

# CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



# Admin

- **Midterm 2 TODAY!**
- No office hours Friday
- **Final project presentations:**
  - Wednesday Dec 18: 1-4pm (block out the entire time, but we may not need all of it)
  - Option to present last day of class (email me)

# Outline for November 21

- Support Vector Machines Review
- Likelihood functions (Bernoulli and Logistic Regression)
- Finish practice problems (Q2-Q4)

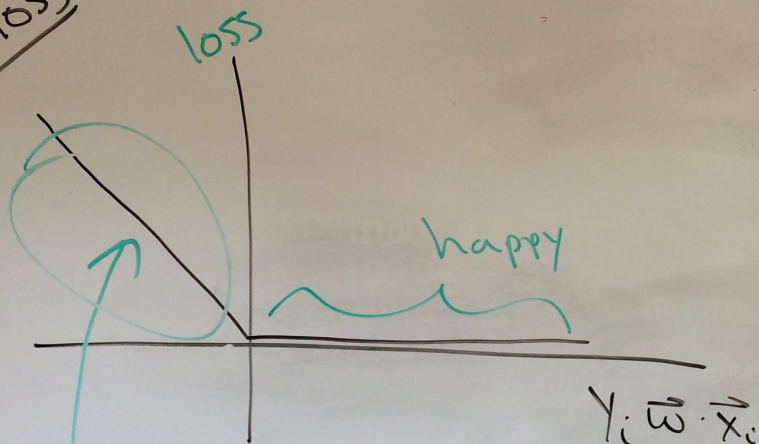
# Outline for November 21

- Support Vector Machines Review
- Likelihood functions (Bernoulli and Logistic Regression)
- Finish practice problems (Q2-Q4)

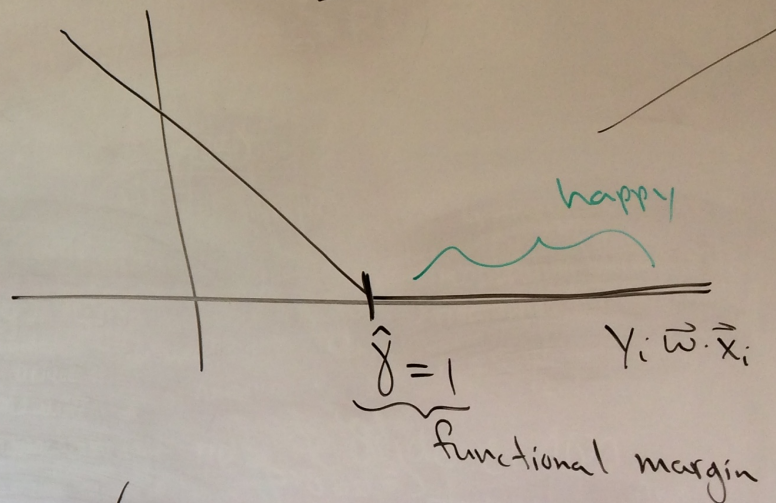


Perceptron & SVM: Same goal: separating hyperplane

hinge loss



SVM



Perceptron:  $J_i(\vec{w}) = \max(0, -y_i \vec{w} \cdot \vec{x}_i)$

if incorrect:  $\vec{w} \leftarrow \vec{w} + y_i \vec{x}_i$  SGD



$$\rightarrow J_i(\vec{w}) = \max(0, 1 - y_i \vec{w} \cdot \vec{x}_i)$$

SVM #1 Goal

maximize geometric margin

$$\max \gamma = \frac{\hat{\gamma}}{\|\vec{w}\|} \leftarrow \min$$

func  $\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2$  ← pts. functional margin.

constraint s.t.  $y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$  ← truth prediction

$(i=1 \dots n)$

Lagrangian

"primal"

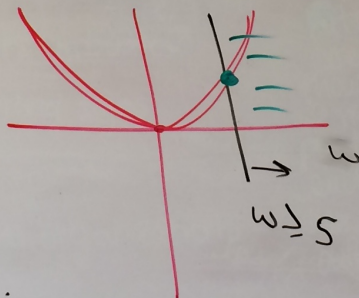
$$L(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1]$$

Lagrange multipliers

$\alpha_i = 0 \Rightarrow$  constraint not used  $\Rightarrow 0$

$\alpha_i > 0 \Rightarrow y_i (\vec{w} \cdot \vec{x}_i + b) - 1 = 0$   
Support vectors

$$L(x, y, \lambda) = \underbrace{f(x, y)}_{\text{min or max}} - \underbrace{\lambda g(x, y)}_{\text{constraint} \Rightarrow g(x, y) = 0}$$



take derivative wrt  $\vec{w}$  &  $b$

"dual"

$$W(\vec{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \vec{x}_i \cdot \vec{x}_j$$

$$\alpha_i \geq 0 \quad \forall i, \quad \sum_{i=1}^n y_i \alpha_i = 0$$

gradient method

$\Rightarrow$  output:  $\vec{q}$

$$\vec{w}^* = \sum_{i=1}^n \alpha_i^* y_i \vec{x}_i$$

weight vector

$$\Rightarrow \text{prediction: } h(\vec{x}) = \text{sign}\left(\sum_{i=1}^n \alpha_i y_i \vec{x}_i \cdot \vec{x} + b^*\right)$$

$$K_{\text{dot}}(\vec{x}, \vec{z}) = \vec{x} \cdot \vec{z}$$

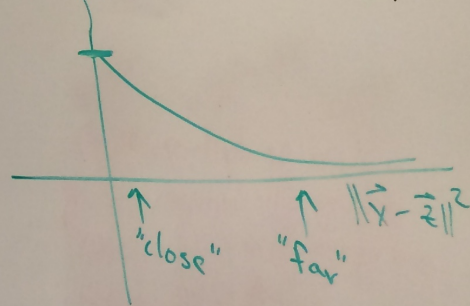
replace with  
any kernel  
I want  
↓  
measure of similarity

$$\vec{x} = \begin{bmatrix} A & G & T & A & C & G \\ A & G & A & A & T & G \end{bmatrix}$$

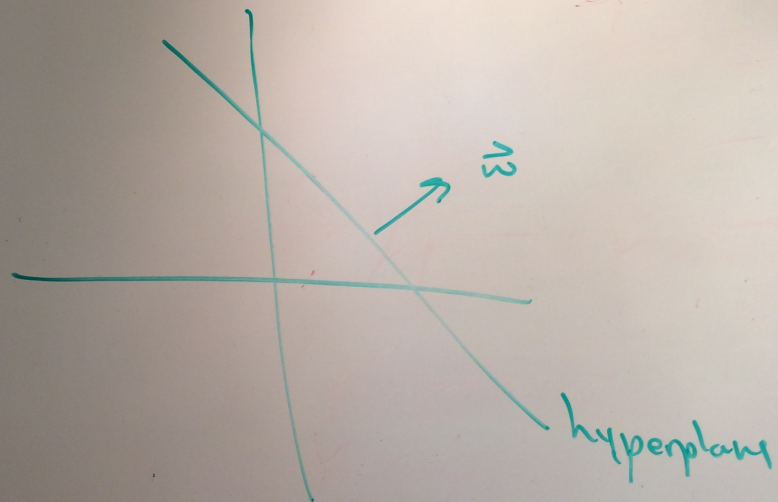
$$K_{\text{dna}}(\vec{x}, \vec{z}) = 4$$

not the margin!

