

# Welcome to CME 250 Introduction to Machine Learning!

Spring 2020 – Online version  
April 28th 2020

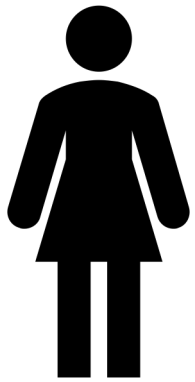


# Today's schedule: Classification

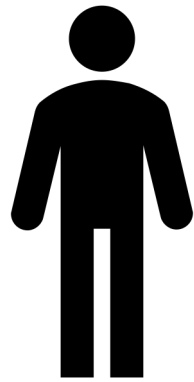
- Why does the distinction between regression and classification matters?
- Classification looking at  $Y$  as a random variable:
  - Logistic regression as a Generalized linear model
- Classification finding boundaries:
  - Support Vector Machines
- How to measure classification success?
  - Confusion Matrix

# Let's get to know each other...

Breakout room



You



Another student

Name

Location

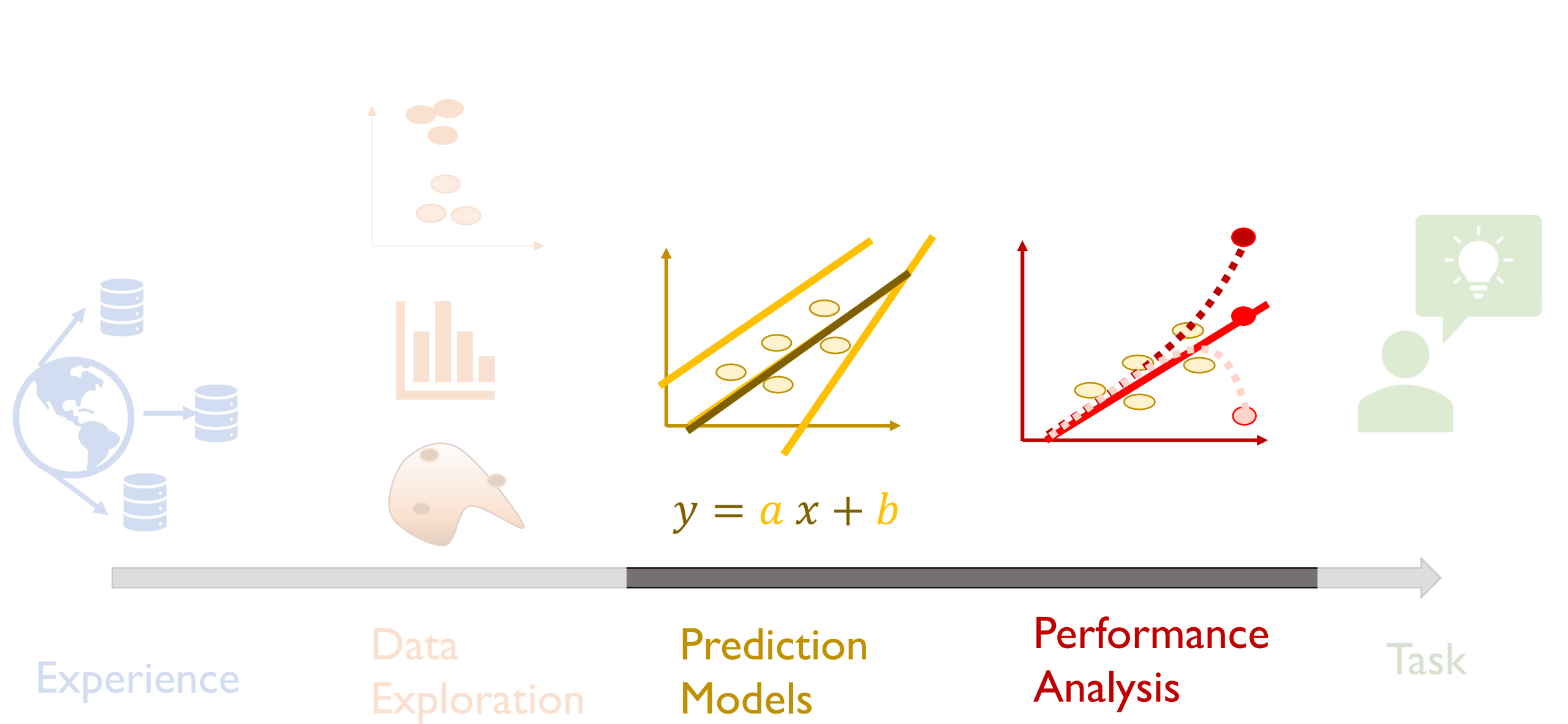
Department

Year

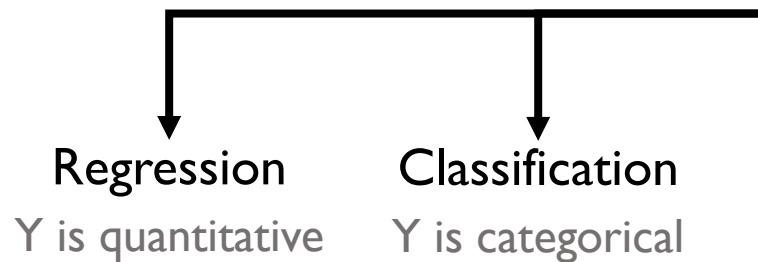
How was Part I Project?  
Interesting/unexpected/  
unforgettable lessons or  
insights.

**3 mins**

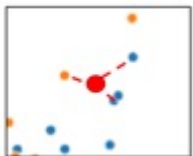
Chat/Audio/Video



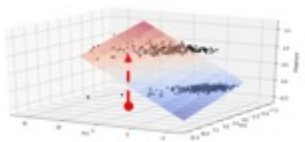
# Last week recap



KNN



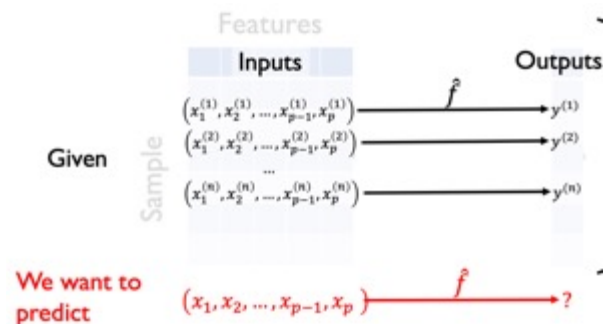
Linear  
Regression



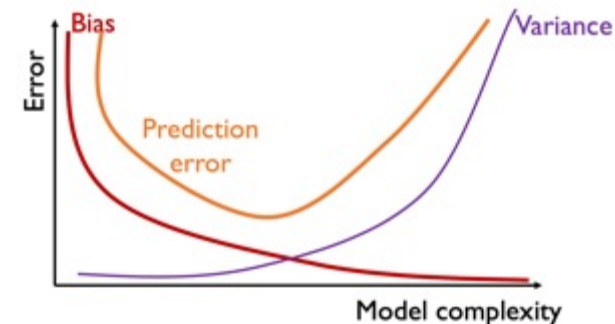
Using dummy  
variables?

## Supervised Learning

Learn from examples

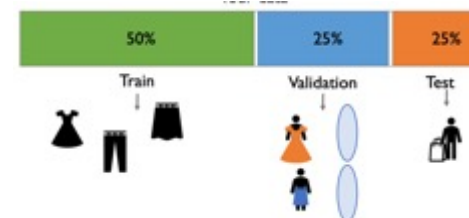


## Model Selection



## Cross-Validation

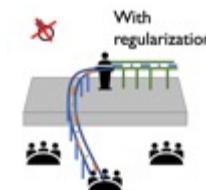
Estimate prediction  
error



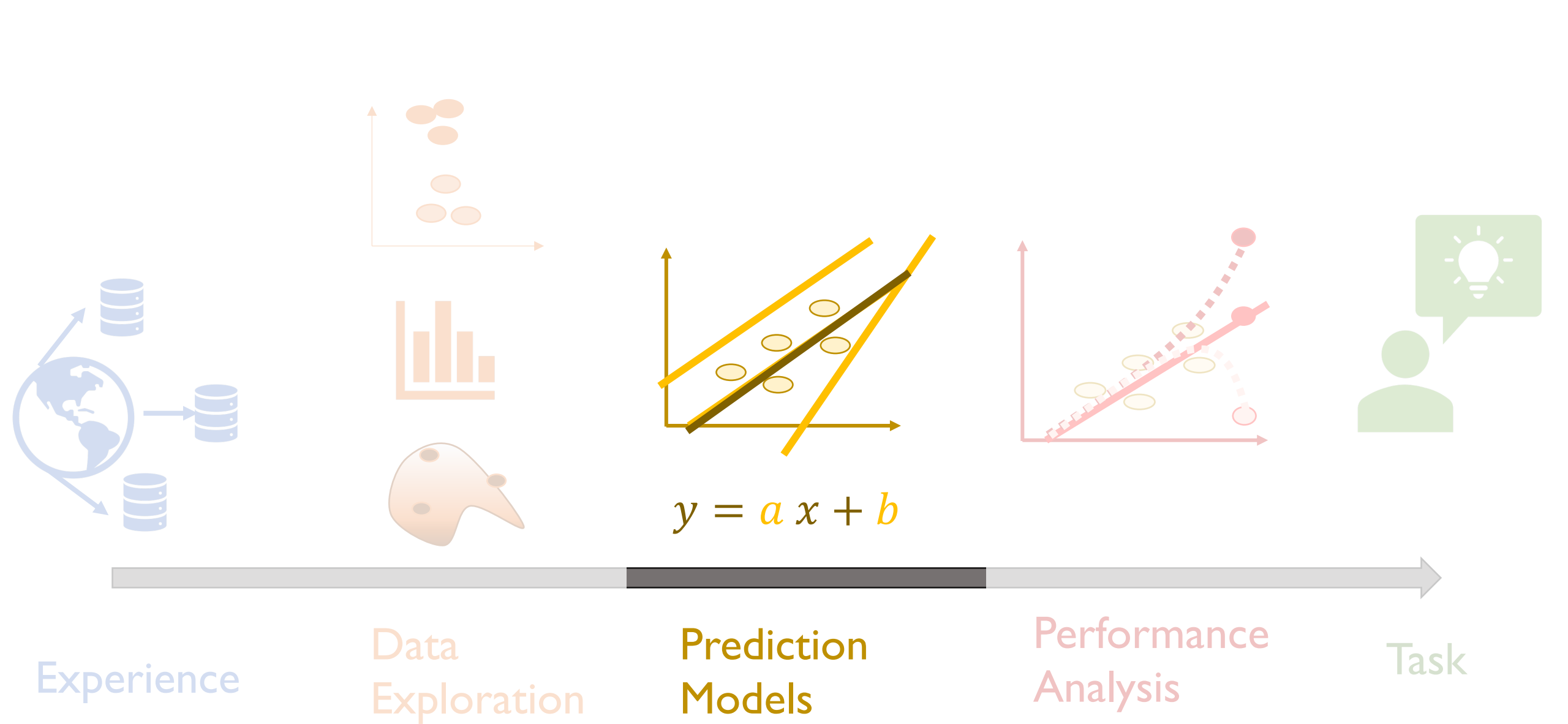
K-fold CV, LOOCV

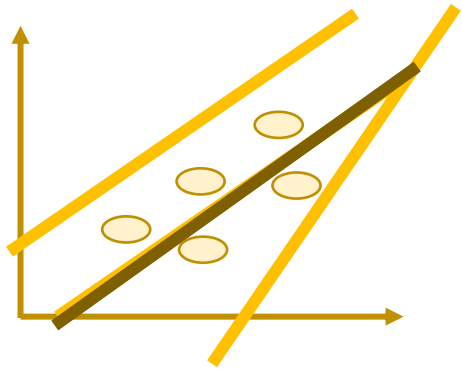
## Regularization

Control complexity:  
Hyperparameters



Ridge, Lasso





$$y = ax + b$$

Prediction  
Models

## Supervised Learning Part II: Prediction Models for Classification

*Introduction to Statistical Learning*

Chapter 4: Classification

Chapter 9: Support Vector Machines

*Elements Statistical Learning*

Chapter 3.2: Linear Methods for Classification

Chapter 12: Support Vector Machines

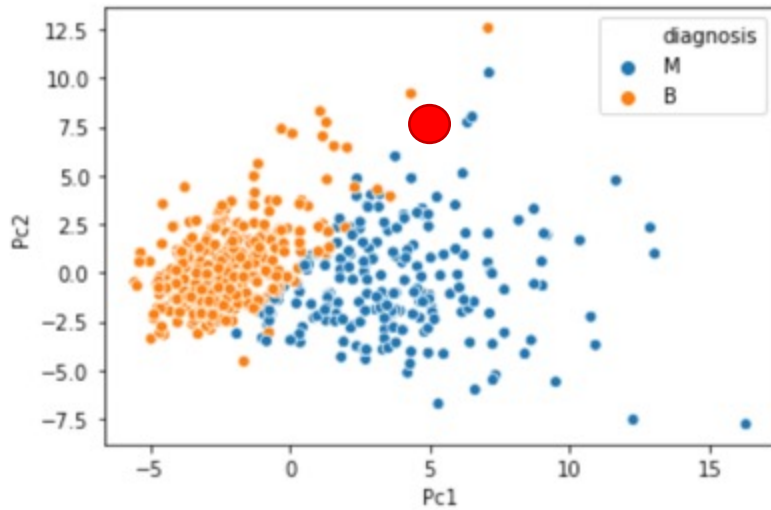
More on Generalized Linear Models

*Bayesian and Frequentist Regression Methods.*

Jon Wakefield, 2013

Chapter 6.3: Generalized Linear Models

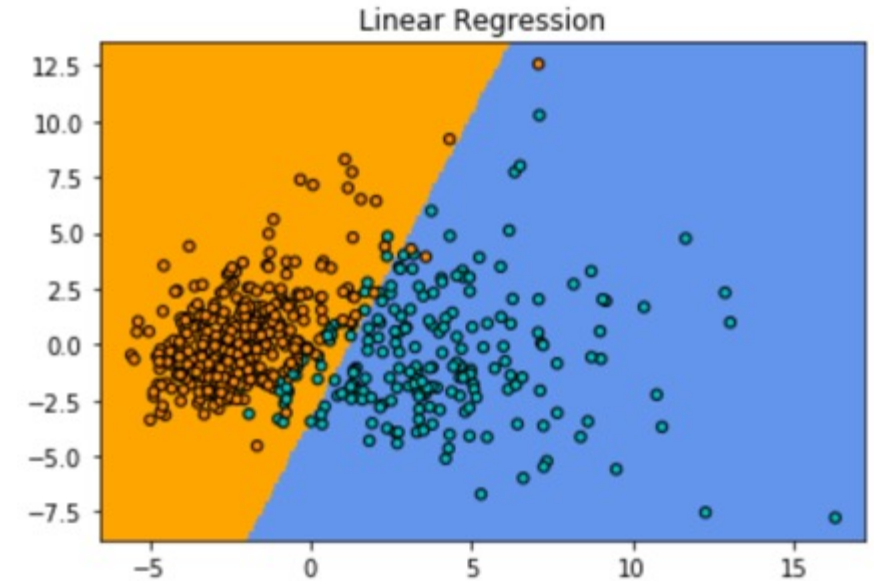
# Breast Cancer Wisconsin (Diagnostic) Dataset



What is the diagnosis for this sample?

Use dummy variable for Y

$$\begin{array}{l} M \rightarrow 1 \\ \text{---} \text{---} \text{---} \text{---} 0.5 \\ B \rightarrow 0 \end{array}$$



$$Y \approx f(x) = \beta_0 + \beta_1 X_{PC1} + \beta_2 X_{PC2}$$

What if we have more than 2 categories?



