FN}$$

False Negative (FN)

True Positive (TP)



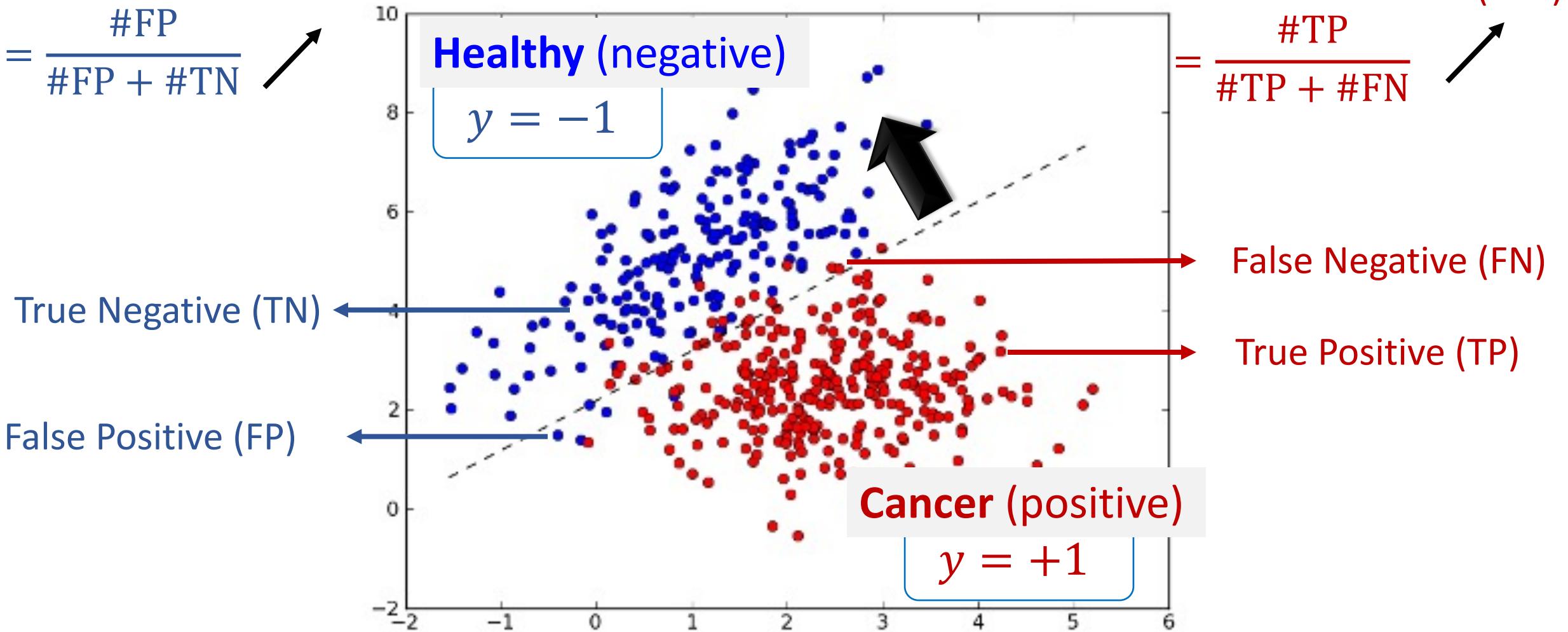
Evaluate Binary Classification

False Positive Rate (FPR)

$$= \frac{\#FP}{\#FP + \#TN}$$

True Negative (TN)

False Positive (FP)



Evaluate Binary Classification

False Positive Rate (FPR)

$$= \frac{\#FP}{\#FP + \#TN}$$

True Positive Rate (TPR)