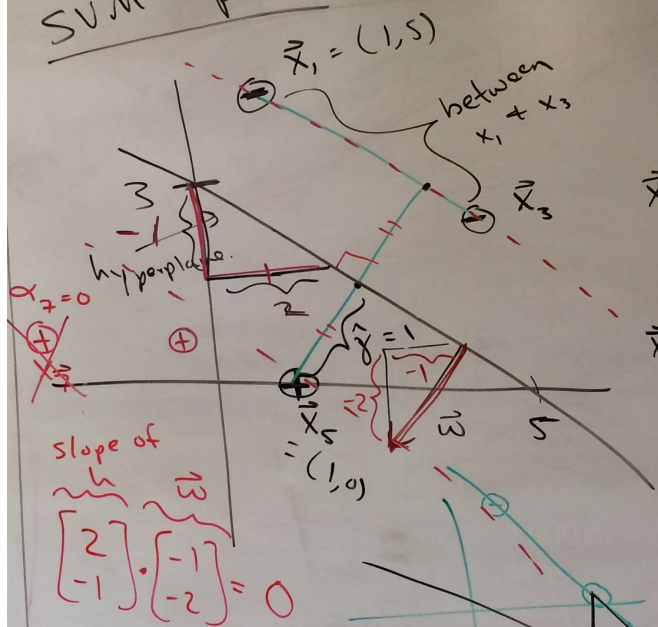
$$K_{\text{rbf}}(\vec{x}, \vec{z}) = \exp\left(-\gamma \|\vec{x} - \vec{z}\|^2\right)$$







SVM pset



$$\vec{w} = a \begin{bmatrix} -1 \\ -2 \end{bmatrix}$$

$$\vec{x}_5: \underbrace{1}_{y_5} \underbrace{\begin{bmatrix} -1 \\ -2 \end{bmatrix}}_{\vec{w}} \cdot \underbrace{\begin{bmatrix} 1 \\ 0 \end{bmatrix}}_{\vec{x}_5} + \underbrace{b}_{b} = 1$$

$$\vec{x}_1: -1 \underbrace{\begin{bmatrix} -1 \\ -2 \end{bmatrix}}_{\vec{w}} \cdot \underbrace{\begin{bmatrix} 1 \\ 5 \end{bmatrix}}_{\vec{x}_1} + \underbrace{b}_{b} = 1$$

$$a = \frac{1}{5}$$

$$b = \frac{6}{5}$$

$$\vec{w} = \begin{bmatrix} -1/5 \\ -2/5 \end{bmatrix}$$

# Perceptron Recap

$$\alpha = 0.2$$

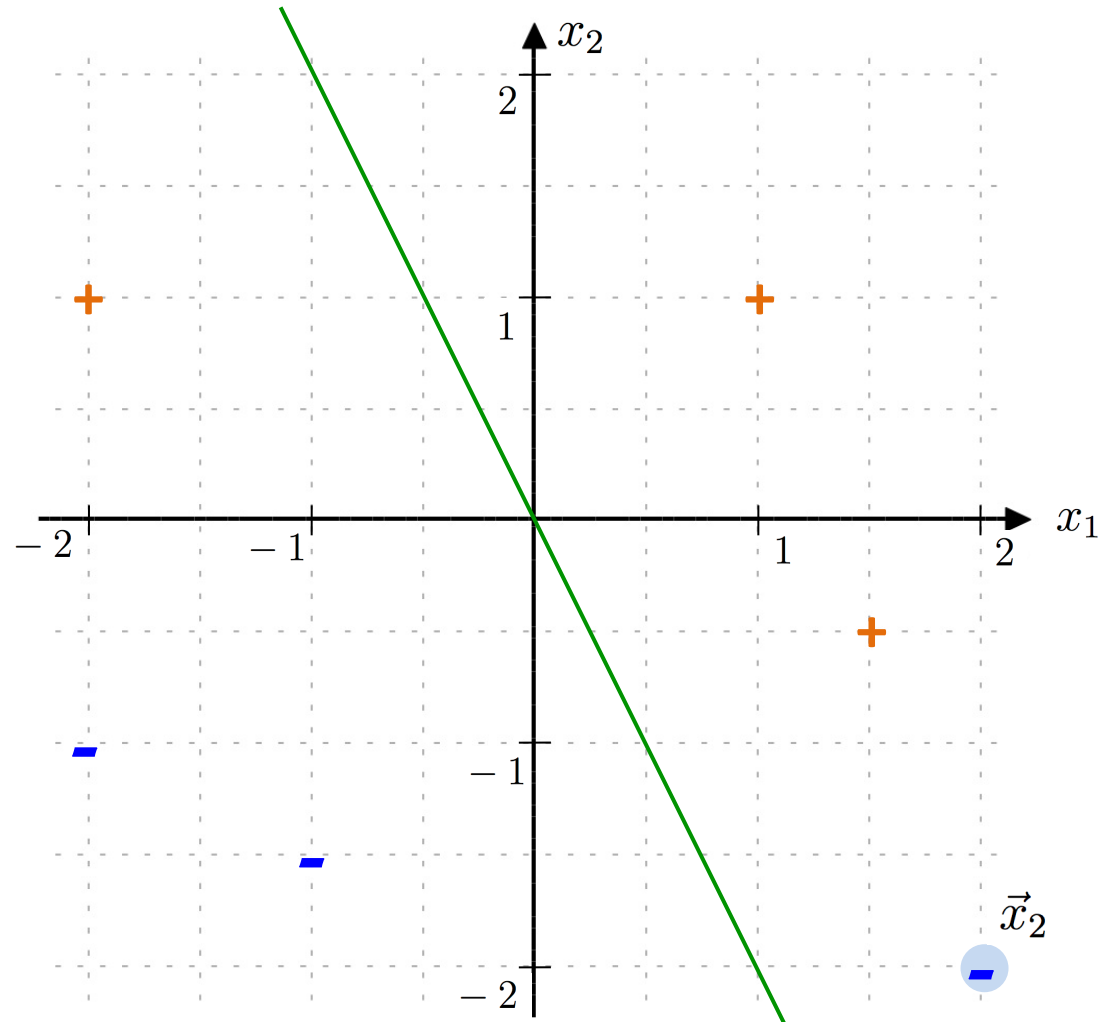
$$\vec{w} = \begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix}$$

Round 2:

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix}$$

$$\vec{w} \cdot \vec{x}_2 > 0$$

Incorrect classification



# Perceptron Recap

$$\alpha = 0.2$$

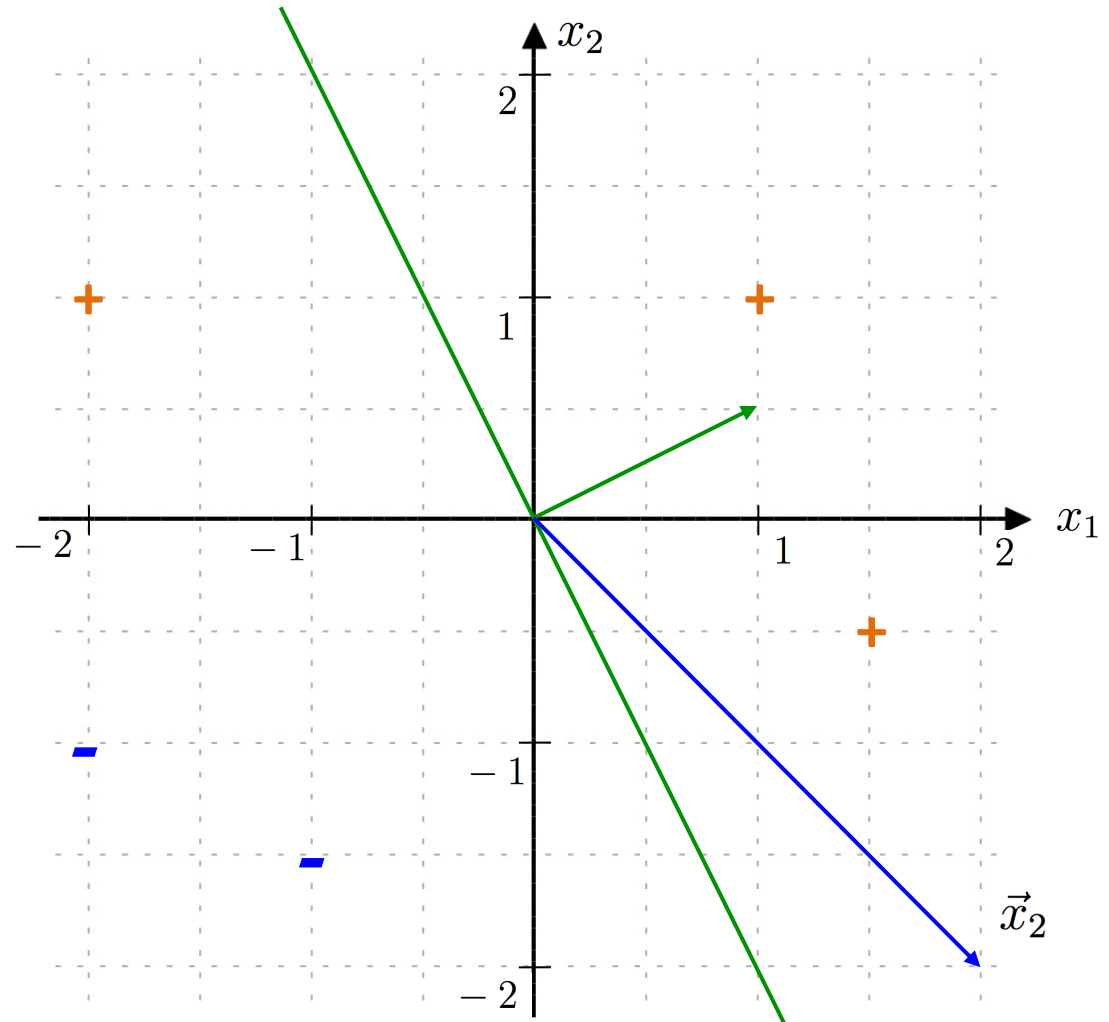
$$\vec{w} = \begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix}$$

Round 2:

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix}$$

$$\vec{w} \cdot \vec{x}_2 > 0$$

Incorrect classification





# Perceptron Recap

$$\alpha = 0.2$$

$$\vec{w} = \begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix}$$

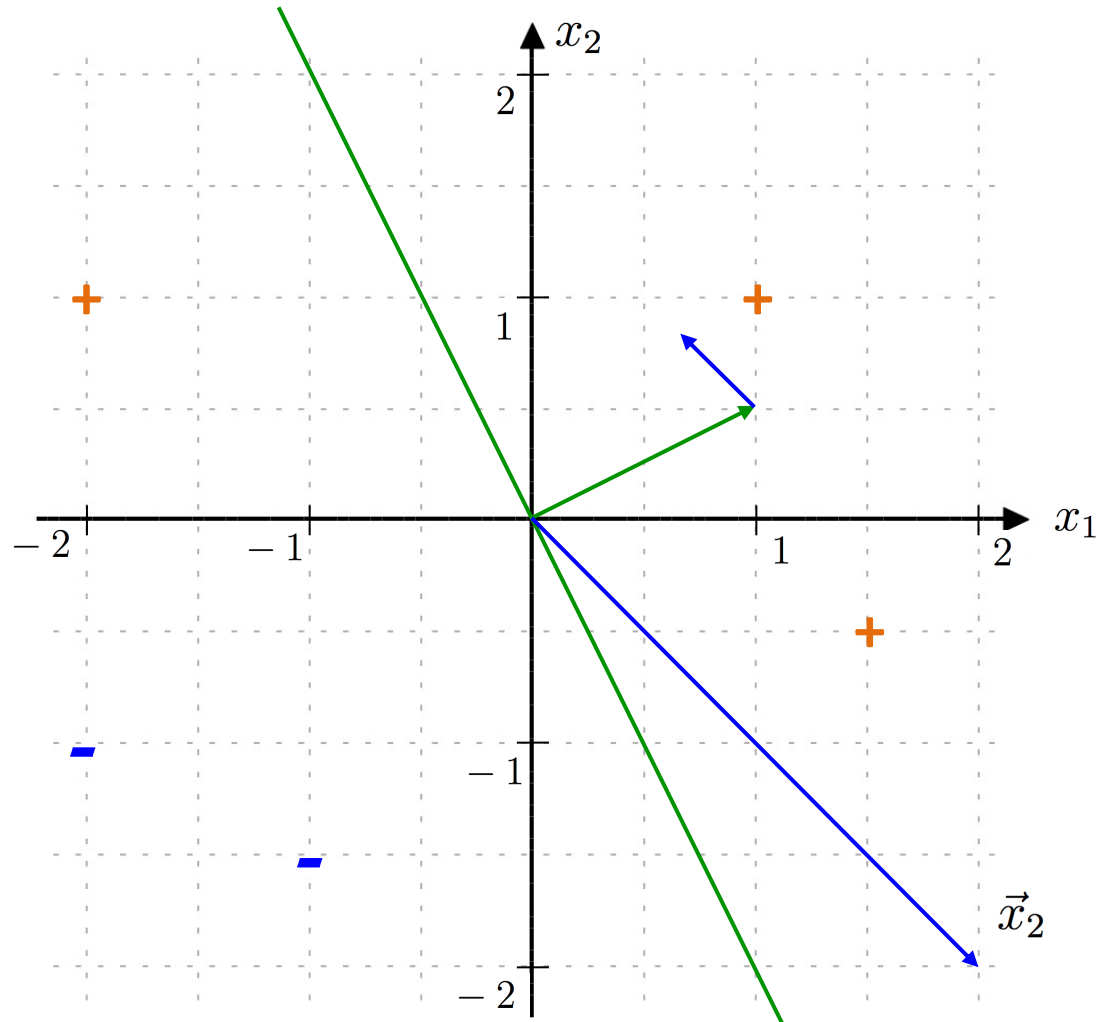
Round 2:

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix}$$

$$\vec{w} \cdot \vec{x}_2 > 0$$

Incorrect classification

“Push”  $\vec{w}$  away from negative point





# Perceptron Recap

$$\alpha = 0.2$$

$$\vec{w} = \begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix}$$

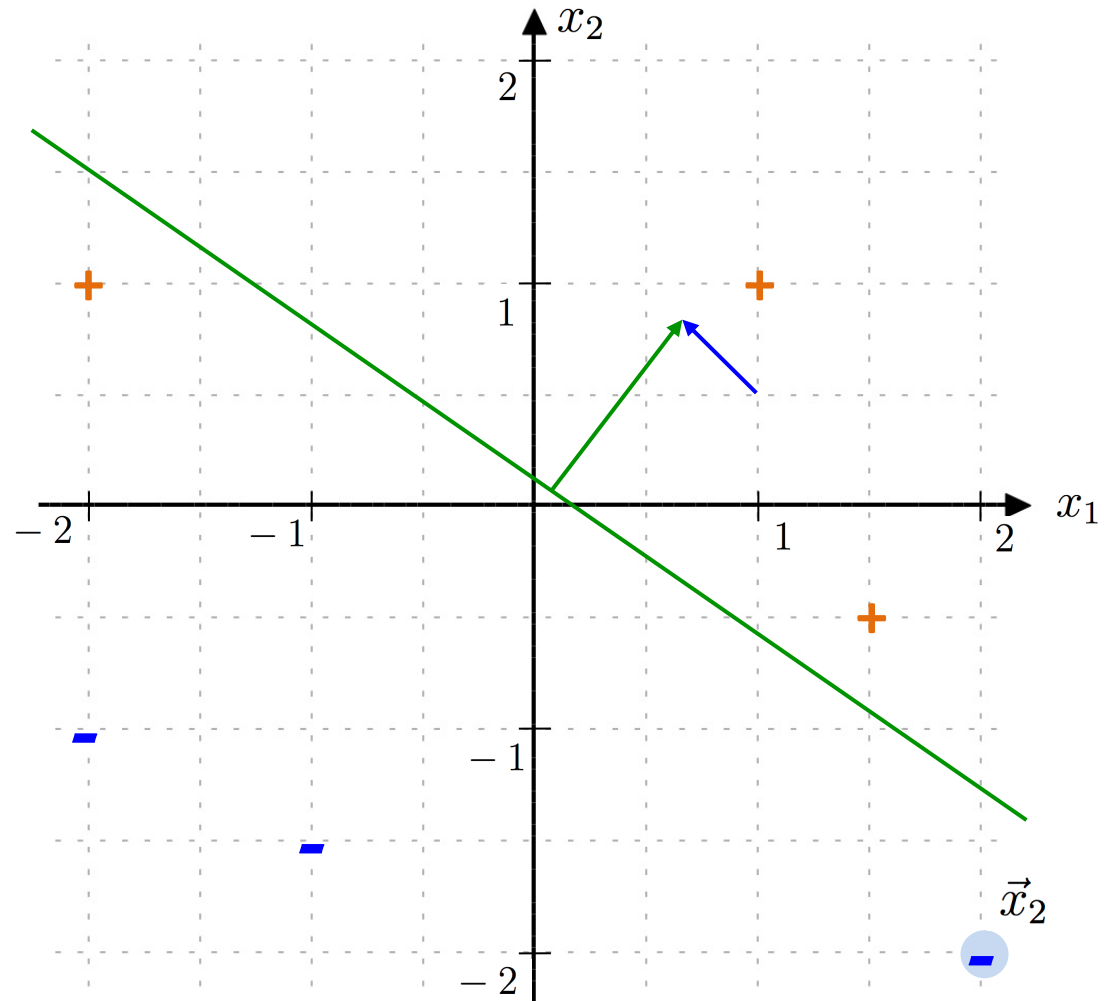
Round 2:

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix}$$

$$\vec{w} \cdot \vec{x}_2 > 0$$

Incorrect classification

“Push”  $\vec{w}$  away from negative point



# Perceptron Recap

$$\alpha = 0.2$$

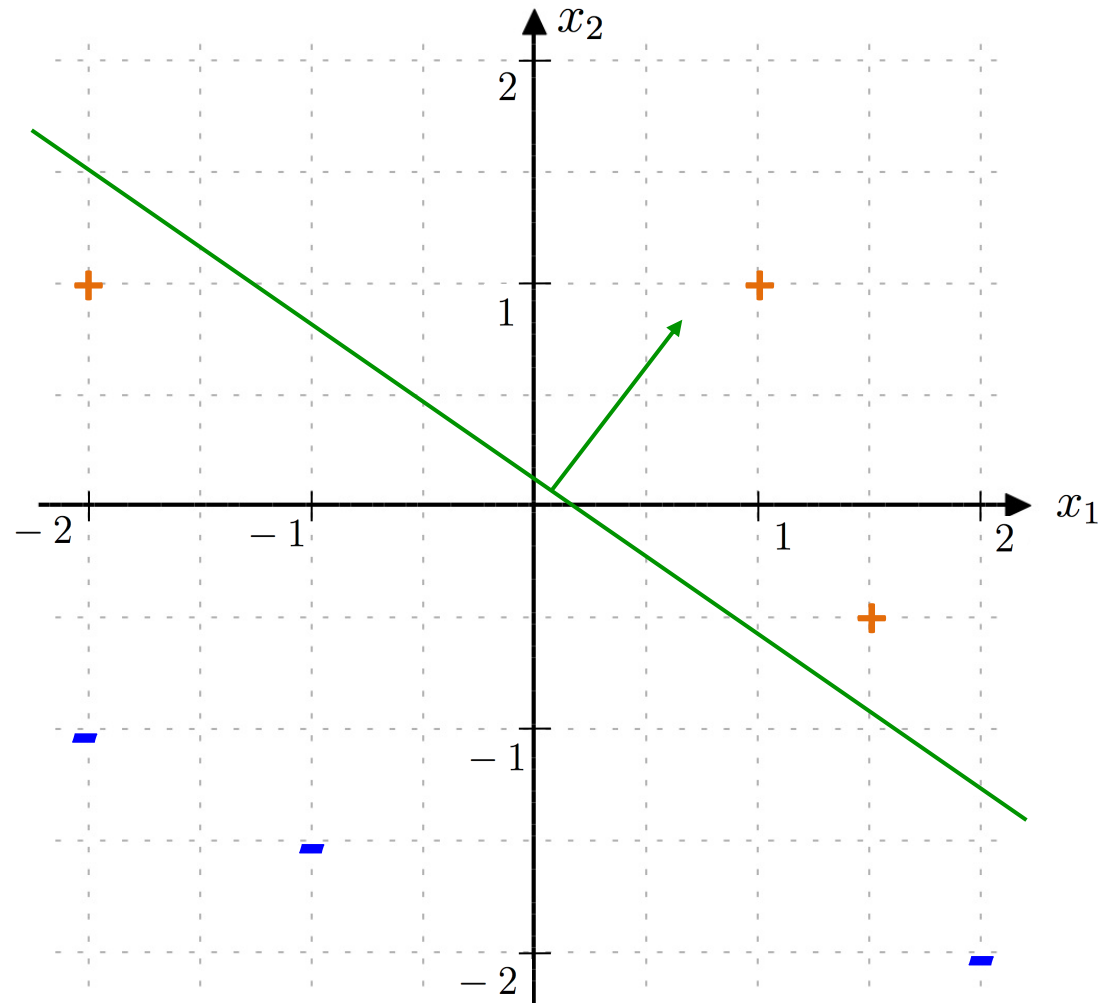
$$\vec{w} = \begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix}$$

Round 2:

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix}$$

$$\vec{w} \cdot \vec{x}_2 > 0$$

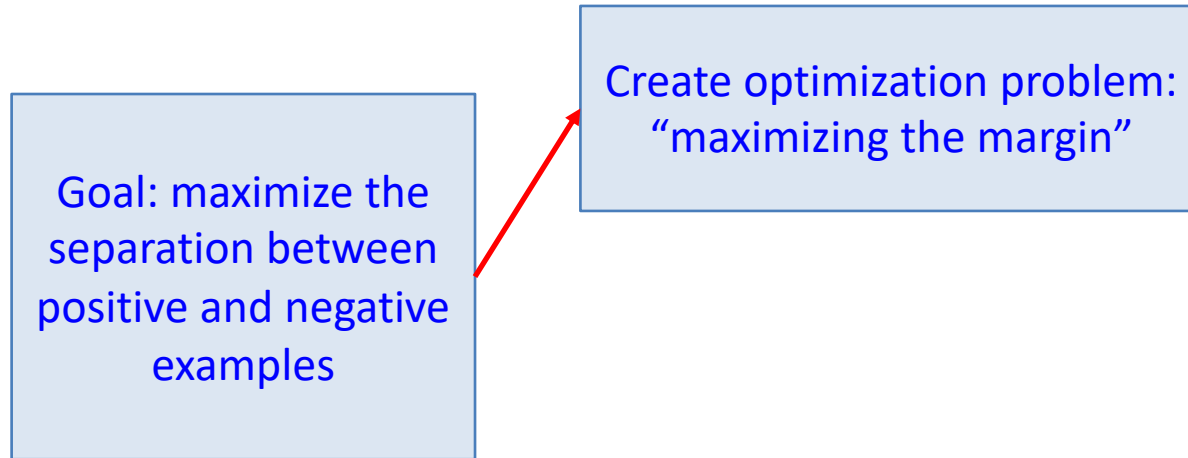
What is the new weight vector?