# Linear regression with more than 2 categories

Option 1)

	Y
Cat. 1	0
-----	0.5
Cat. 2	1
-----	1.5
Cat. 3	2

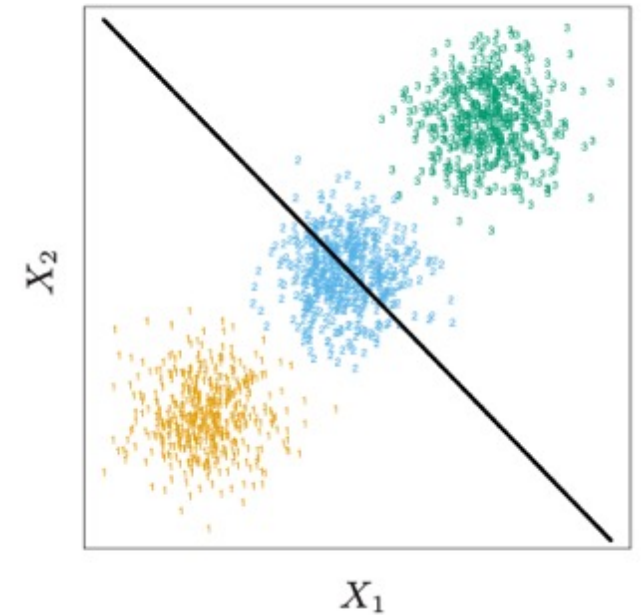
✗ Imposes order in categories

Option 2)

	$Y_1$	$Y_2$	$Y_3$
Cat. 1	1	0	0
Cat. 2	0	1	0
Cat. 3	0	0	1

✗ Ignores category 2

Linear Regression

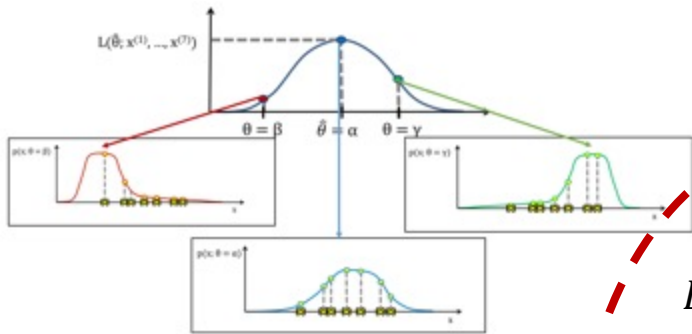


ESL Fig 4.2

**We need a different approach!**

# How can we extend Linear Regression?

LR is Maximum Likelihood estimator



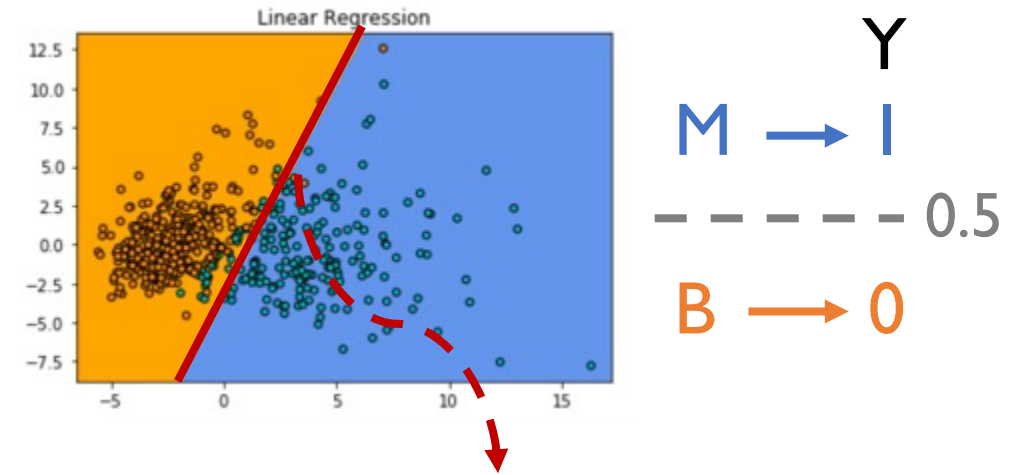
$Y \sim \text{Normal}$

$$E[Y|X] = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

Find a better distribution for Y  
categorical

Logistic Regression

LR creates separating hyperplanes



Optimize the hyperplane

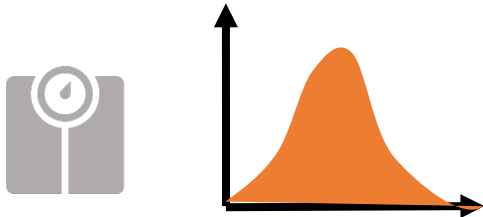
Support Vector Machines

# Extend LR: Generalized Linear Models

## Exponential Family

$$p(y | \theta, \alpha) = \exp \left( \frac{y\theta - b(\theta)}{\alpha} + c(y, \alpha) \right)$$

Normal distribution  
 $N(\mu, \sigma^2)$



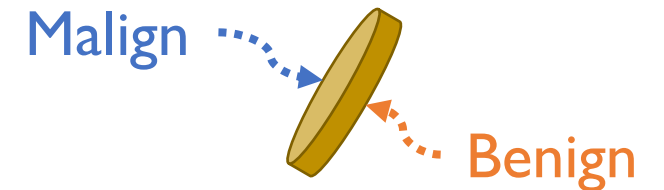
$$p(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(y - \mu)^2}{2\sigma^2} \right)$$

Poisson distribution  
 $Poisson(\lambda)$



$$p(y | \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!}$$

Bernoulli distribution  
 $Bernoulli(p)$



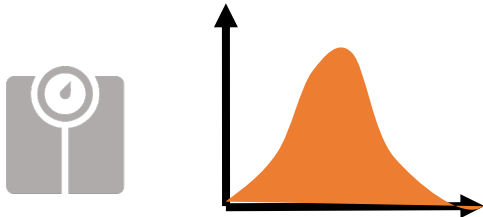
$$p(y | \lambda) = p^y (1 - p)^{1-y}$$

# Extend LR: Generalized Linear Models

Exponential Family

$$E(y|\theta, \alpha) = b'(\theta)$$

Normal distribution  
 $N(\mu, \sigma^2)$



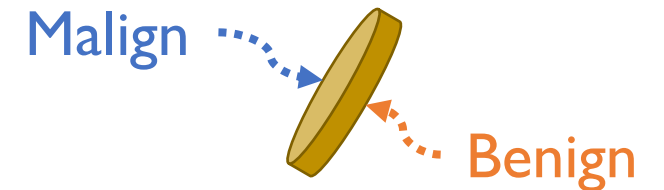
$$b(\theta) = \theta^2/2$$
$$\mu = E[Y|\theta, \alpha] = \theta$$

Poisson distribution  
 $Poisson(\lambda)$



$$b(\theta) = \exp(\theta)$$
$$\lambda = E[Y|\theta, \alpha] = \exp(\theta)$$

Bernoulli distribution  
 $Bernoulli(p)$



$$b(\theta) = \log(1 + \exp(\theta))$$
$$p = E[Y|\theta, \alpha] = \frac{\exp(\theta)}{1 + \exp(\theta)}$$

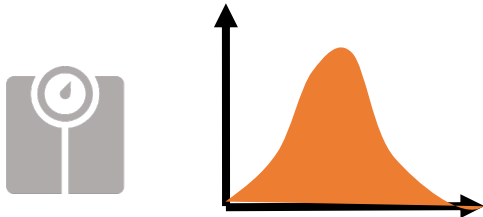
Calculate a linear regression of  $\theta = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$

# Extend LR: Generalized Linear Models

Exponential Family

$$E(y | \theta, \alpha) = b'(\theta)$$

Normal distribution  
 $N(\mu, \sigma^2)$



Linear Regression

$$y \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

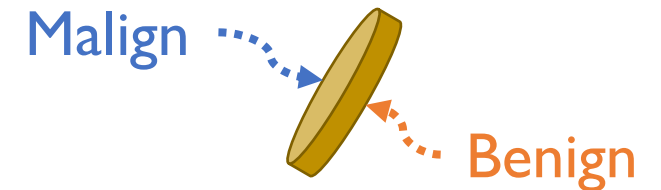
Poisson distribution  
 $Poisson(\lambda)$



Log-Linear Regression

$$\log(y) \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

Bernoulli distribution  
 $Bernoulli(p)$



Logistic Regression

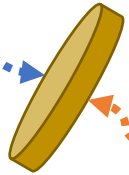
$$\log\left(\frac{p}{1-p}\right) \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

# Logistic Regression

Bernoulli distribution

$Bernoulli(p)$

Malign



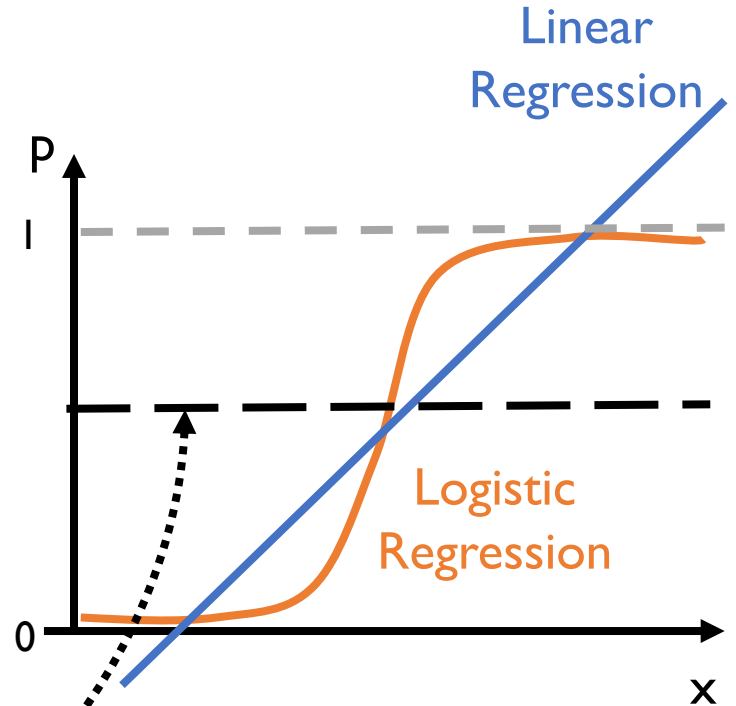
Benign

Log-odds / Logit

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$$p = \frac{\exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}{1 + \exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}$$

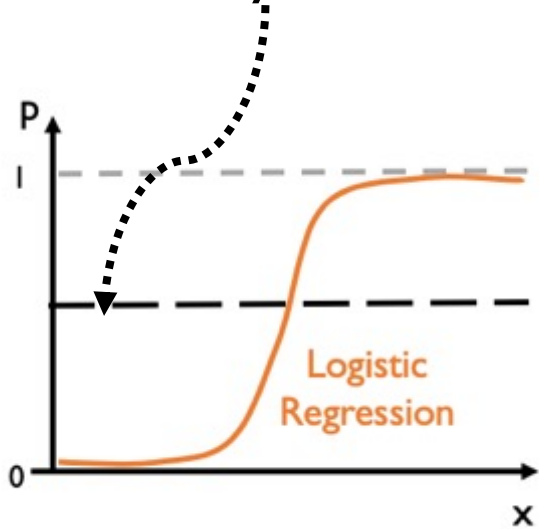
Sigmoid



To classify, we specify probability threshold

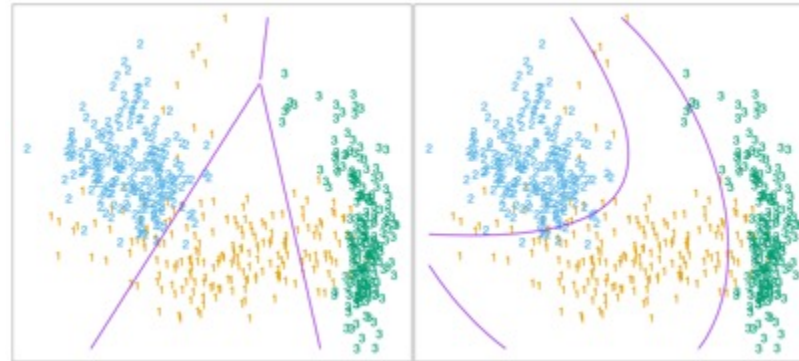
# Challenges of Logistic Regression

## Probability Threshold



Hyperparameter  
Usually 0.5 (not always)

## Linear Decision Boundary



ESL Fig 4.1

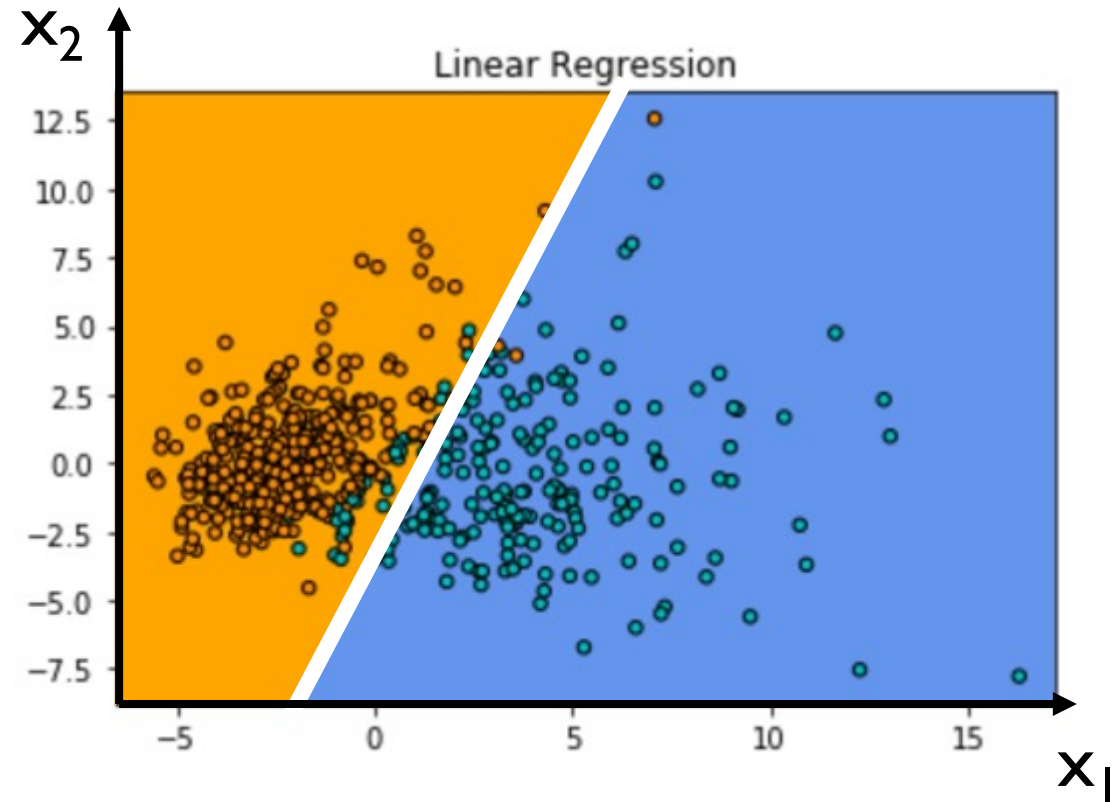
As in Linear regression:  
Add additional features  
 $X_1, X_2, X_1X_2, X_1^2, X_2^2$