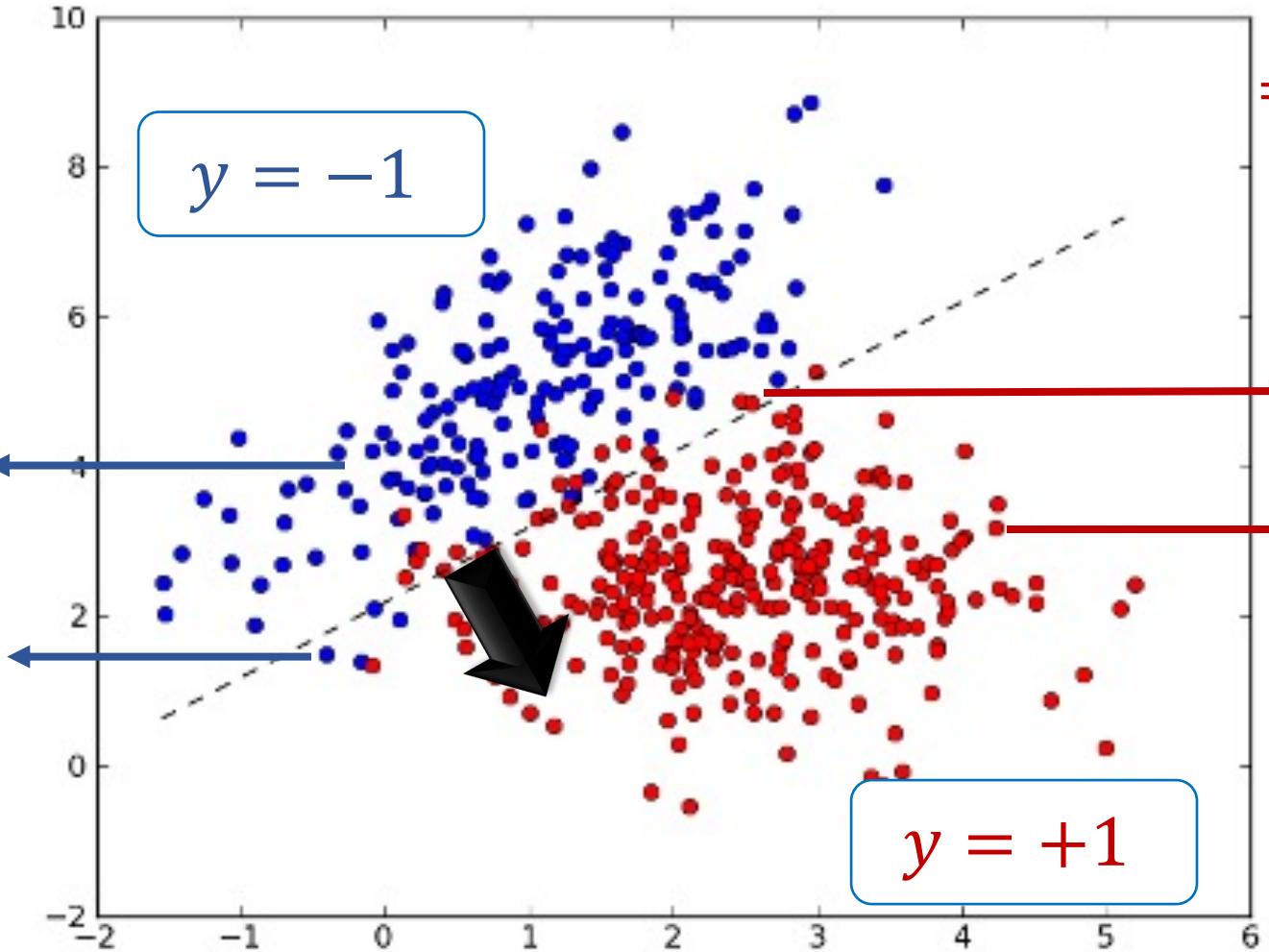
$$= \frac{\#TP}{\#TP + \#FN}$$

True Negative (TN)

False Positive (FP)

False Negative (FN)