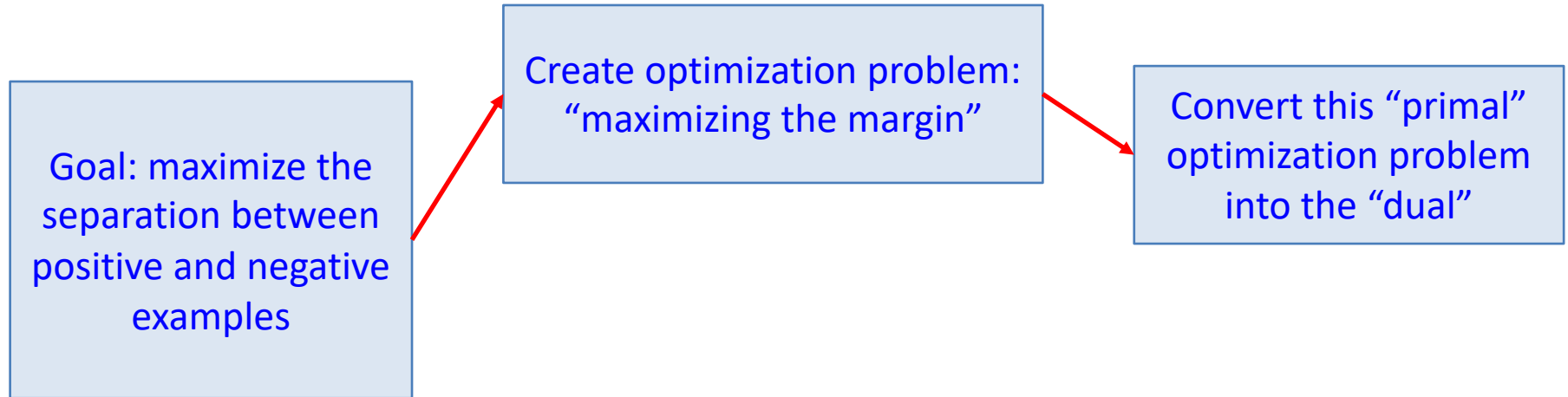
# SVM flowchart

Goal: maximize the  
separation between  
positive and negative  
examples

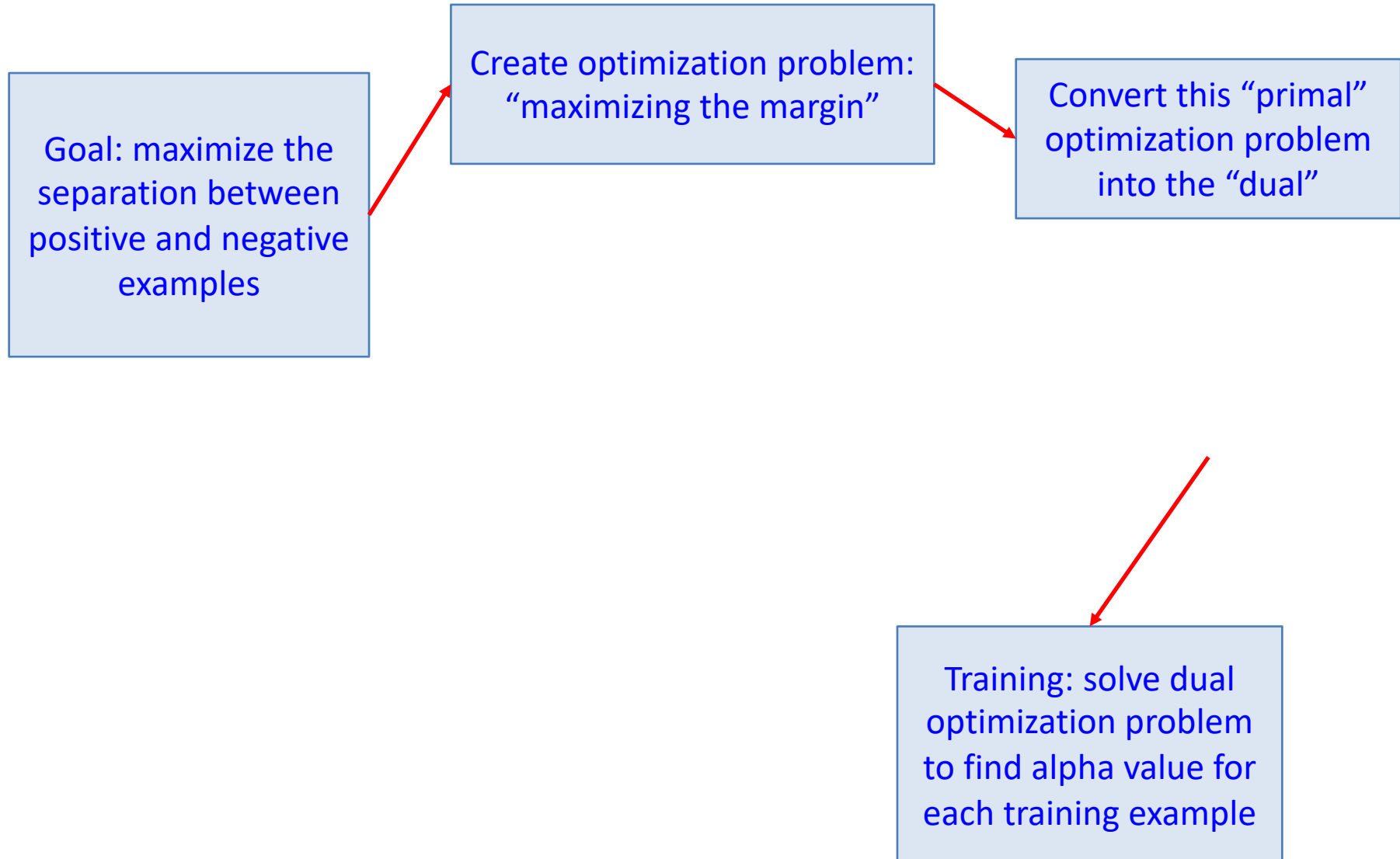
# SVM flowchart



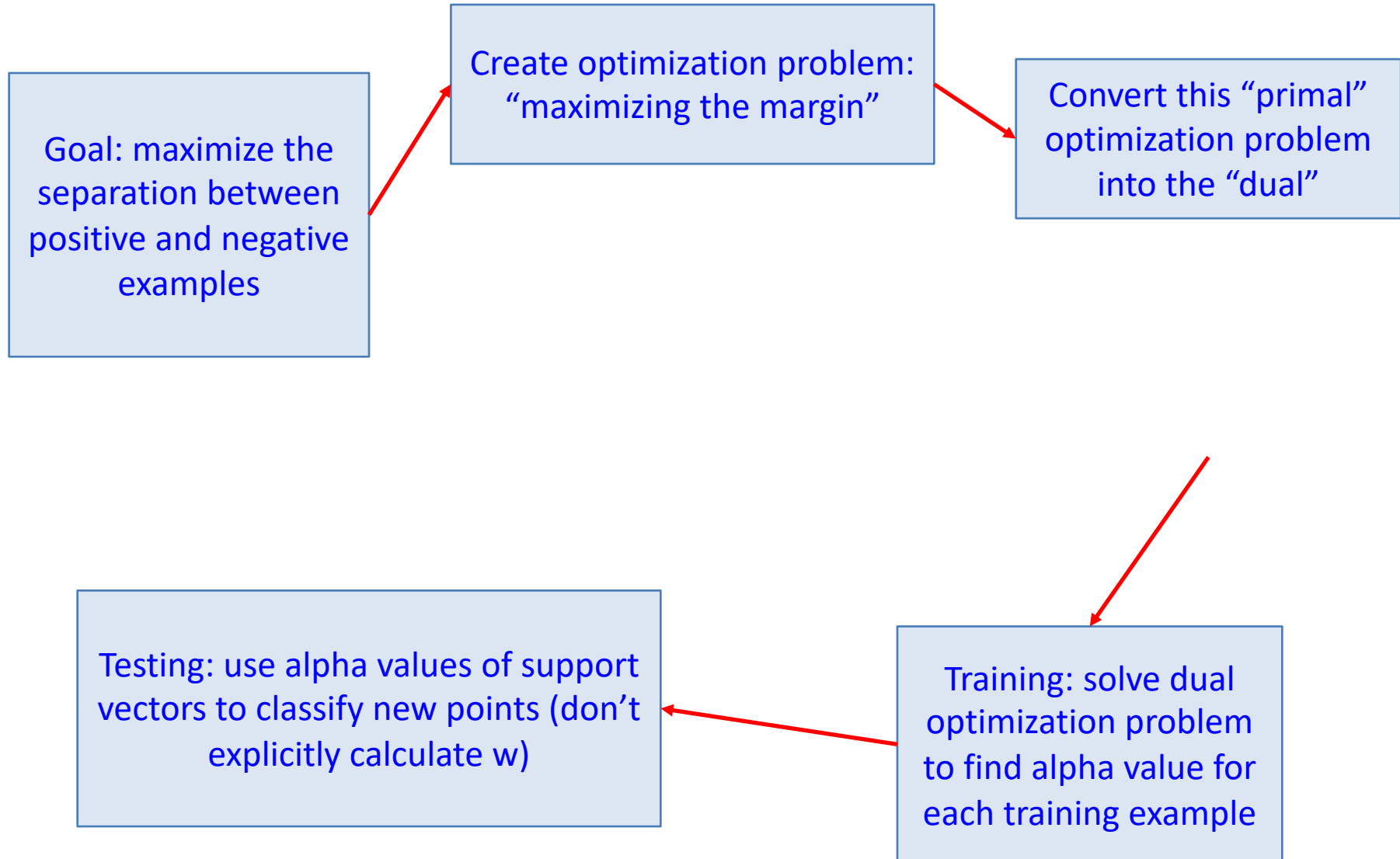
# SVM flowchart



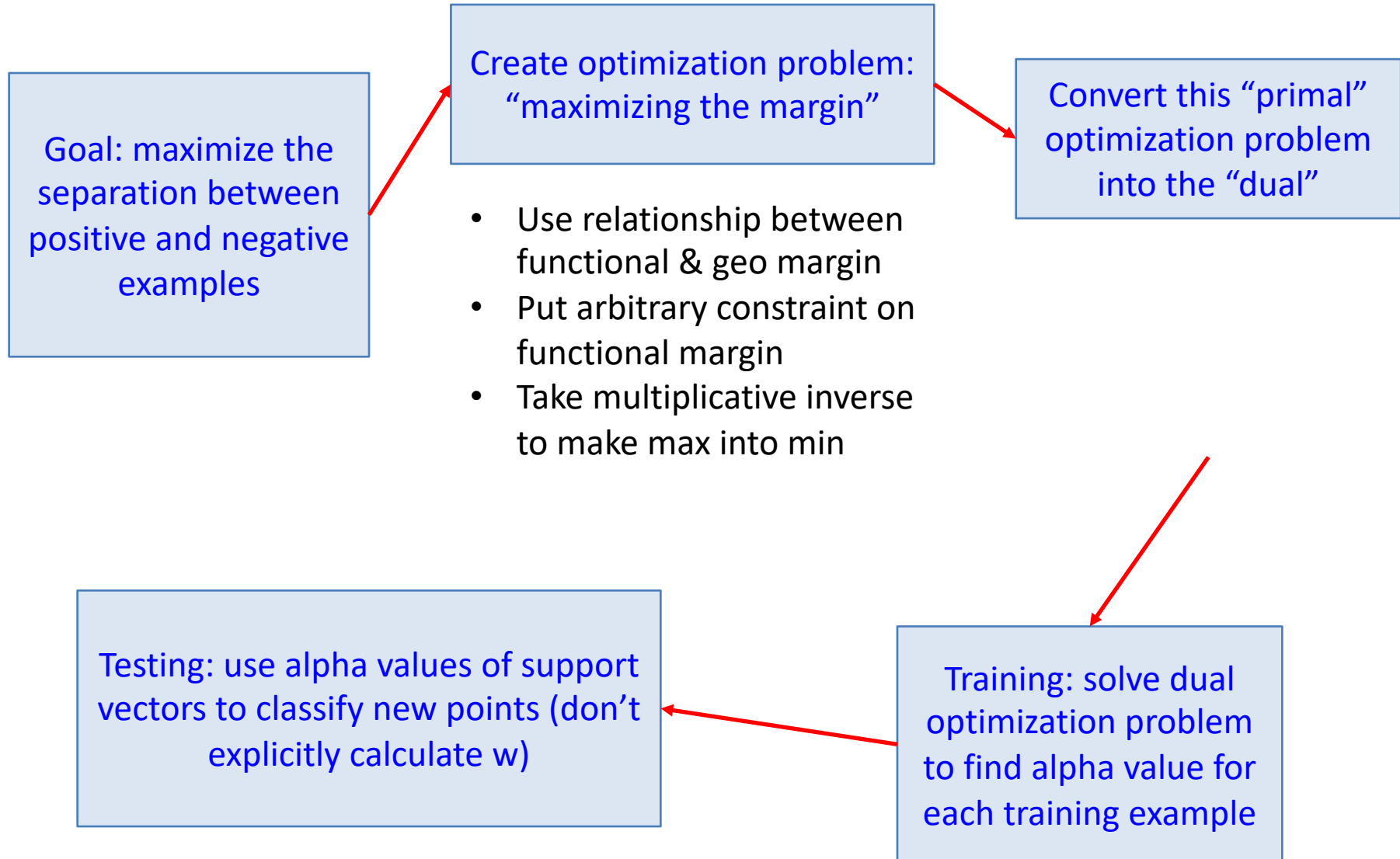
# SVM flowchart



# SVM flowchart

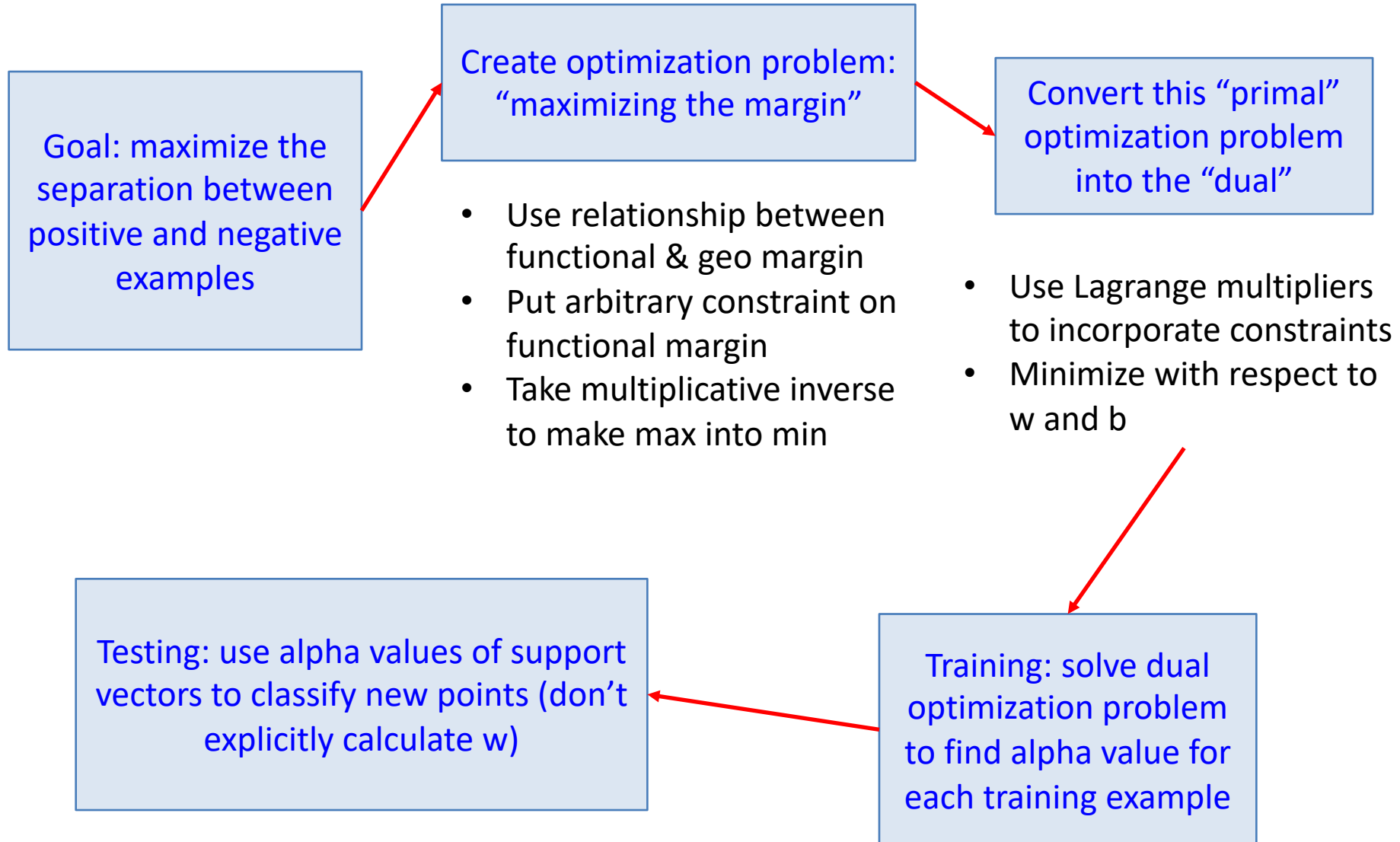


# SVM flowchart





# SVM flowchart



# Functional and Geometric Margins

SVM classifier:  
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

# Functional and Geometric Margins

SVM classifier:  
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional Margin:  $\hat{\gamma}_i = y_i(\vec{w} \cdot \vec{x}_i + b)$

# Functional and Geometric Margins

SVM classifier:  
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional Margin:  $\hat{\gamma}_i = y_i(\vec{w} \cdot \vec{x}_i + b)$

Geometric Margin:  
(distance between  
example and hyperplane)

$$\gamma_i = y_i \left( \frac{\vec{w}}{\|\vec{w}\|} \cdot \vec{x}_i + \frac{b}{\|\vec{w}\|} \right)$$

# Functional and Geometric Margins

SVM classifier:  
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional Margin:  $\hat{\gamma}_i = y_i(\vec{w} \cdot \vec{x}_i + b)$