## Multiple Categories

$$p_1 = \frac{\exp(\beta_{01} + \beta_1^T X)}{1 + \sum_{l=1}^{K-1} \exp(\beta_{0l} + \beta_l^T X)}$$
$$\dots$$
$$p_{K-1} = \frac{\exp(\beta_{0K-1} + \beta_{K-1}^T X)}{1 + \sum_{l=1}^{K-1} \exp(\beta_{0l} + \beta_l^T X)}$$
$$p_K = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(\beta_{0l} + \beta_l^T X)}$$

# Do we need probabilities or decision boundaries?

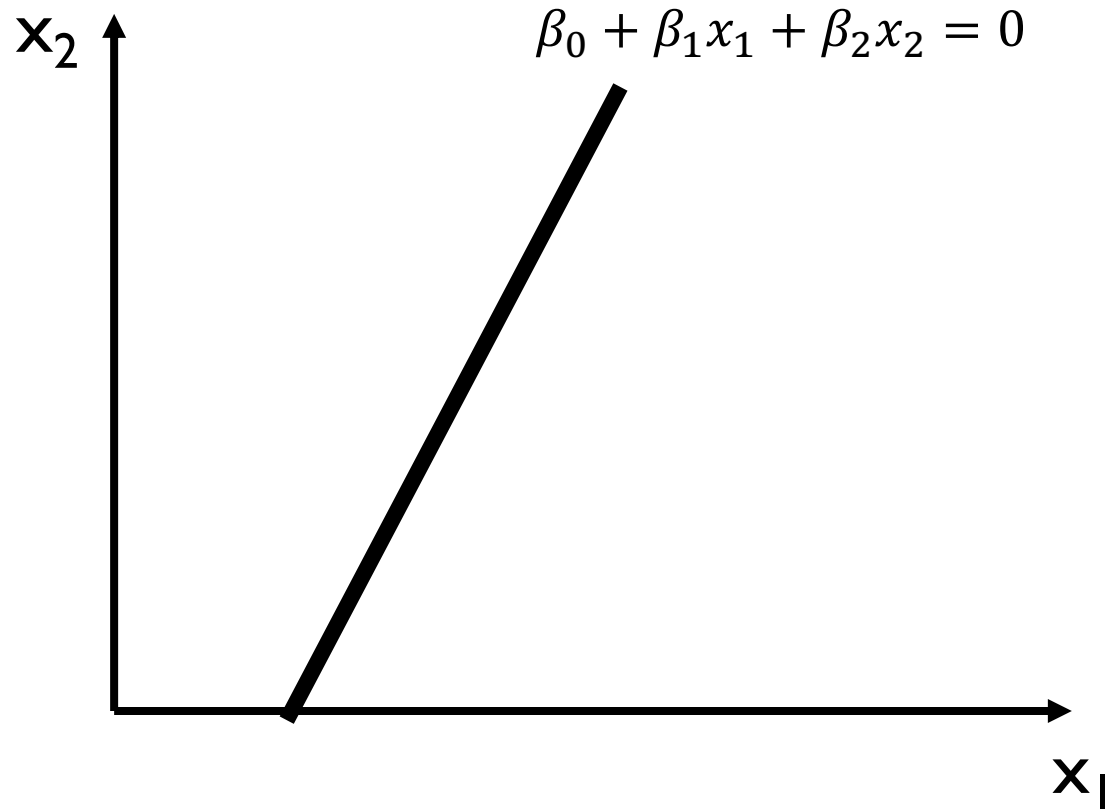
Find “optimal”  
boundary to  
separate classes





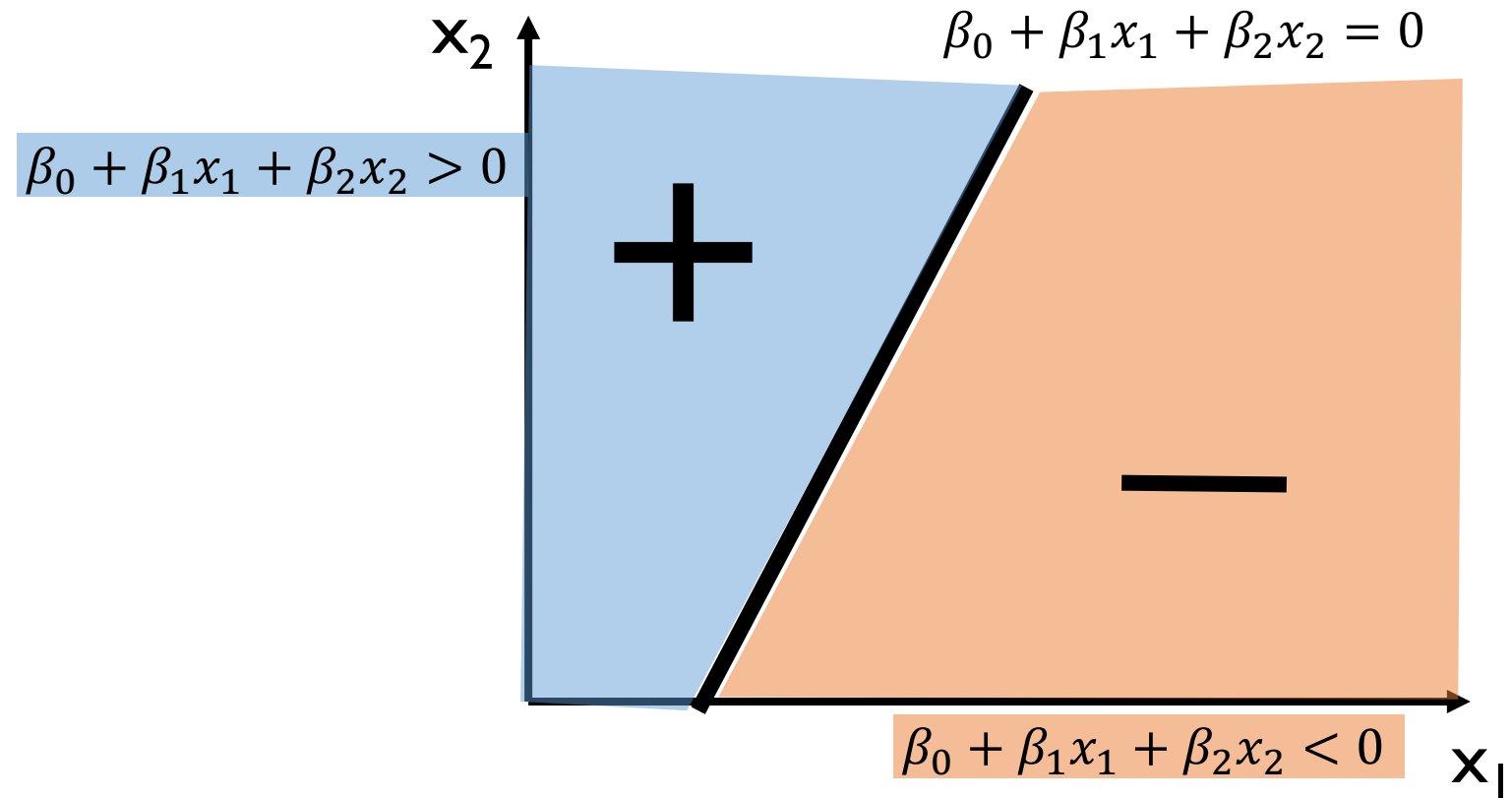
# Do we need probabilities or decision boundaries?

Find “optimal”  
boundary to  
separate classes

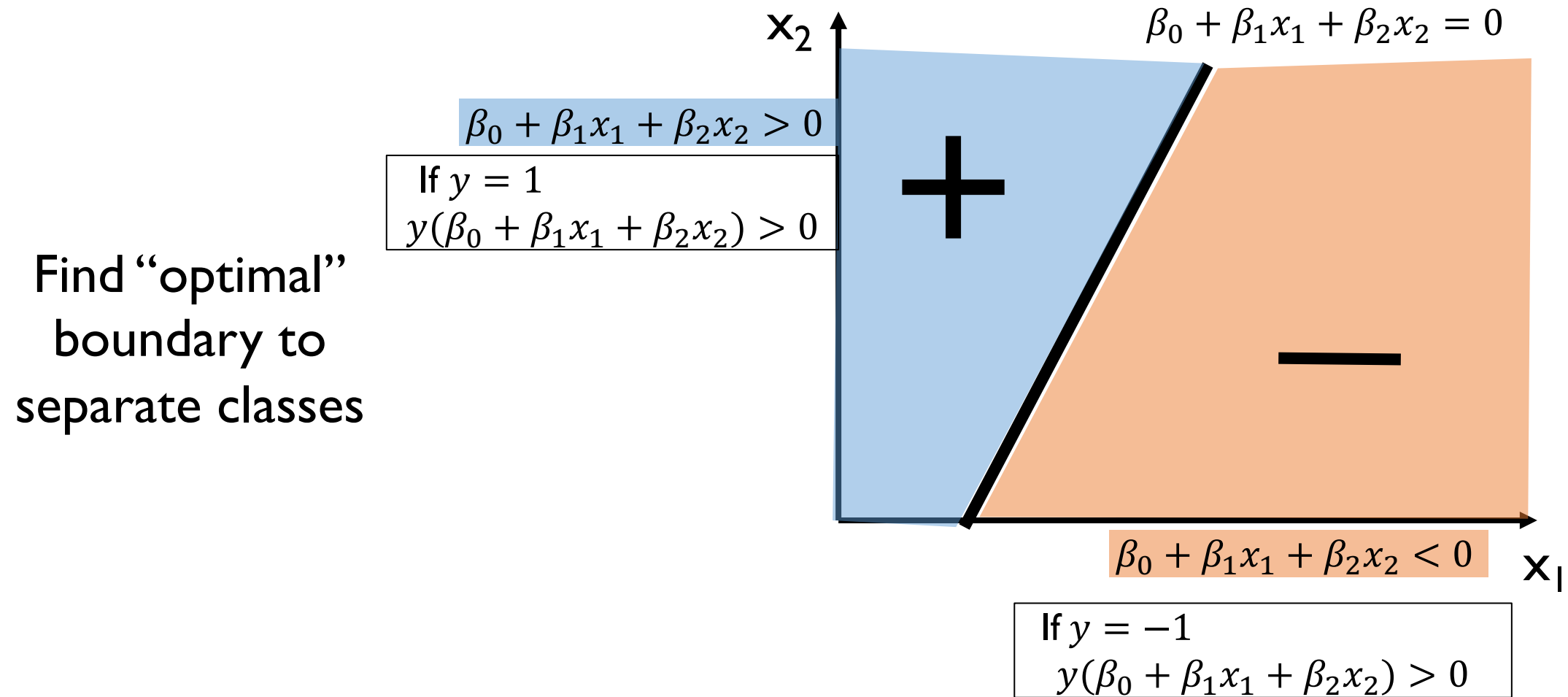


# Do we need probabilities or decision boundaries?

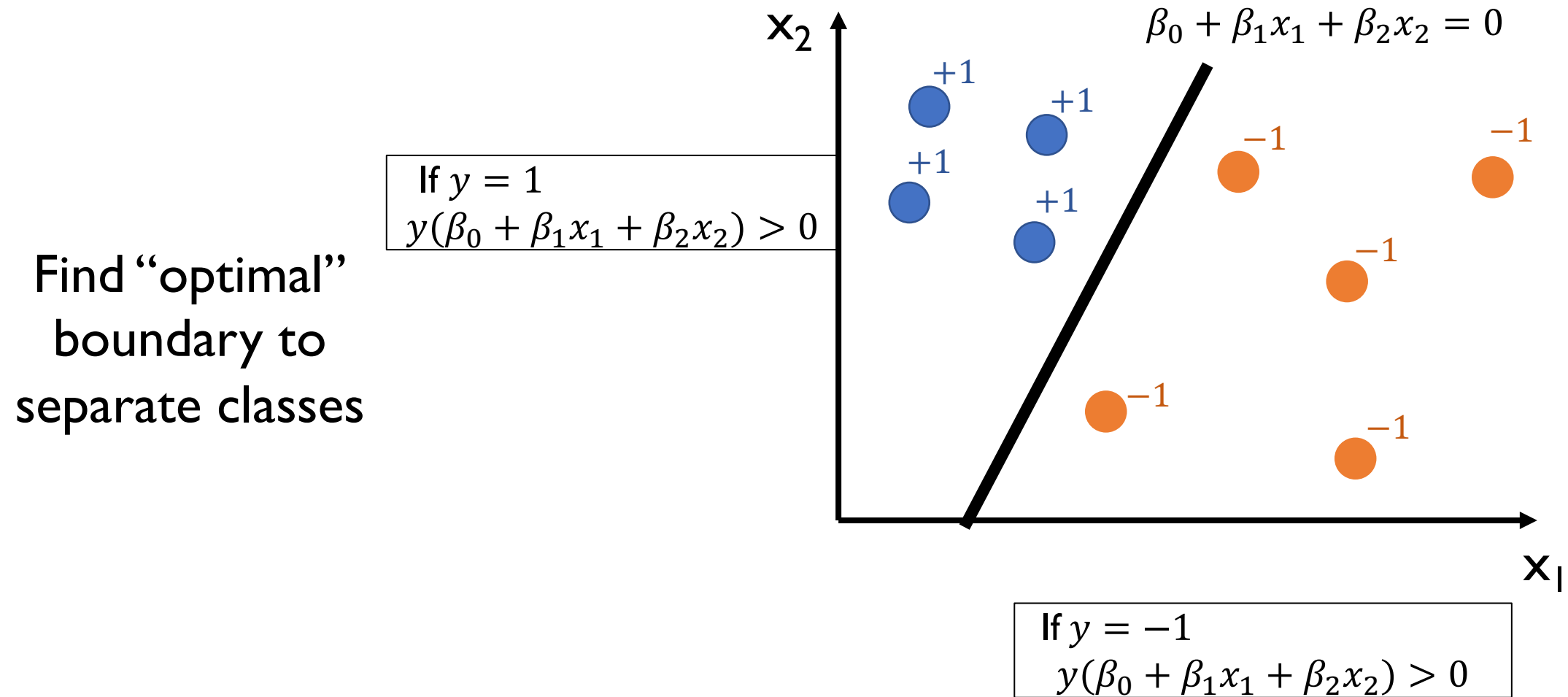
Find “optimal”  
boundary to  
separate classes