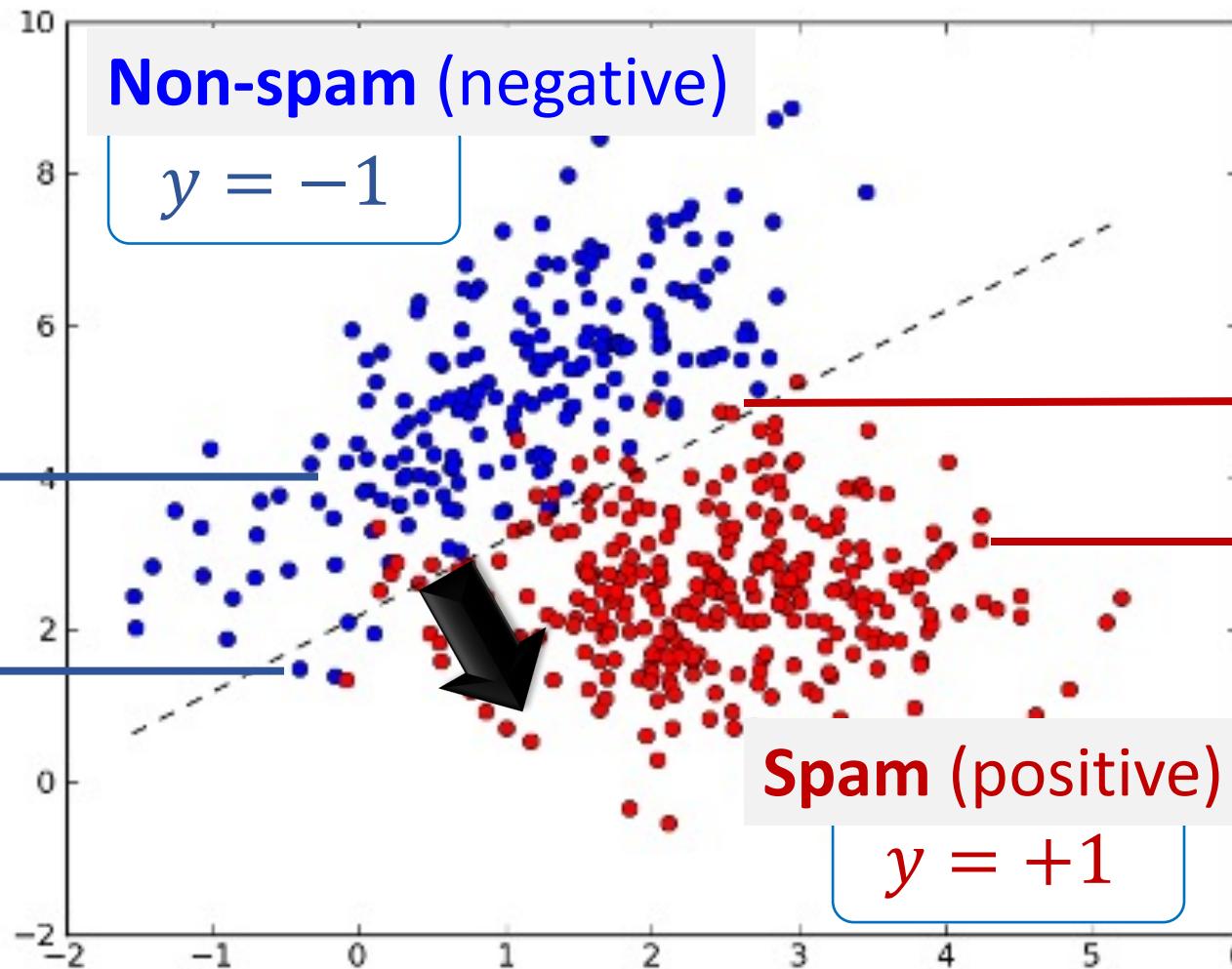
True Positive (TP)



Evaluate Binary Classification

False Positive Rate (FPR)