Geometric Margin:  
(distance between  
example and hyperplane)

$$\gamma_i = y_i \left( \frac{\vec{w}}{\|\vec{w}\|} \cdot \vec{x}_i + \frac{b}{\|\vec{w}\|} \right)$$

Note:

$$\gamma_i = \frac{\hat{\gamma}_i}{\|\vec{w}\|}$$

# Optimization Problem: try 1

Goal: maximize the minimum distance  
between example and hyperplane

$$\gamma = \min_{i=1, \dots, n} \gamma_i$$

# Optimization Problem: try 1

Goal: maximize the minimum distance  
between example and hyperplane

$$\gamma = \min_{i=1, \dots, n} \gamma_i$$

Formulation: optimize a function with  
respect to a constraint

$$\max_{\gamma, \vec{w}, b} \quad \gamma$$

$$\text{s.t.} \quad y_i(\vec{w} \cdot \vec{x}_i + b) \geq \gamma, \quad i = 1, \dots, n$$

$$\text{and} \quad \|\vec{w}\| = 1$$

(force functional and geometric  
margin to be equal)

# Optimization Problem: try 2

Idea: substitute functional margin  
divided by magnitude of weight vector

$$\begin{aligned} \max_{\hat{\gamma}, \vec{w}, b} \quad & \frac{\hat{\gamma}}{\|\vec{w}\|} \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq \hat{\gamma}, \quad i = 1, \dots, n \end{aligned}$$