# Do we need probabilities or decision boundaries?

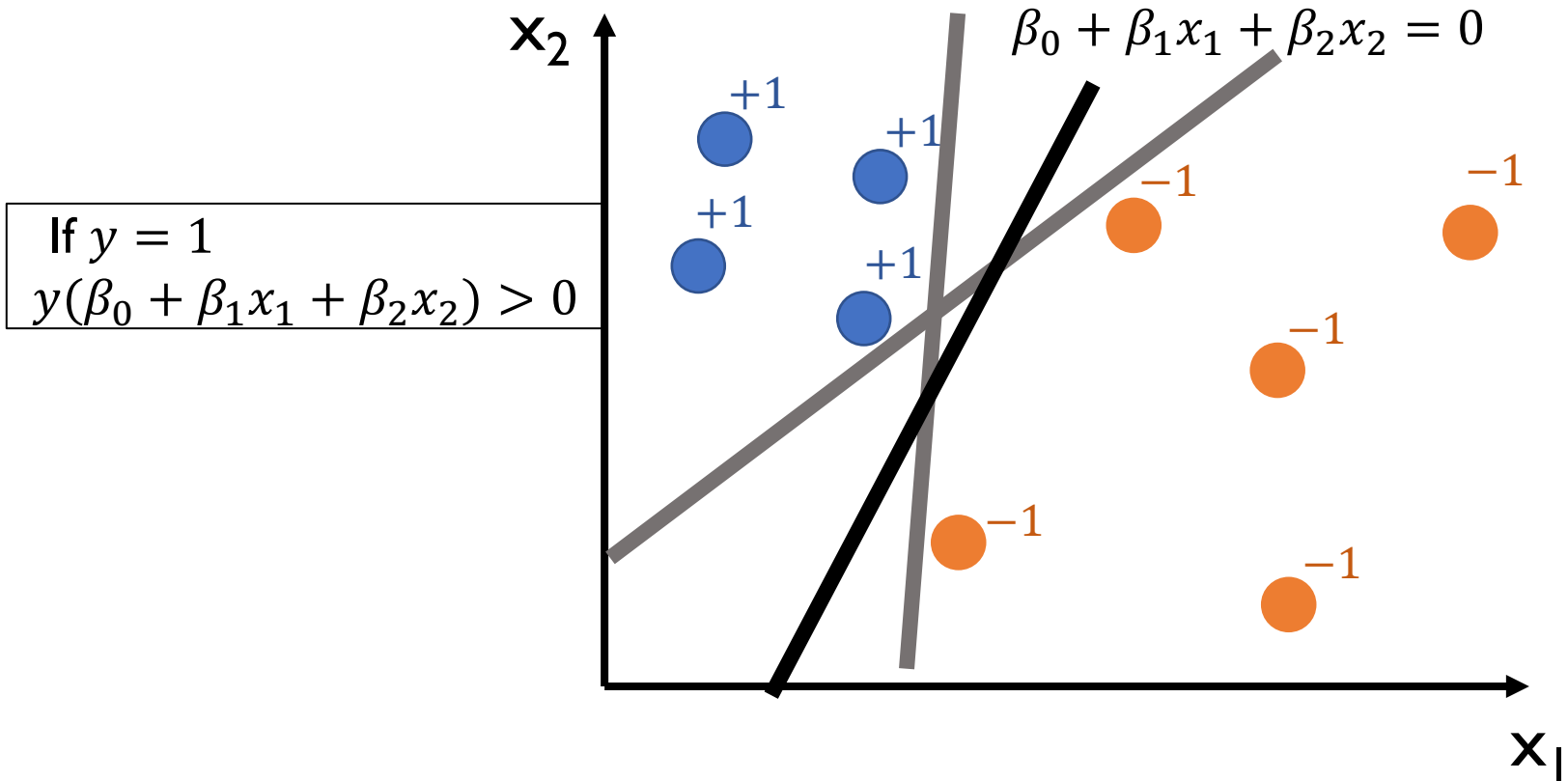


# Do we need probabilities or decision boundaries?



# Do we need probabilities or decision boundaries?

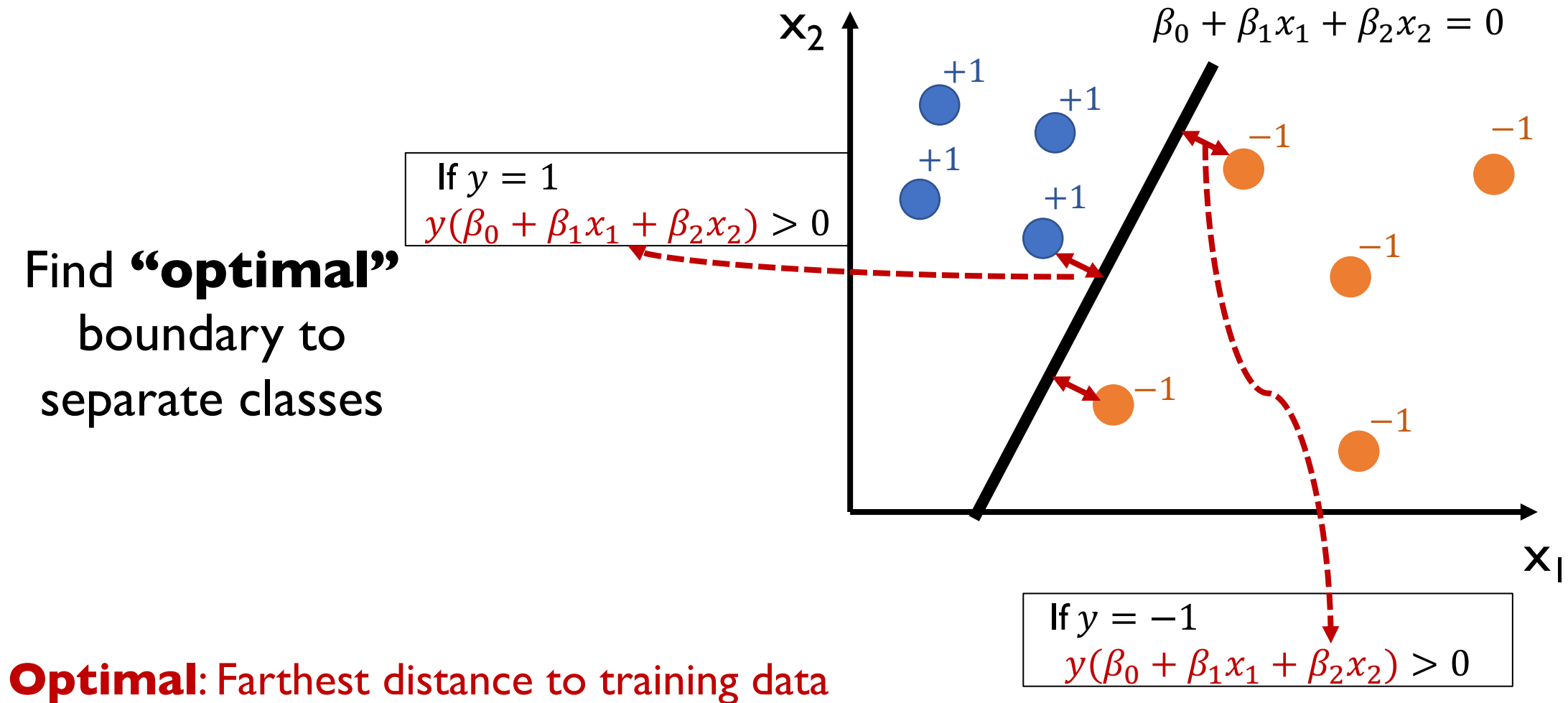
Find “**optimal**”  
boundary to  
separate classes



**Optimal:** Farthest distance to training data

$$\text{If } y = -1 \\ y(\beta_0 + \beta_1 x_1 + \beta_2 x_2) > 0$$

# Maximal Margin Classifier



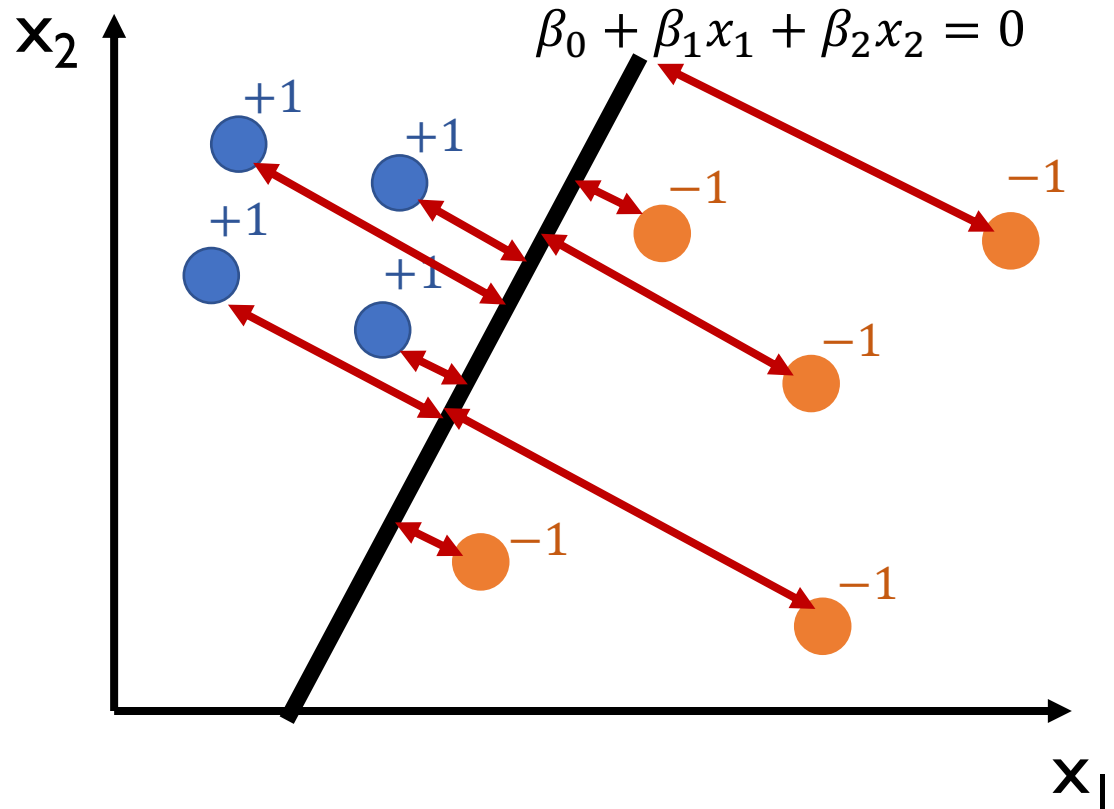
# Maximal Margin Classifier

$$\max_{\beta_0, \beta_1, \beta_2} M$$

$$\text{such that } \beta_0^2 + \beta_1^2 + \beta_2^2 = 1$$

For all training data

$$y^{(i)} (\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) \geq M$$



**Optimal:** Farthest distance to training data

# Maximal Margin Classifier

$$\max_{\beta_0, \beta_1, \beta_2} M$$

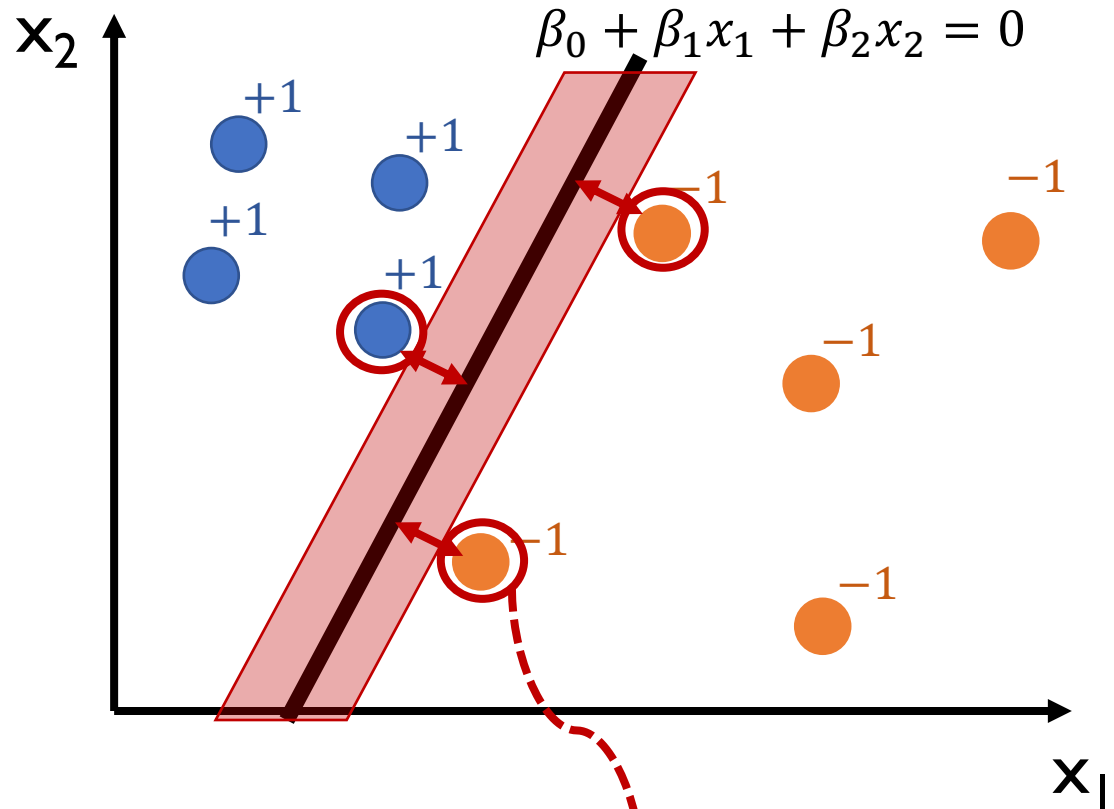
$$\text{such that } \beta_0^2 + \beta_1^2 + \beta_2^2 = 1$$