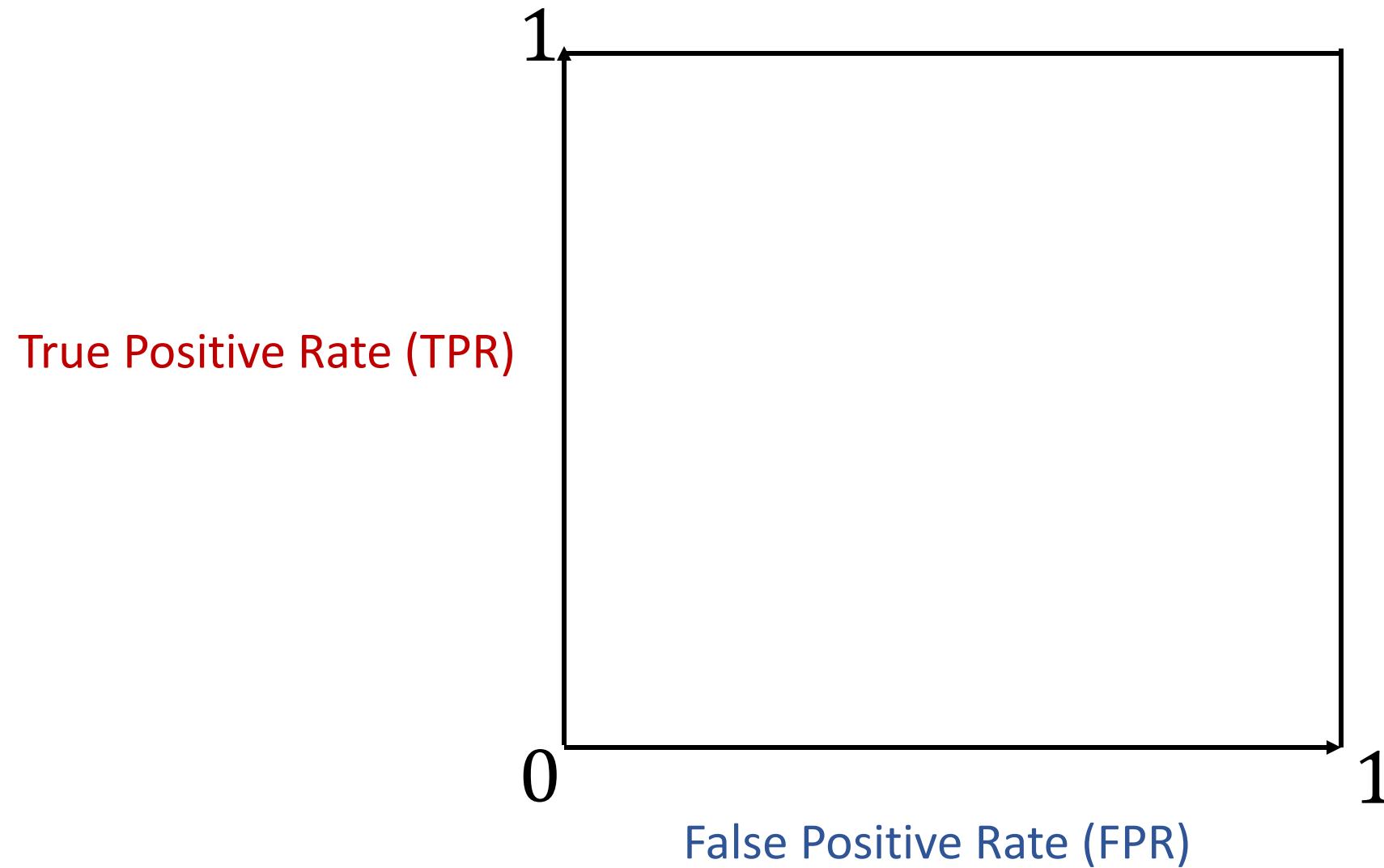
$$= \frac{\#FP}{\#FP + \#TN}$$



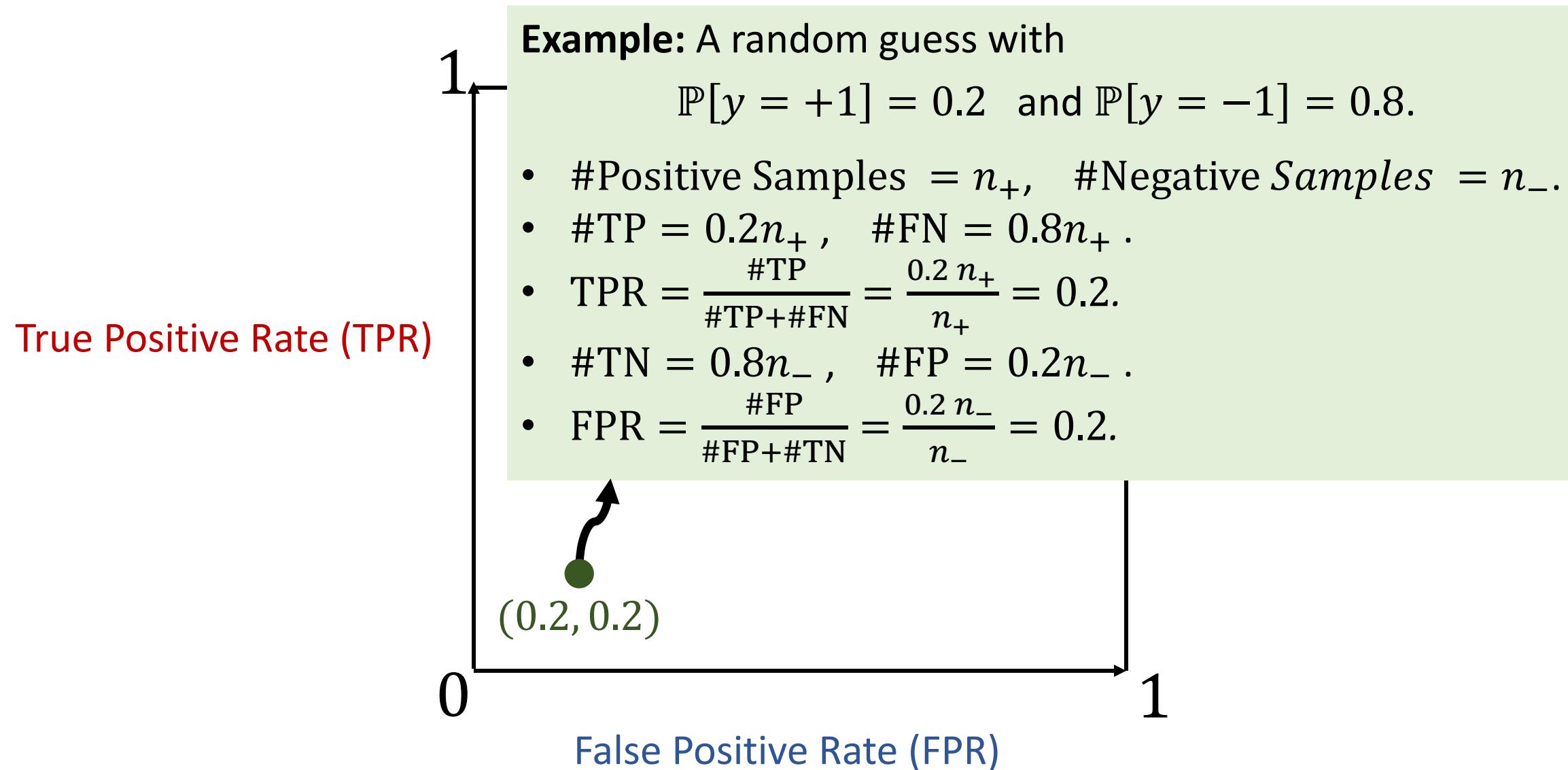
True Positive Rate (TPR)