(gets rid of non-convex constraint)



# Optimization Problem: try 3

Idea: put arbitrary constraint on functional margin

$$\hat{\gamma} = 1$$

$$\begin{array}{ll} \min_{\vec{w}, b} & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t.} & y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1, \quad i = 1, \dots, n \end{array}$$

# Optimization Problem: try 3

Idea: put arbitrary constraint on functional margin

$$\hat{\gamma} = 1$$

$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t.} \quad & -y_i(\vec{w} \cdot \vec{x}_i + b) + 1 \leq 0, \quad i = 1, \dots, n \end{aligned}$$

# Lagrangian

- The alpha values are our Lagrange multipliers
- We don't care about our constraint if it is not *active*

$$\mathcal{L}(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i(\vec{w} \cdot \vec{x}_i + b) - 1]$$

# Lagrangian

- The alpha values are our Lagrange multipliers
- We don't care about our constraint if it is not *active*

$$\mathcal{L}(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i(\vec{w} \cdot \vec{x}_i + b) - 1]$$

- First minimize with respect to w & b, becomes W(alpha)

$$\mathcal{L}(\vec{w}, b, \vec{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \vec{x}_i \cdot \vec{x}_j$$

# Lagrangian

- The alpha values are our Lagrange multipliers
- We don't care about our constraint if it is not *active*

$$\mathcal{L}(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i(\vec{w} \cdot \vec{x}_i + b) - 1]$$

- First minimize with respect to w & b, becomes W(alpha)

$$\mathcal{L}(\vec{w}, b, \vec{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \boxed{\vec{x}_i \cdot \vec{x}_j}$$

# Kernel Trick

- Now we can replace dot products with any kernel!

$$W(\vec{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j K(\vec{x}_i, \vec{x}_j)$$

# Final Goal: classification

- After using Kernel Trick with dual optimization problem, we have:

$$\vec{w}^* = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$$

- To classify, we use:

$$h(\vec{x}) = \text{sign} (\vec{w}^* \cdot \vec{x} + b^*)$$

# Final Goal: classification

- After using Kernel Trick with dual optimization problem, we have:

$$\vec{w}^* = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$$

- To classify, we use:

$$h(\vec{x}) = \text{sign} (\vec{w}^* \cdot \vec{x} + b^*)$$

$$h(\vec{x}) = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i \vec{x}_i \cdot \vec{x} + b^* \right)$$



# Final Goal: classification

- After using Kernel Trick with dual optimization problem, we have:

$$\vec{w}^* = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$$

- To classify, we use:

$$h(\vec{x}) = \text{sign} (\vec{w}^* \cdot \vec{x} + b^*)$$

$$h(\vec{x}) = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i \vec{x}_i \cdot \vec{x} + b^* \right)$$

$$h(\vec{x}) = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i K(\vec{x}_i, \vec{x}) + b^* \right)$$

# Outline for November 21

- Support Vector Machines Review
- Likelihood functions (Bernoulli and Logistic Regression)
- Finish practice problems (Q2-Q4)

$$\vec{y} = \underbrace{1}_p \underbrace{0}_{1-p} \underbrace{0}_{1-p} \underbrace{1}_p \leftarrow \text{given; want } (p)$$

$$L(p | \vec{y}) = p^{\text{\#heads}} (1-p)^{\text{\#tails}}$$

$$L(p | \vec{y}) = \prod_{i=1}^n \underbrace{p^{y_i}}_{y_i=1} \underbrace{(1-p)^{1-y_i}}_{y_i=0}$$

MLE  $\left\{ \begin{array}{l} \cdot \log \\ \cdot \text{derivative wrt } p \\ \cdot \text{set to } 0 \\ \cdot \text{solve for } \hat{p} \end{array} \right.$

log reg

$$h(\vec{w} | \vec{x}_i, y_i, i=1 \dots n)$$

$$= \prod_{i=1}^n \underbrace{h_{\vec{w}}(\vec{x}_i)^{y_i}}_{\substack{\text{prob} \\ (+)}} \underbrace{(1-h_{\vec{w}}(\vec{x}_i))^{1-y_i}}_{\substack{\text{prob} \\ (-)}}$$

# Outline for November 21

- Support Vector Machines Review
- Likelihood functions (Bernoulli and Logistic Regression)
- Finish practice problems (Q2-Q4)

# Follow ups on questions from class

- **ANOVA**: considered a special case of linear regression
- **AdaBoost**: why  $\frac{1}{2}$  in front of the score?
  - Comes out in the derivation
  - Main idea: solve for the classifier scores that minimize exponential loss

# Handout 19, Question 2

- $r = 1/3$ , probability of one classifier being wrong
- $T = 5$ , number of classifiers
- $R$  = number of votes for the wrong class
- If  $R=3,4,5$  then we will vote for the wrong class overall



Q2  $T=5$  # classifiers

$R = \# \text{ wrong}$

$r = \text{prob of one being wrong.} \left. \vphantom{\text{prob of one being wrong.}} \right\} \frac{1}{3}$

$$R > \frac{T}{2} \quad \underbrace{\hspace{1cm}}_{2.5}$$

$$\left\{ \begin{array}{l} R=5 \rightarrow \overset{\text{prob}}{\binom{5}{5}} \left(\frac{1}{3}\right)^5 \\ R=4 \rightarrow \binom{5}{4} \left(\frac{1}{3}\right)^4 \left(\frac{2}{3}\right) \\ R=3 \rightarrow \binom{5}{3} \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^2 \end{array} \right\} \oplus \boxed{\approx 0.21}$$

prob that ensemble is wrong

Q3  $n=2$

3 unique datasets  $\left\{ \begin{array}{l} \{x_1, x_1\} \\ \{x_1, x_2\} \\ \{x_2, x_2\} \end{array} \right.$

$n=3$

$\Rightarrow \underline{10}$  datasets

$$3 + 3 \cdot 2 + 1$$

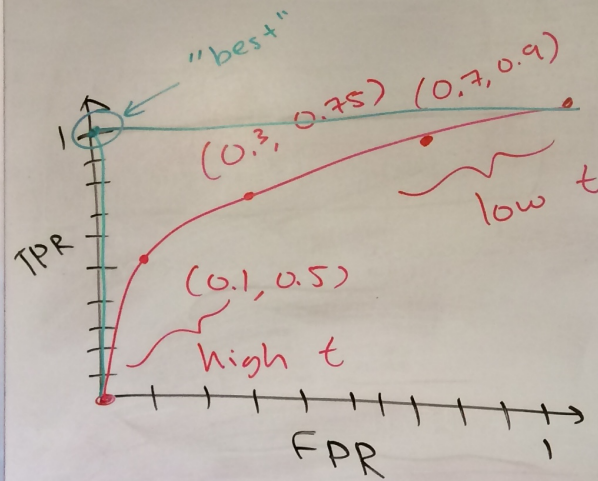
# Handout 19, Question 2

- This analysis assumed classifiers were independent!
- What if they are not? How did Random Forests help us decorrelate classifiers?



# Handout 19, Question 2

- This analysis assumed classifiers were independent!
- What if they are not? How did Random Forests help us decorrelate classifiers?
- Note about Bagging: choosing  $n$  with resampling actually does produce a very different dataset
  - As  $n$  increases, roughly 0.37 not chosen each time



ideal conf matrix

10	0
0	20

Q4