For all training data

$$y^{(i)} (\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) \geq M$$

Only points with equality  
matter

Support Vectors



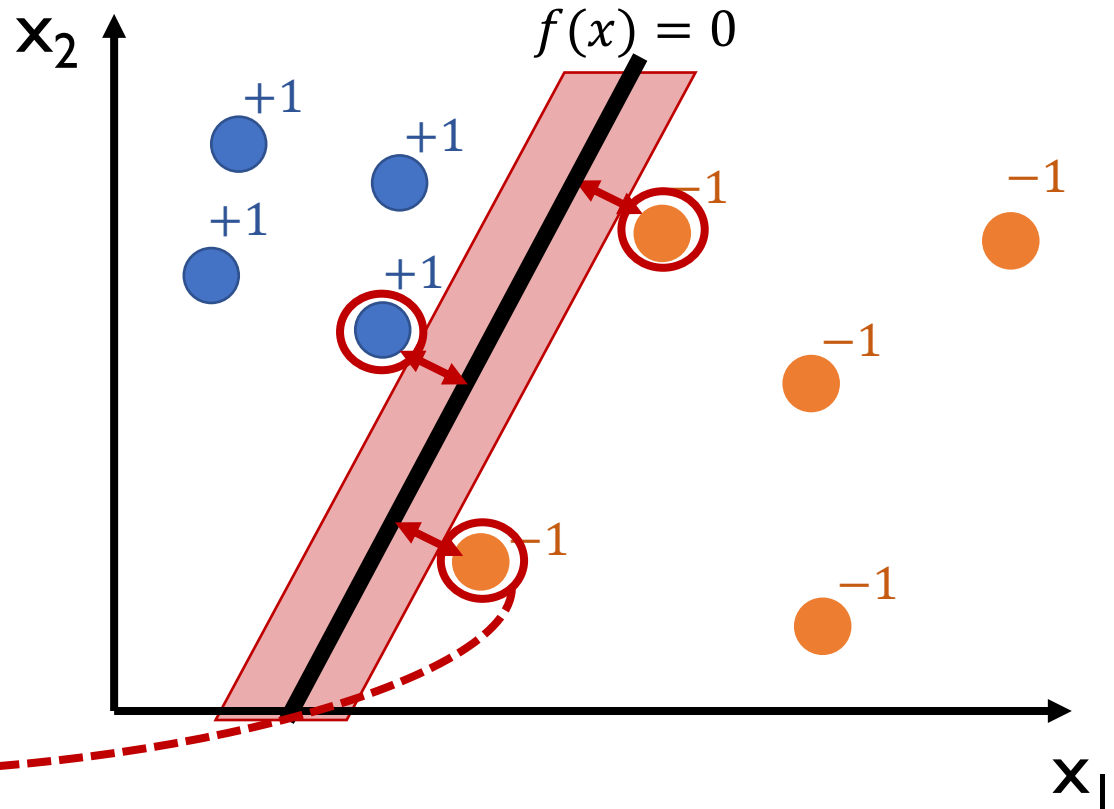


# Maximal Margin Classifier

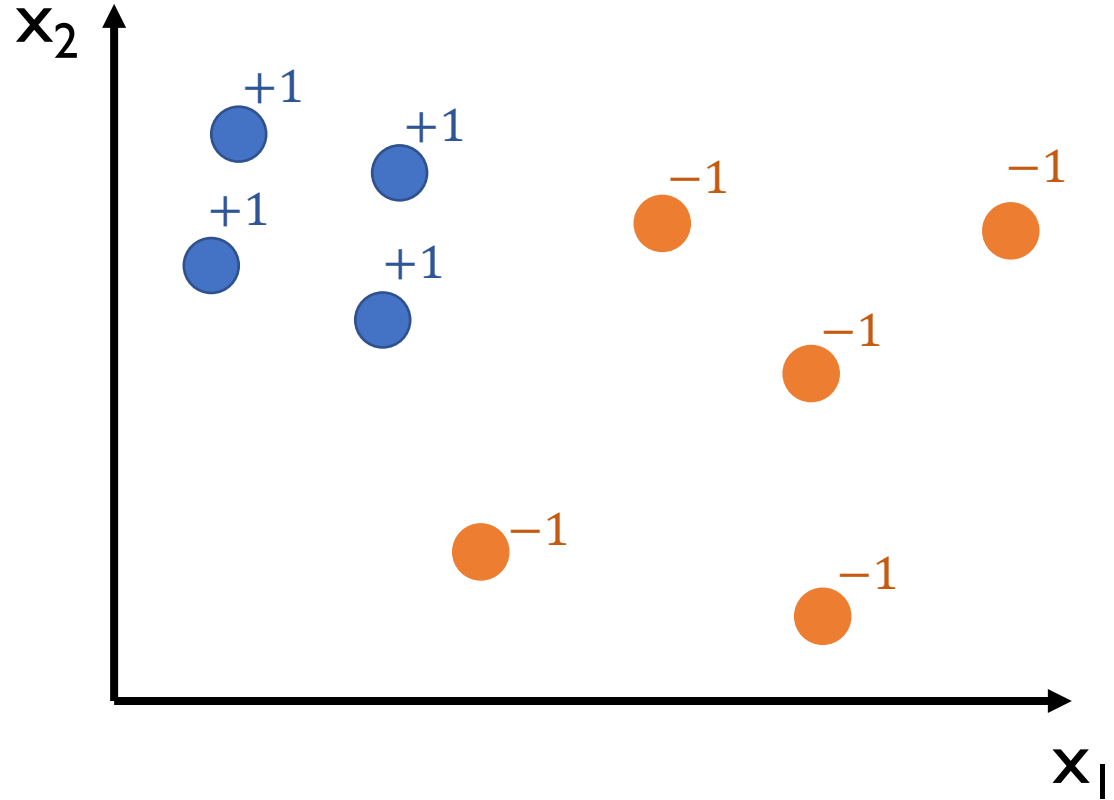
Solving optimization problem we find  $\beta_0, \alpha_1, \dots, \alpha_N$  such that hyperplane:

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i x^T x^{(i)}$$

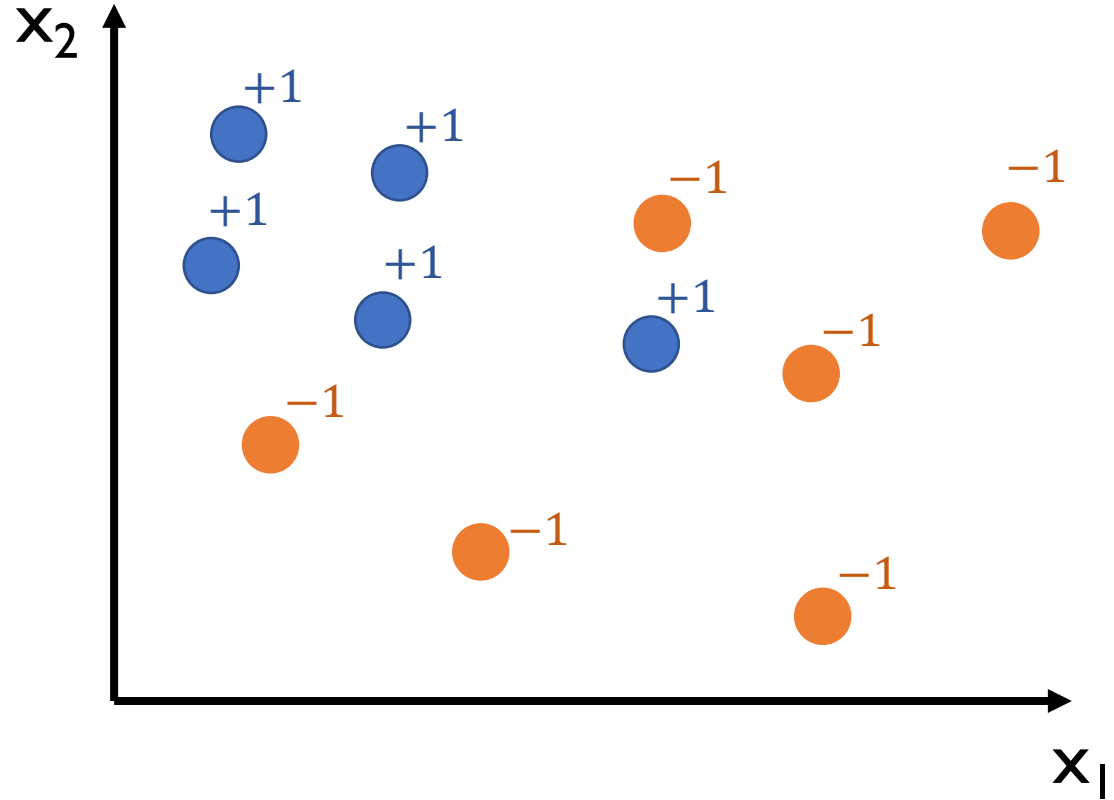
Support Vectors



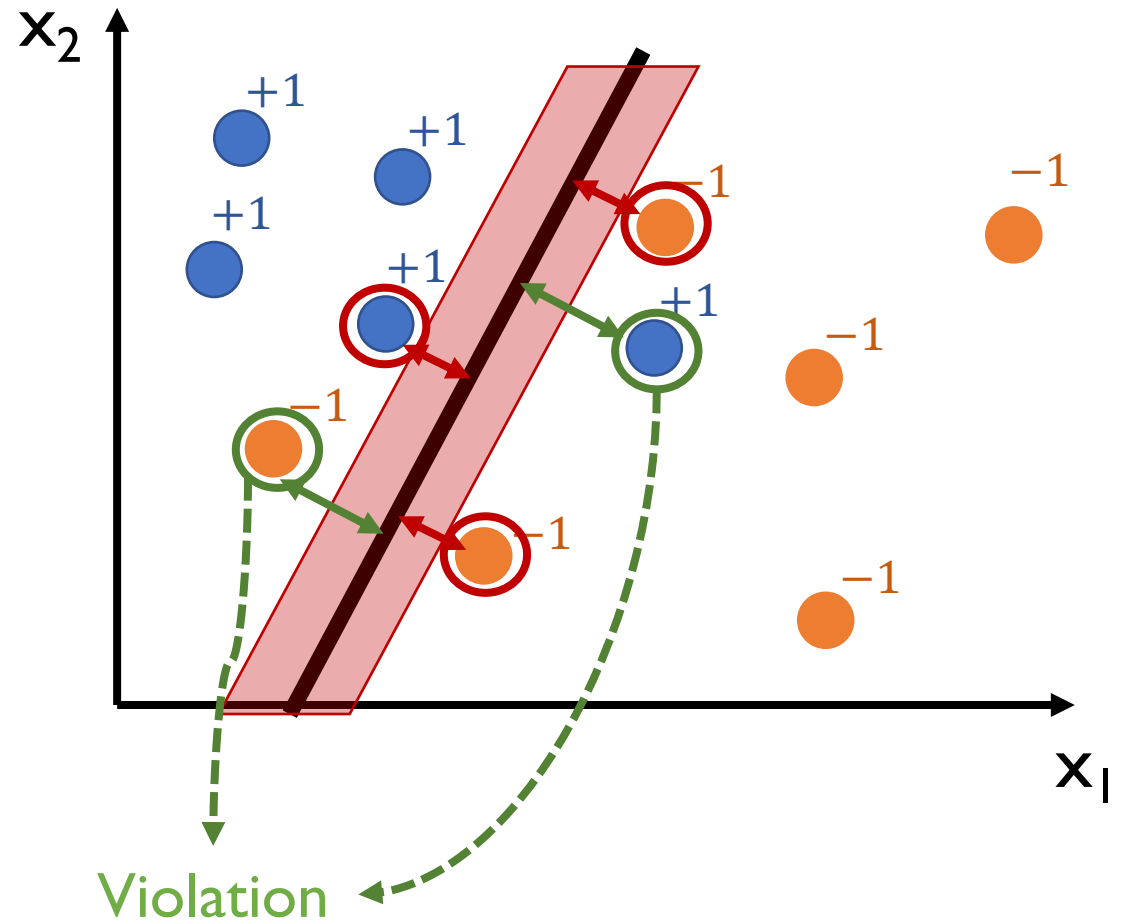
# What if there is no separating hyperplane?



# What if there is no separating hyperplane?



# Support Vector Classifier



# Support Vector Classifier

$$\max_{\beta_0, \beta_1, \beta_2} M$$

$$\text{such that } \beta_0^2 + \beta_1^2 + \beta_2^2 = 1$$

For all training data

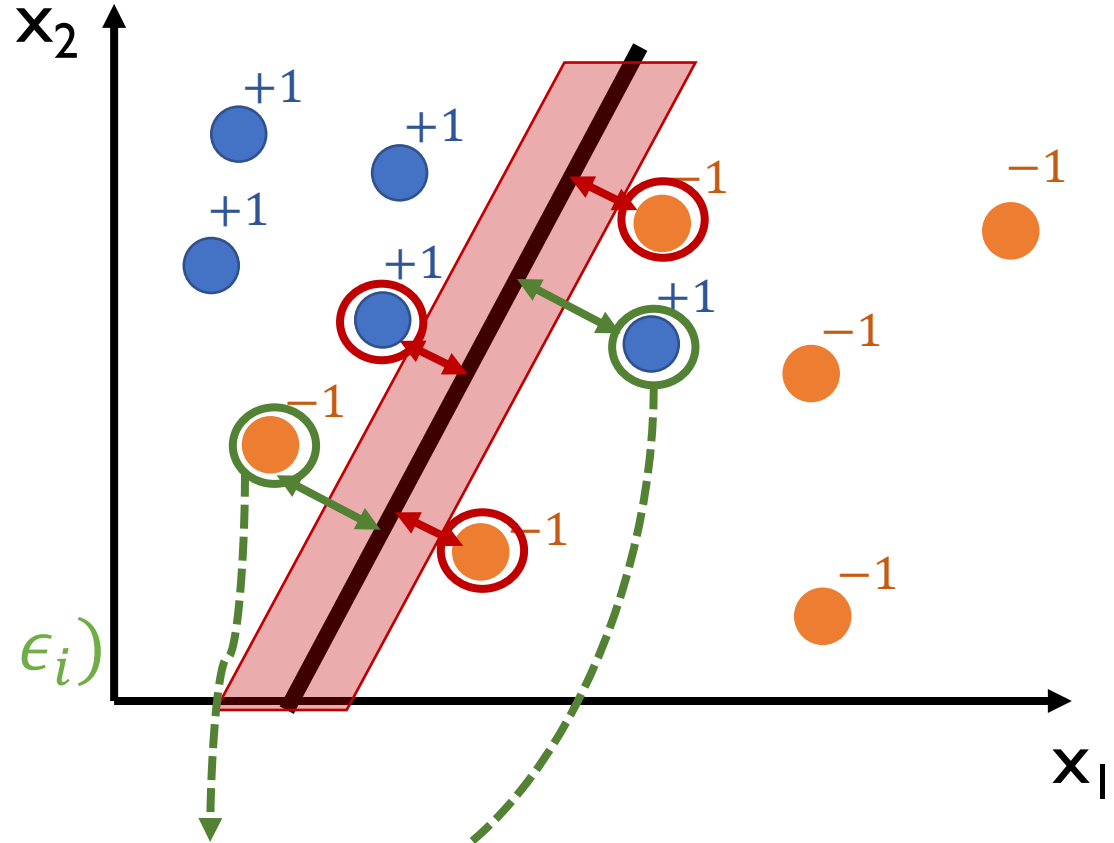
$$y^{(i)} (\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) \geq M(1 - \epsilon_i)$$