$$= \frac{\#TP}{\#TP + \#FN}$$

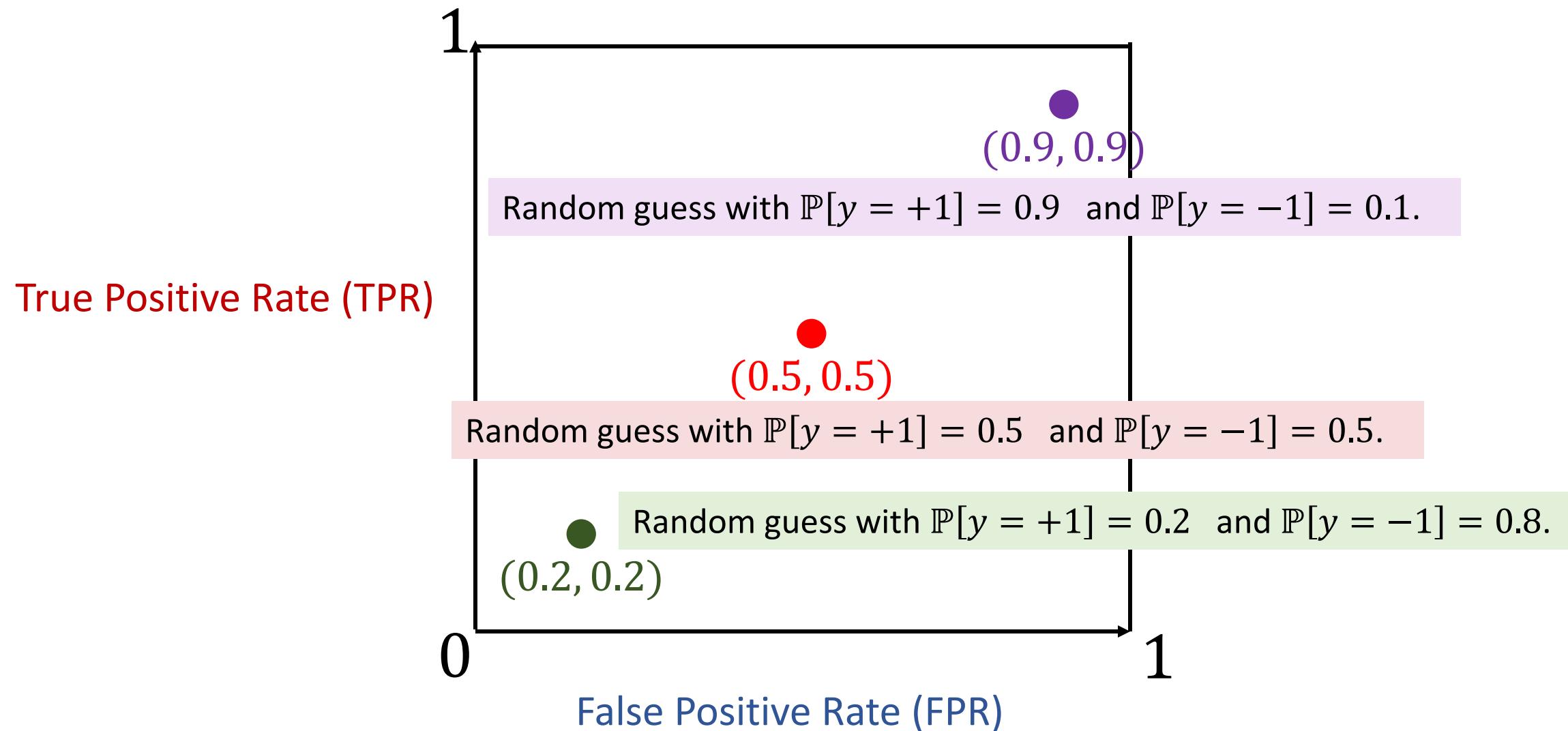
Receiver Operating Characteristic (ROC) Curve