such that

$$\sum_{i=1}^N \epsilon_i \leq C,$$

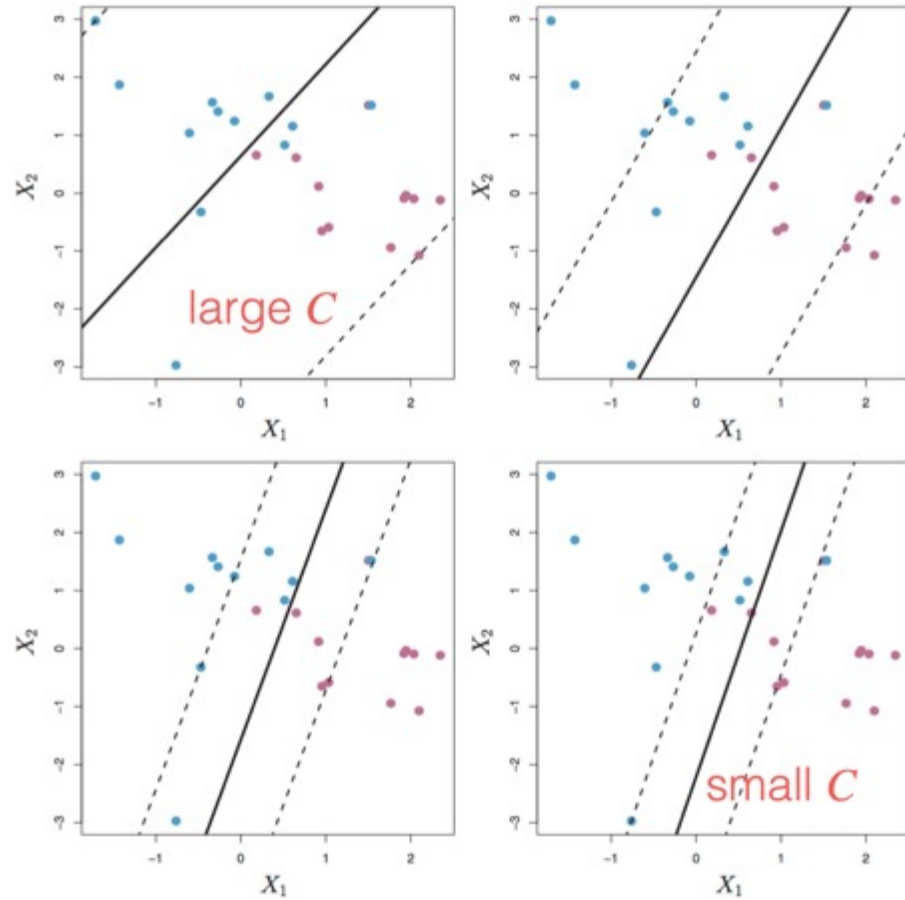
$$\epsilon_i \geq 0$$

Violation



# Support Vector Classifier

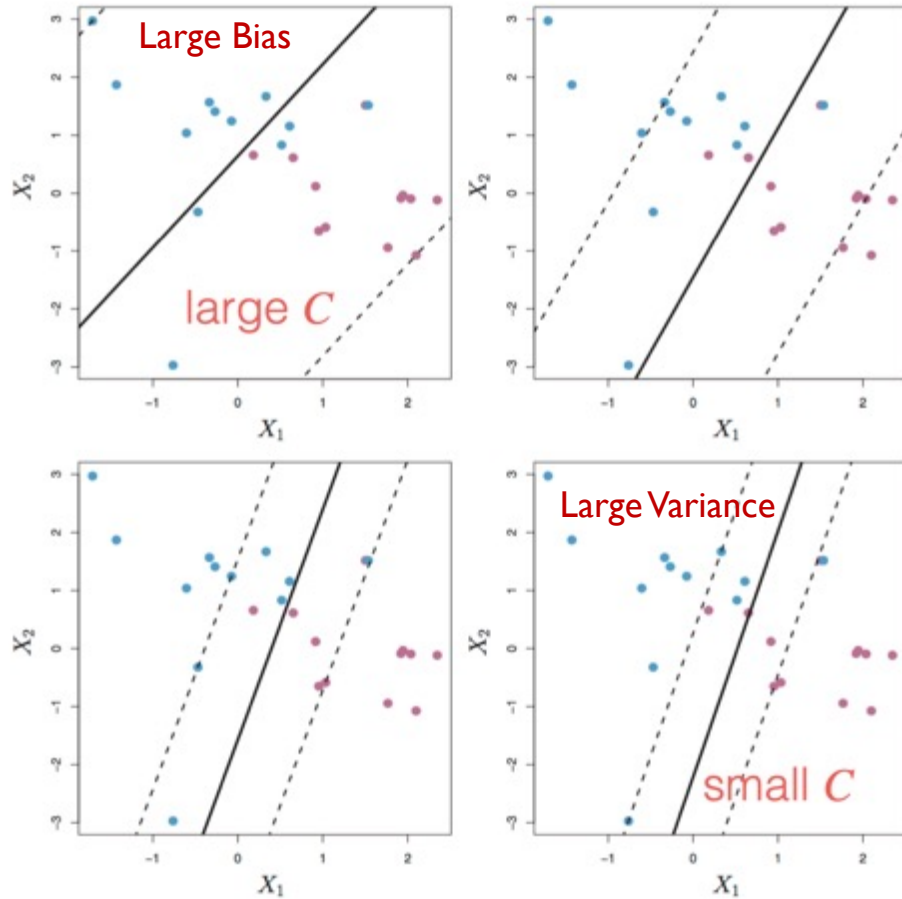
$C$  is the “budget” for violations



ISL Fig 9.7

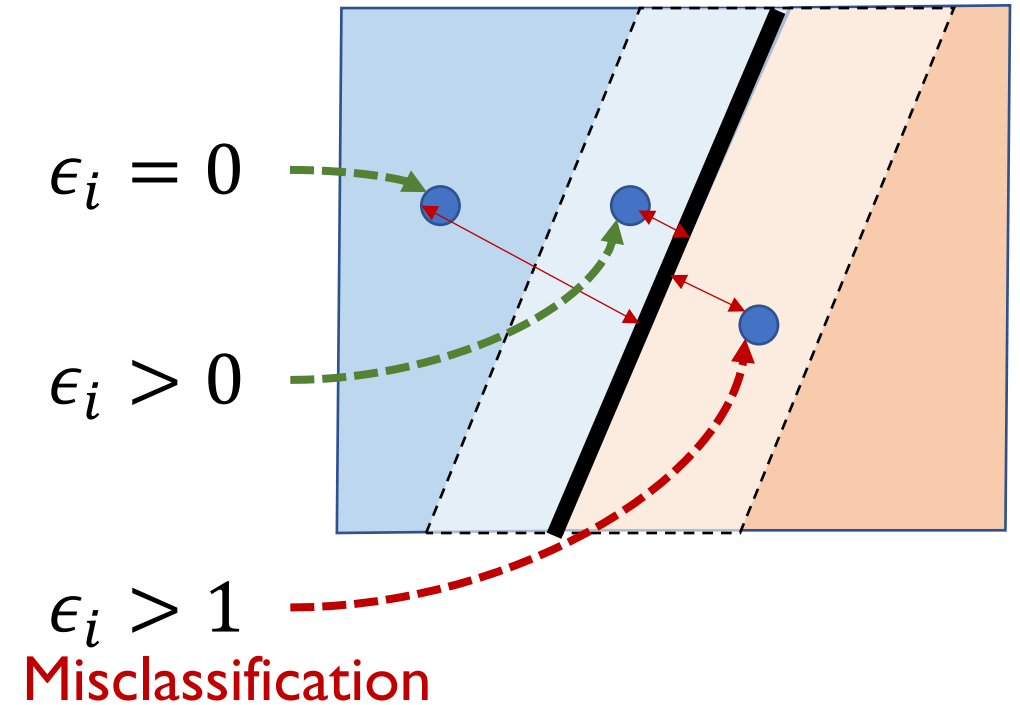
# Support Vector Classifier

$C$  is the “budget” for violations

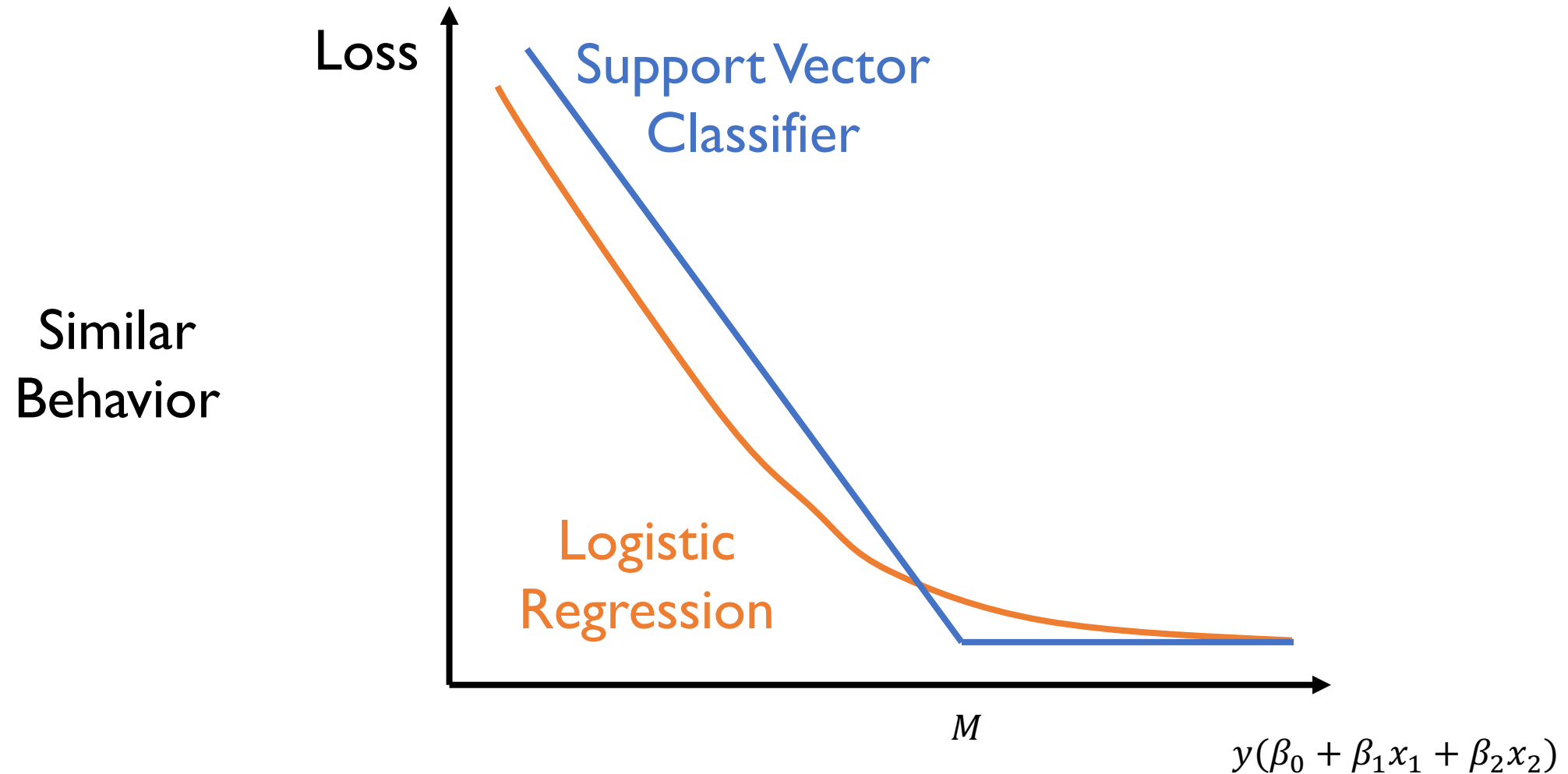


ISL Fig 9.7

Slack variables  $\epsilon_i$



# Support Vector Classifier vs Logistic Regression



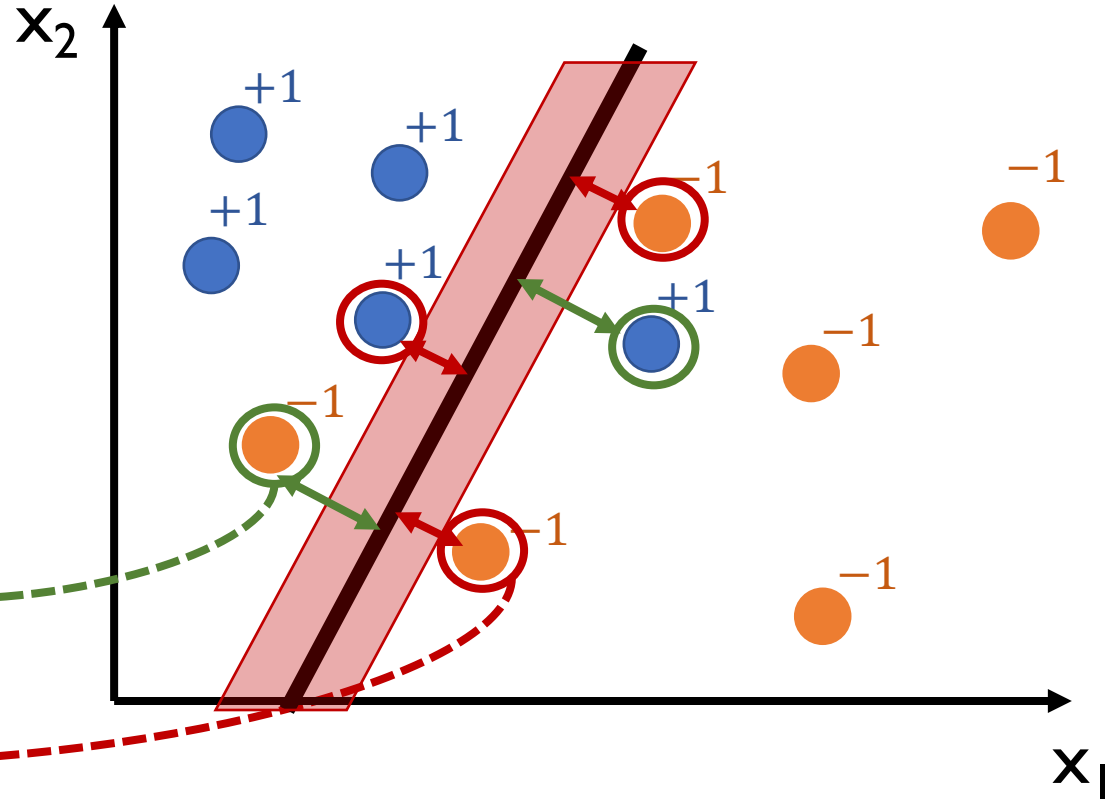


# Support Vector Classifier

Solving optimization problem we find  $\beta_0, \alpha_1, \dots, \alpha_N$  such that hyperplane:

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i x^T x^{(i)}$$

Support Vectors



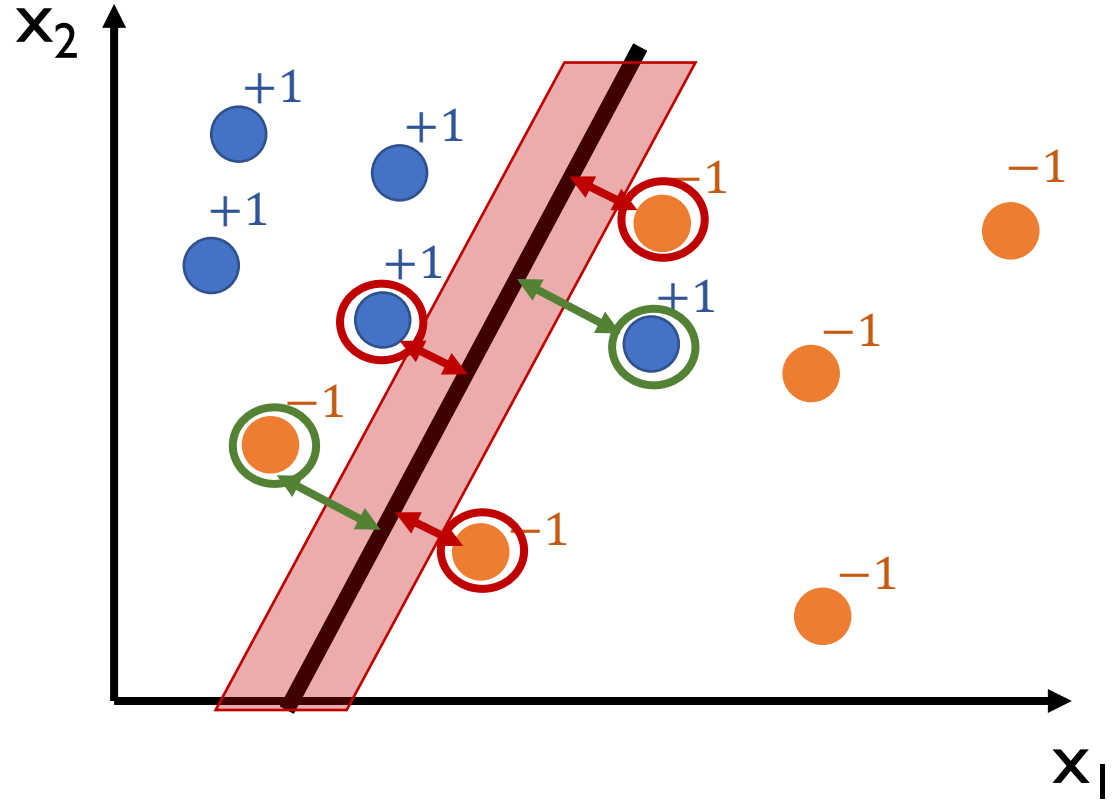
# Support Vector Classifier

Solving optimization problem we find  $\beta_0, \alpha_1, \dots, \alpha_N$  such that hyperplane:

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \boxed{x^T x^{(i)}}$$



**We only get linear boundaries**



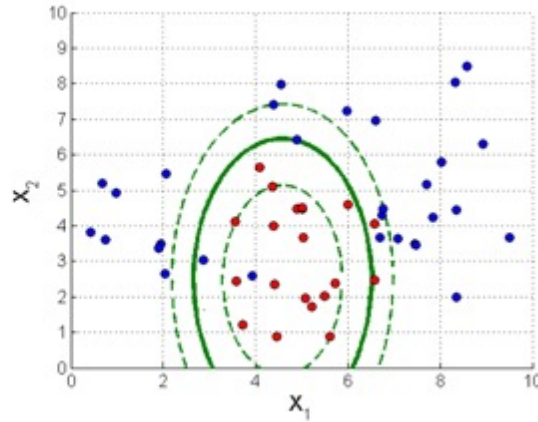
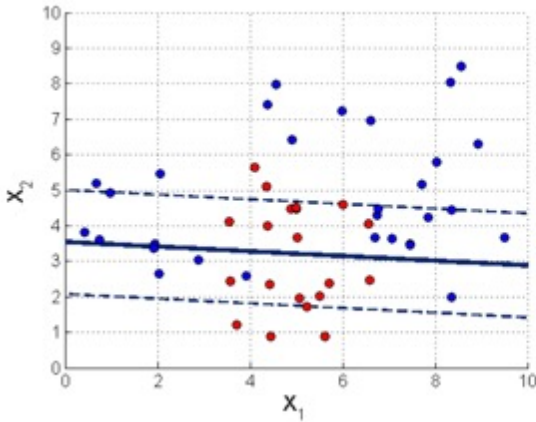
# Beyond linear decision boundaries

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \boxed{x^T x^{(i)}}$$

Option 1)

Add additional features

$$X_1, X_2, X_1 X_2, X_1^2, X_2^2, \dots$$

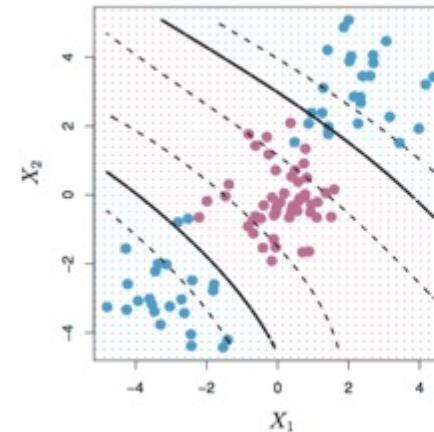


Option 2)

Generalize inner product

$$\text{Kernels } x^T x^{(i)} \rightarrow K(x, x^{(i)})$$

Cubic polynomial kernel



Radial kernel

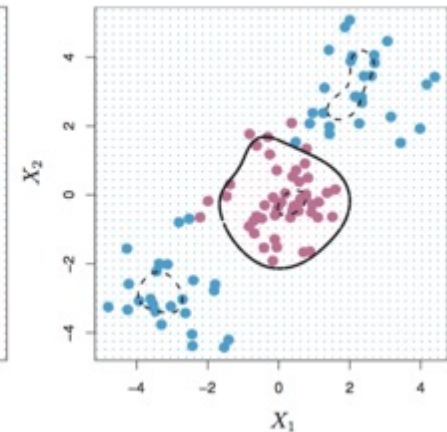


FIGURE 9.9. ISL (8th printing 2017)

# Support Vector Machines

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \boxed{K(x, x^{(i)})}$$

Linear Kernel  $K(x, x^{(i)}) = x^T x^{(i)}$

Polynomial Kernel  $K(x, x^{(i)}) = (1 + x^T x^{(i)})^p \longrightarrow$  Includes  $x_1^k x_2^l$ ,  
 $k + l \leq p$

Radial Basis Kernel  $K(x, x^{(i)}) = \exp(-\gamma \|x - x^{(i)}\|^2) \longrightarrow$  Includes infinite  
# features

Kernel  $\approx$  Similarity Measure

# Challenges Support Vector Machines

Hyperparameters:  
C and Kernel

✓ Flexibility

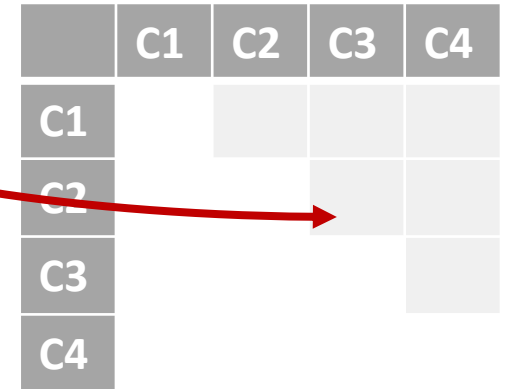
✗ Delicate tuning

✗ Difficult to Interpret

✓ Convex Optimization  
(easy)

Extend to more than 2 classes

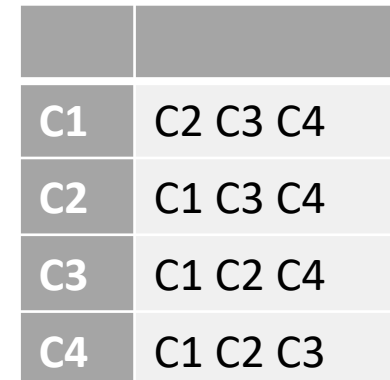
One vs one  
1 SVM per pair



A 4x4 matrix representing pairwise SVMs for 4 classes (C1, C2, C3, C4). The diagonal cells are empty, and the off-diagonal cells are shaded gray. A red arrow points from the text 'One vs one' to the cell at row C2, column C3.

	C1	C2	C3	C4
C1				
C2				
C3				
C4				

One vs all  
1 SVM per class

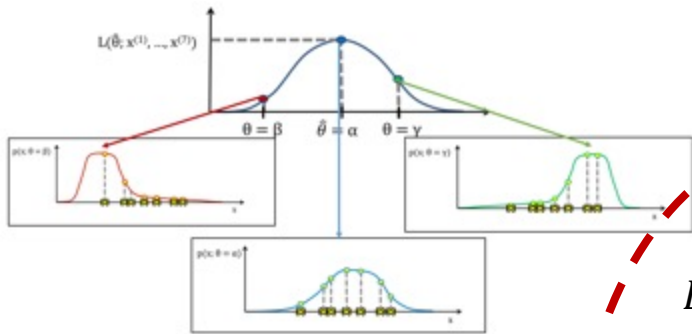


A 4x2 matrix representing One vs All SVMs for 4 classes (C1, C2, C3, C4). The first column contains the class labels, and the second column contains the other three classes.

C1	C2 C3 C4
C2	C1 C3 C4
C3	C1 C2 C4
C4	C1 C2 C3

# How can we extend Linear Regression?

LR is Maximum Likelihood estimator



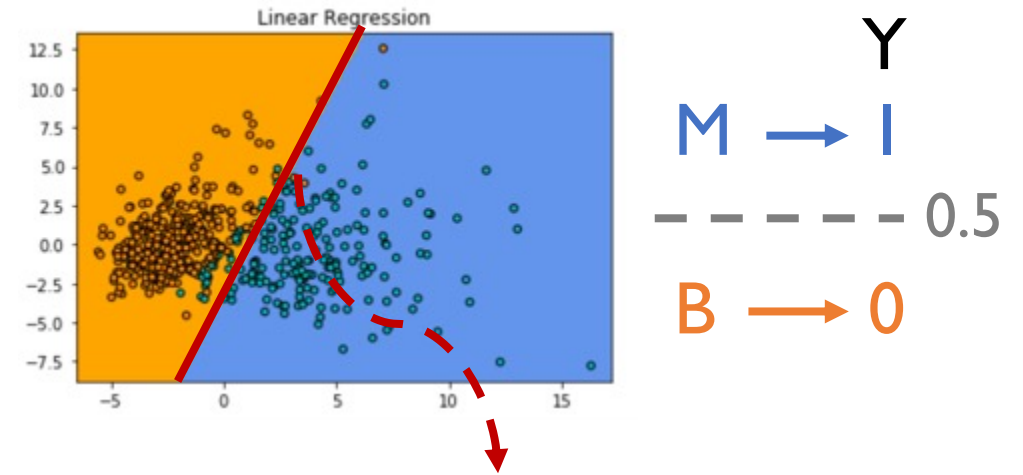
$Y \sim \text{Normal}$

$$E[Y|X] = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

Find a better distribution for Y  
categorical

Logistic Regression

LR creates separating hyperplanes




Optimize the hyperplane

Support Vector Machines

How do we measure the error?

# Most common approach: Error rate

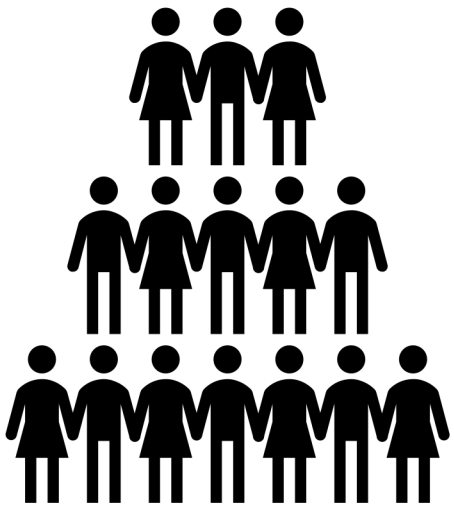
$$error = \frac{1}{N} \sum_{i=1}^N I(\hat{y}^{(i)} \neq y^{(i)})$$

 1 if  $\hat{y}^{(i)} \neq y^{(i)}$   
0 if  $\hat{y}^{(i)} = y^{(i)}$

# Most common approach: Error rate

$$error = \frac{1}{N} \sum_{i=1}^N I(\hat{y}^{(i)} \neq y^{(i)})$$

1 if  $\hat{y}^{(i)} \neq y^{(i)}$   
0 if  $\hat{y}^{(i)} = y^{(i)}$



998 healthy

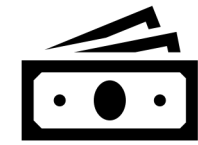
+



2 sick

Design test with  
 $error \leq 0.2\%$

Option 1

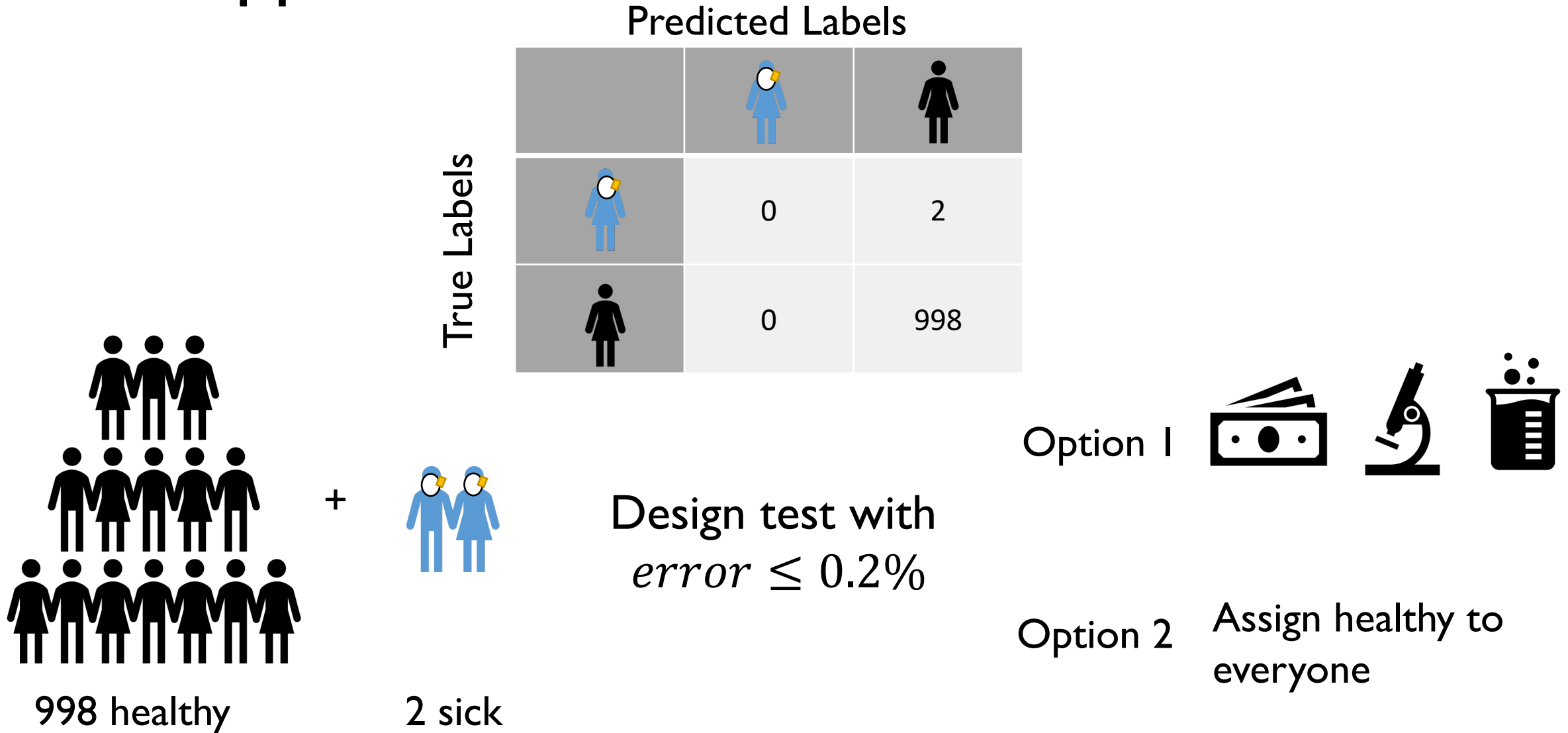


Option 2





Assign healthy to  
everyone



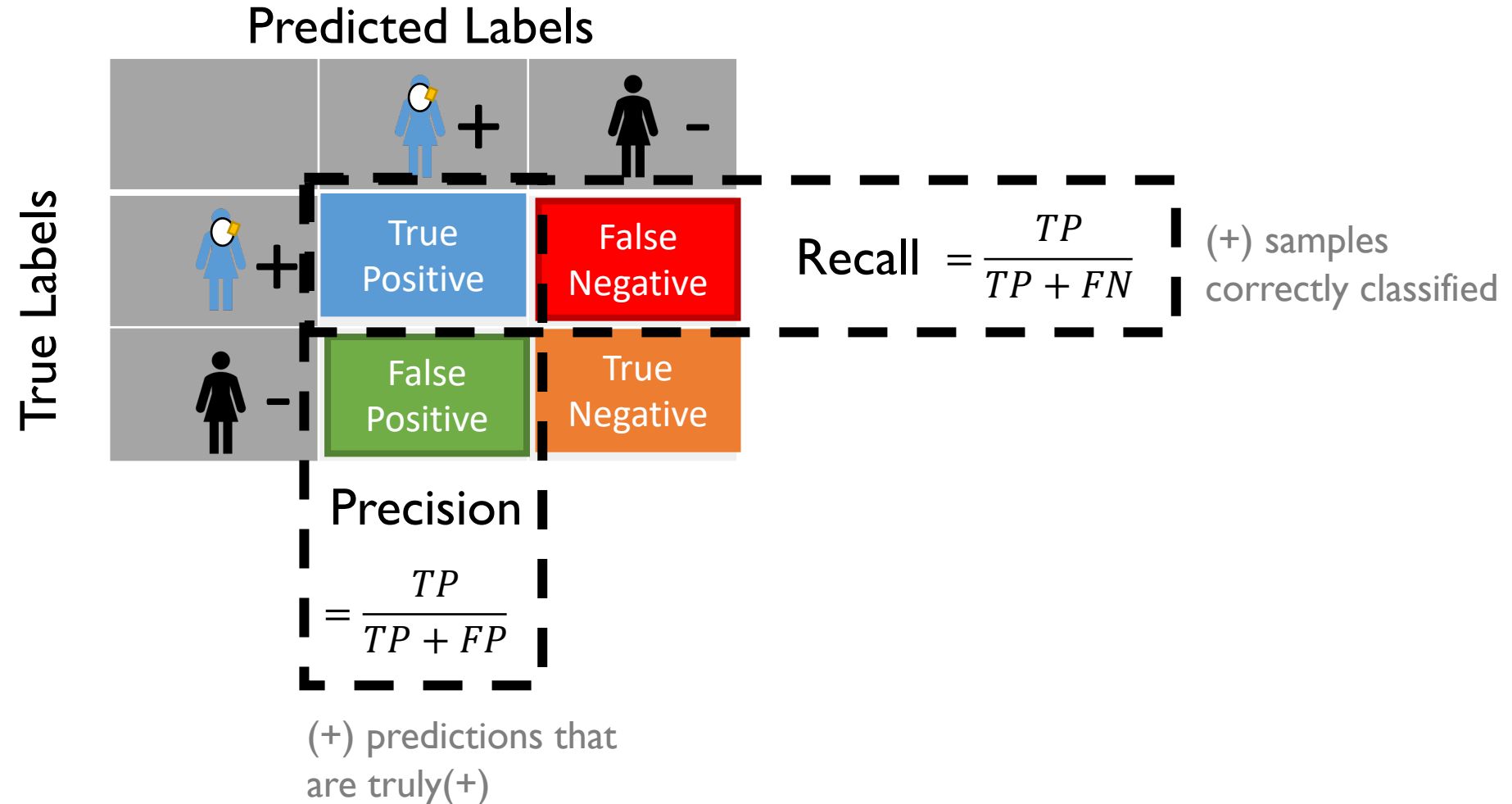
# Useful approach: Confusion Matrix







# Useful approach: Confusion Matrix

		Predicted Labels	
True Labels		 +	 -
	 +	True Positive	False Negative
	 -	False Positive	True Negative





# Useful approach: Precision vs Recall



# Useful approach: Precision vs Recall

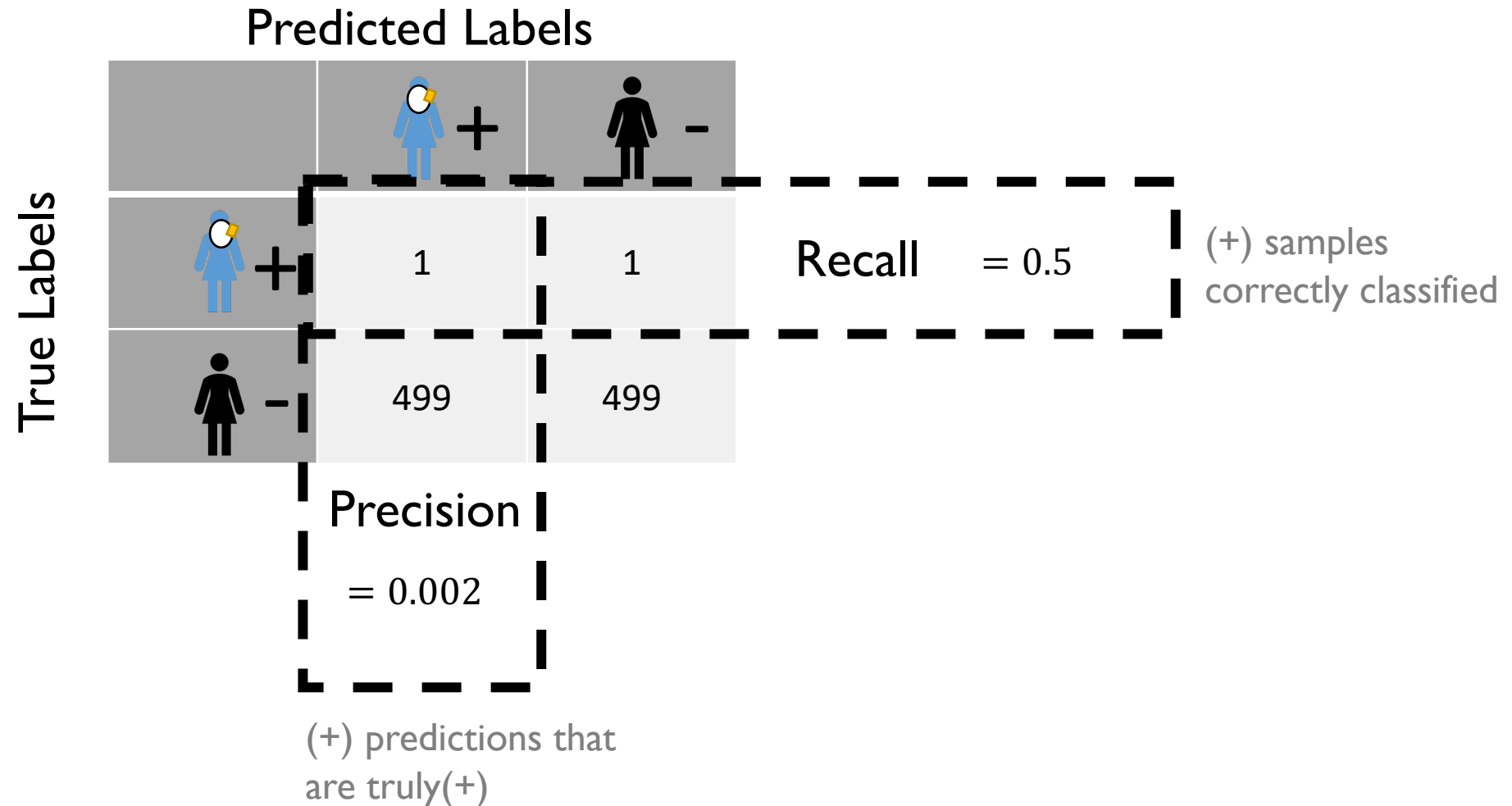
		Predicted Labels				
		 +	 -			
True Labels	 +	0	2	Recall = 0	(+) samples correctly classified	
	 -	0	998			
			Precision = ?			
		(+) predictions that are truly(+)				

# Useful approach: Precision vs Recall

		Predicted Labels				
		 +	 -			
True Labels	 +	2	0	Recall = 1	(+) samples correctly classified	
	 -	0	998			
		Precision = 1				
		(+) predictions that are truly(+)				

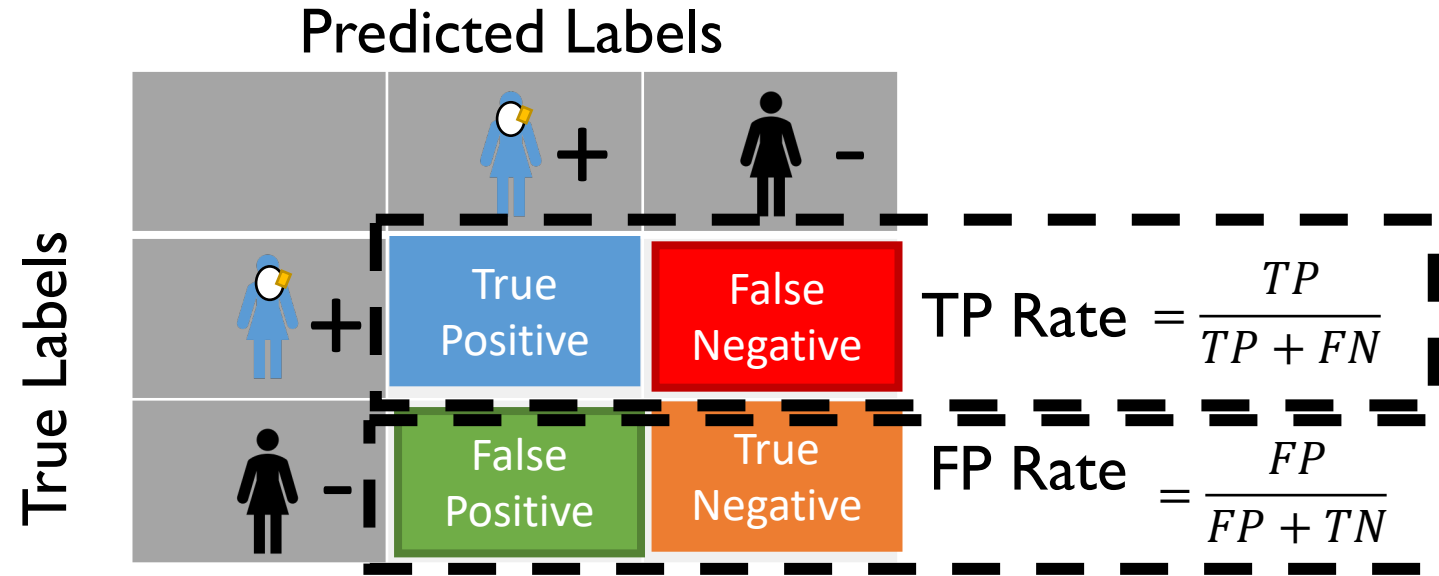
Perfect Classifier

# Useful approach: Precision vs Recall

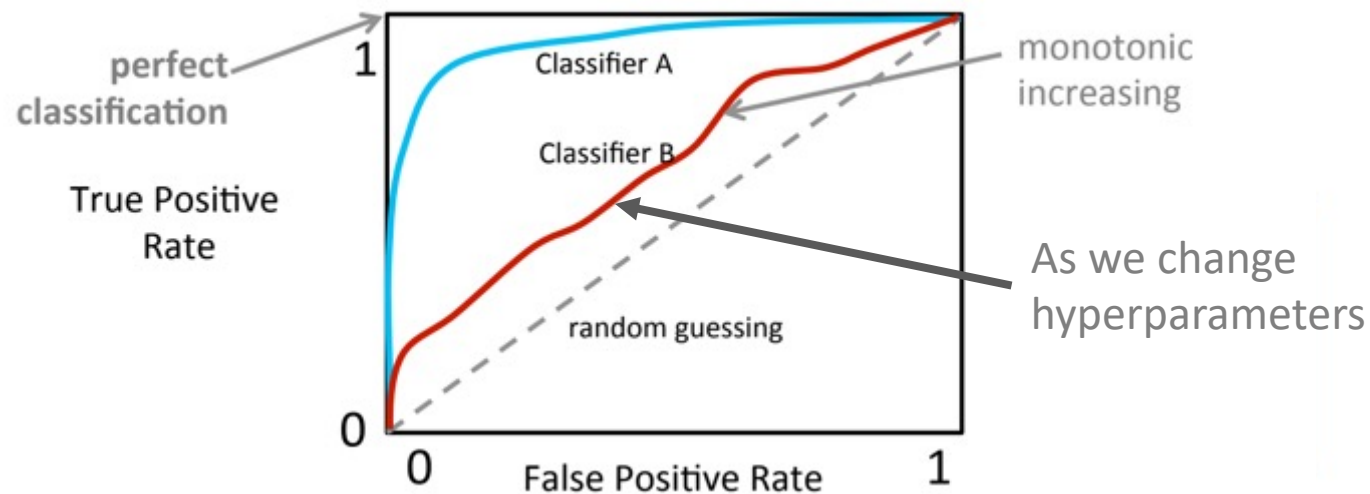
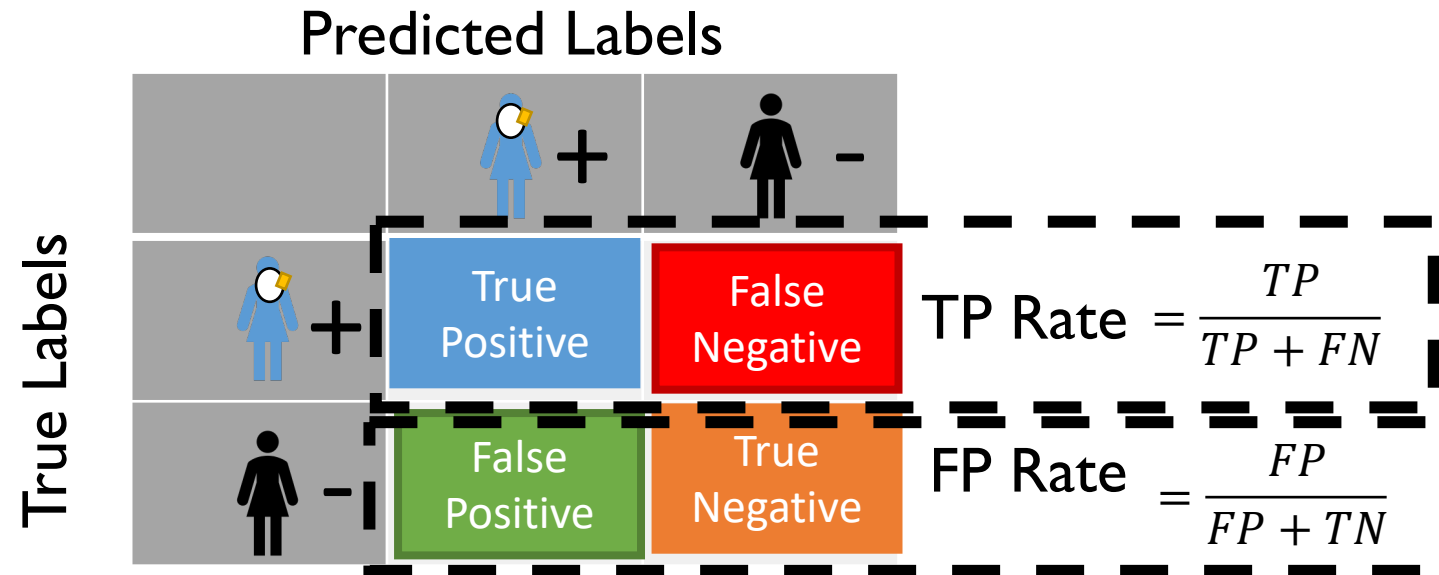


Random Guessing 50/50

# Useful approach: ROC curve



# Useful approach: ROC curve



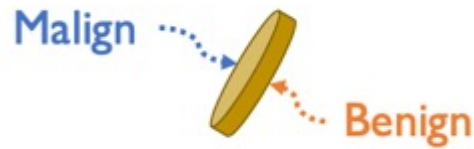
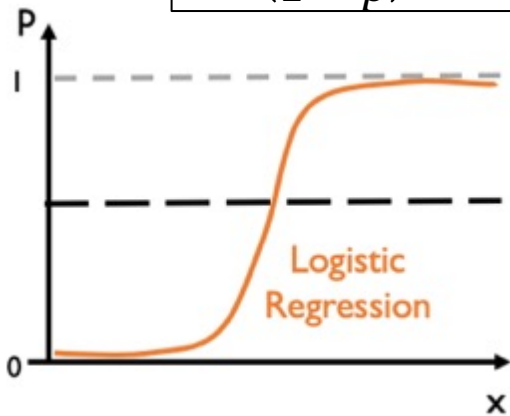


# Today's Recap

## Classification

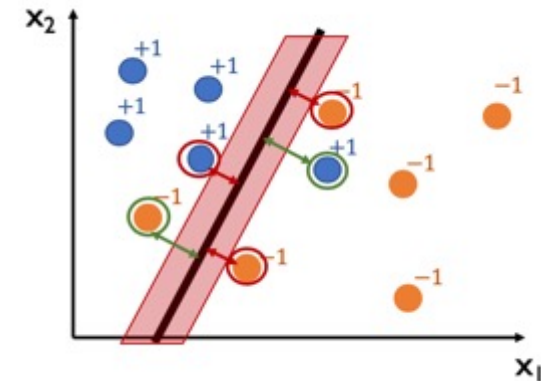
### Logistic Regression

$$\log\left(\frac{p}{1-p}\right) \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$







### Support Vector Machines

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x^{(i)})$$



### Evaluation: Confusion Matrix

		Predicted Labels	
		 +	 -
True Labels	 +	True Positive	False Negative
	 -	False Positive	True Negative