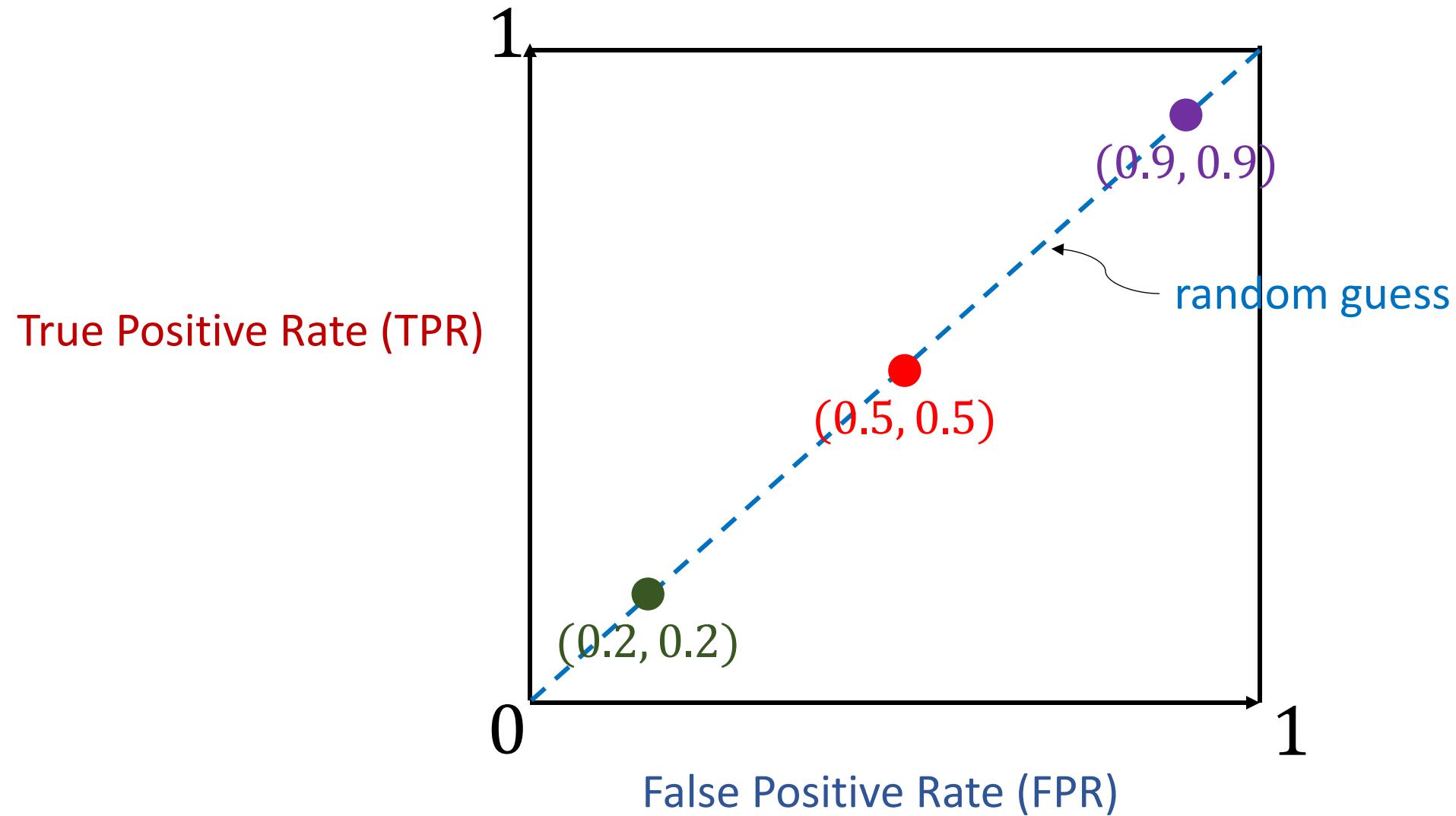
Receiver Operating Characteristic (ROC) Curve



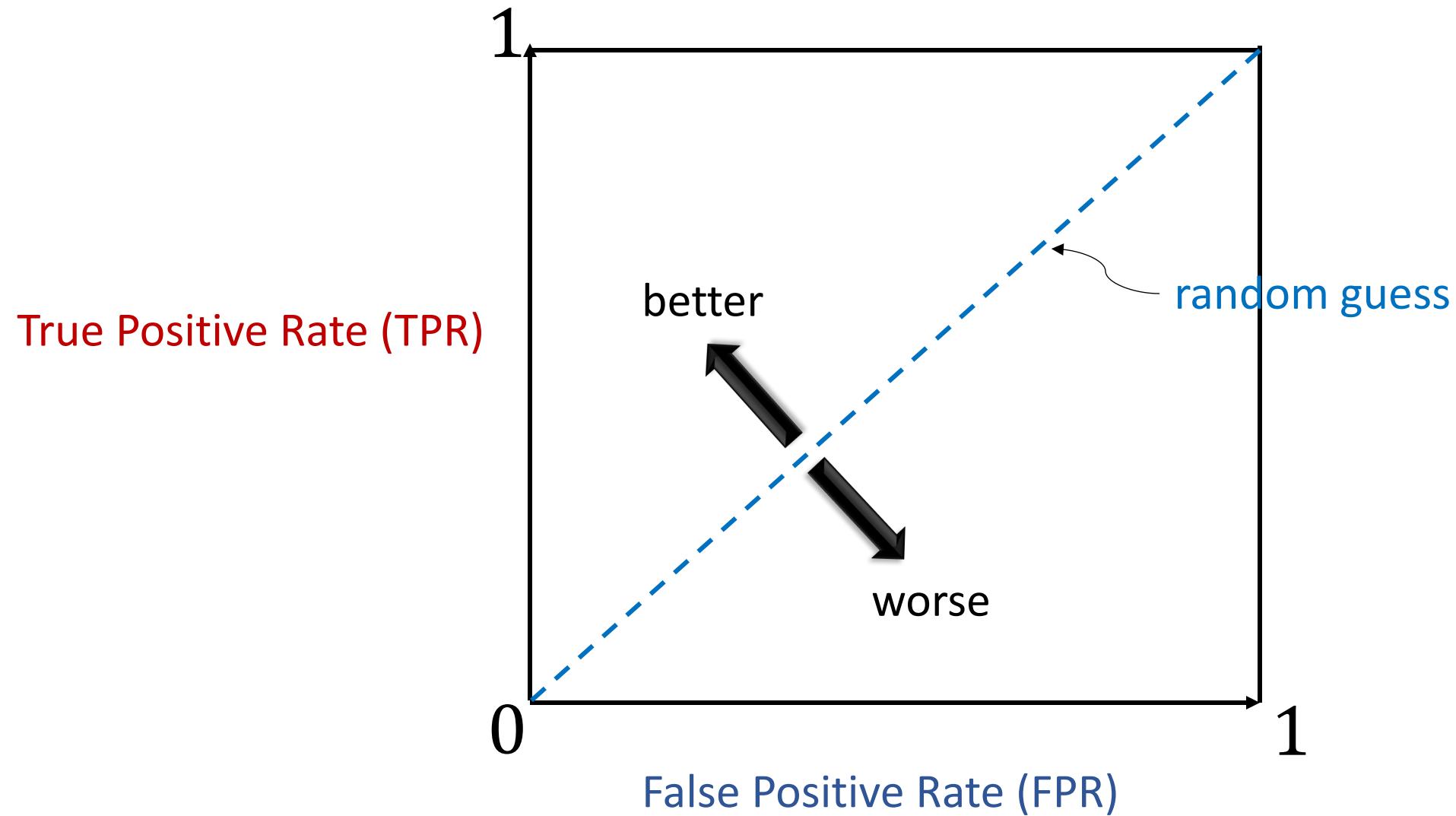
Receiver Operating Characteristic (ROC) Curve



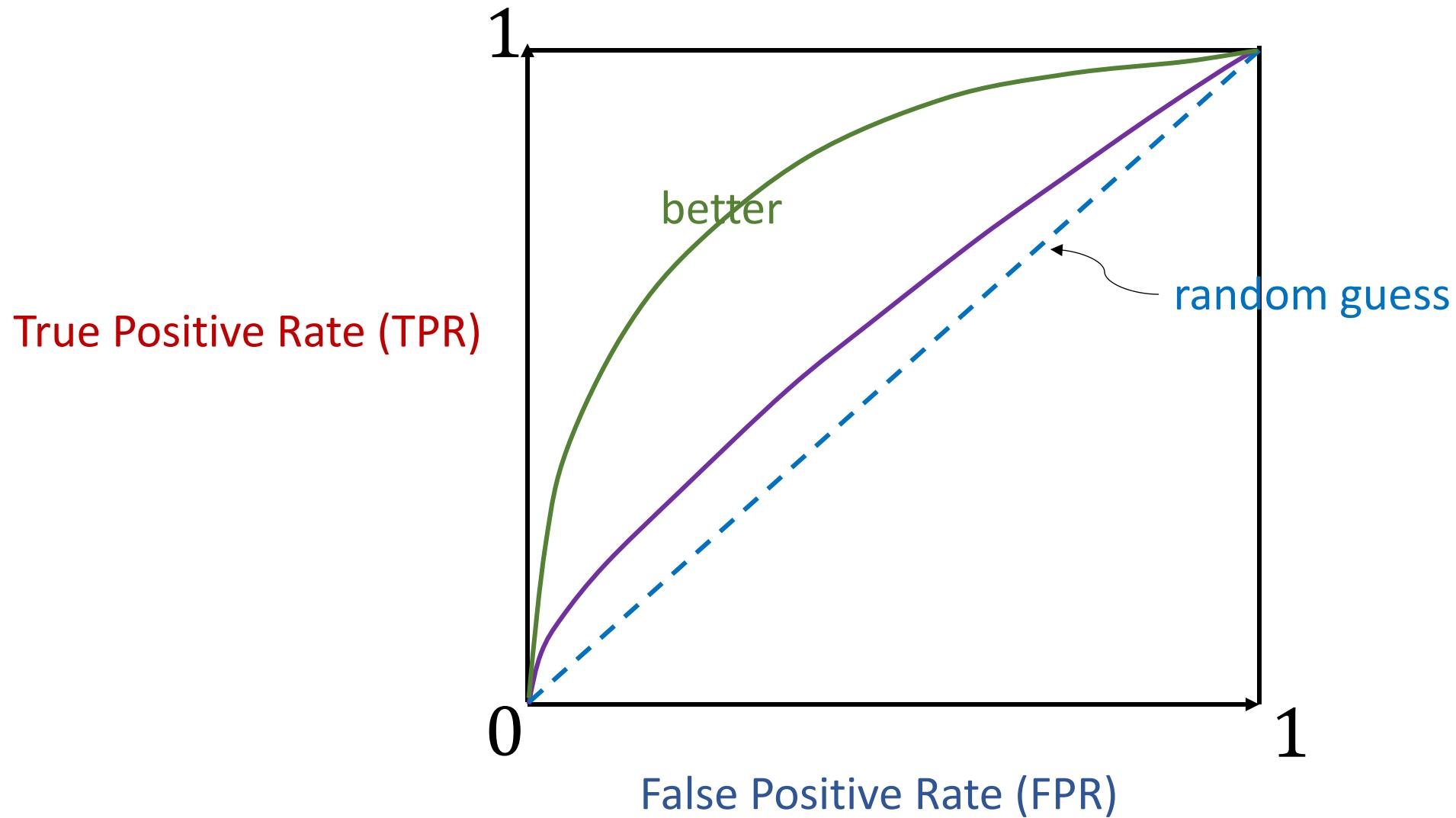
Receiver Operating Characteristic (ROC) Curve



Receiver Operating Characteristic (ROC) Curve



Receiver Operating Characteristic (ROC) Curve



Thank you!