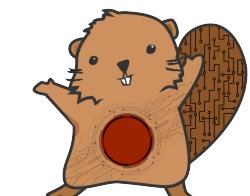




# Machine Learning and Data Mining

Lecture 4.2: Hard-Margin Support Vector Machines



CS 434

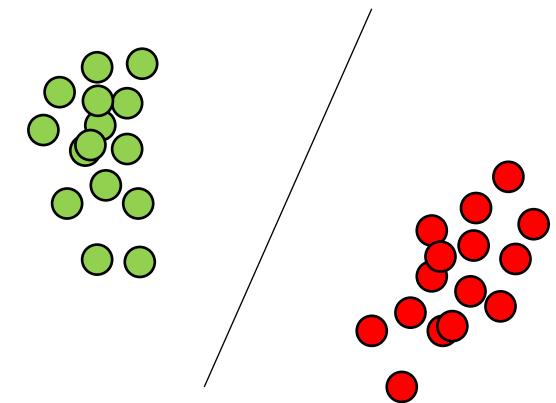


# RECAP

## From Last Lecture

## Discriminative Classifiers:

- Learn  $P(y|x)$  directly
- Logistic regression is one example
- *Nomenclature note -- people will also refer to algorithms that model no distribution as discriminative (such as kNN).*



# Today's Learning Objectives

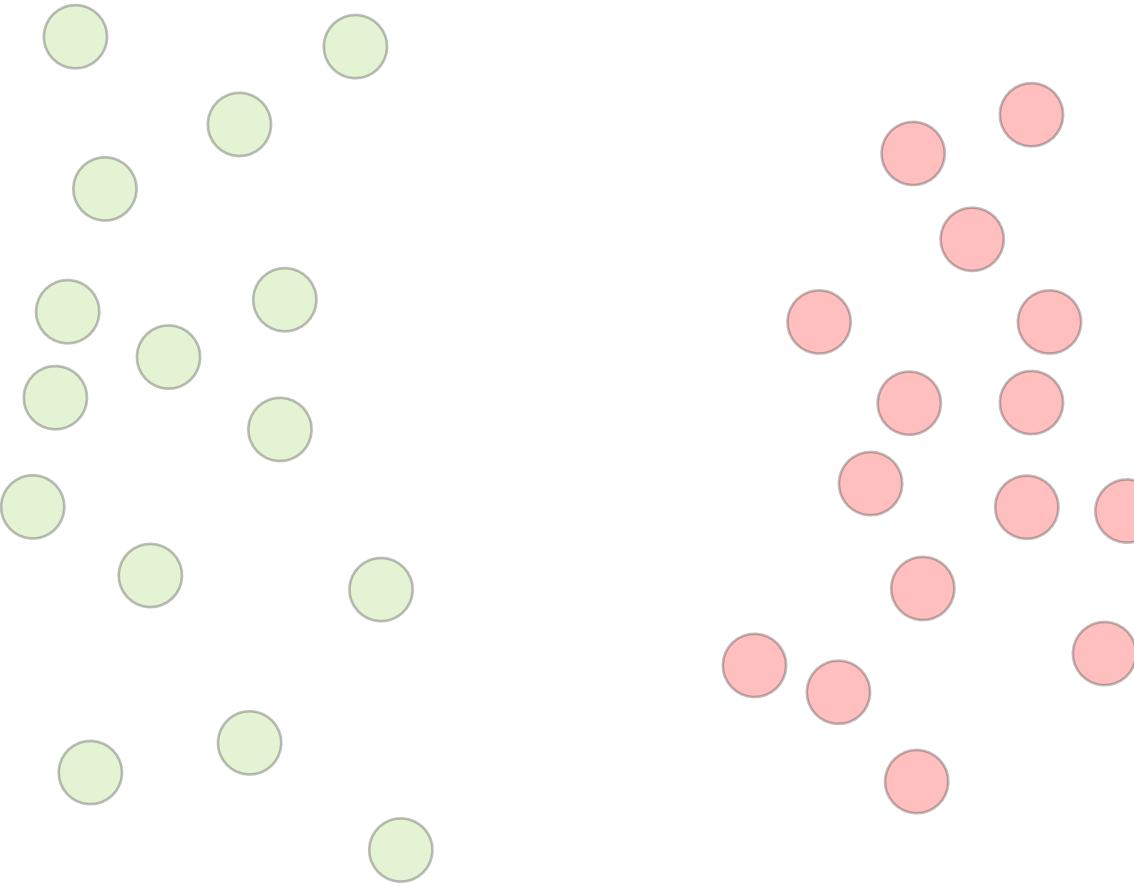


Be able to answer:

- What is the margin of a classifier?
- What is a Support Vector Machine?
  - How is a hard-margin SVM formalized?
  - What is the dual formalization?
  - What is its decision boundary?
  - What is a support vector?
- We'll have to touch on a few math tools to do this:
  - Constrained Optimization with Lagrange Multipliers
  - Quadratic Programming Solvers
  - **You do not need to have an operational understanding of these two techniques at all for this course.**

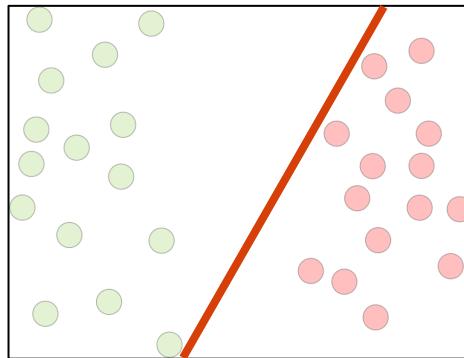


# A Motivating Discussion

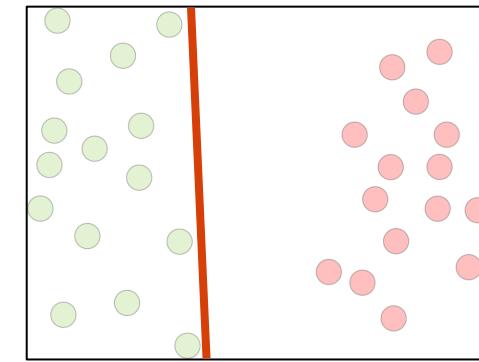




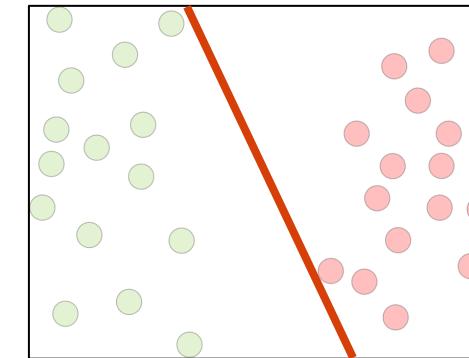
# Question Break!



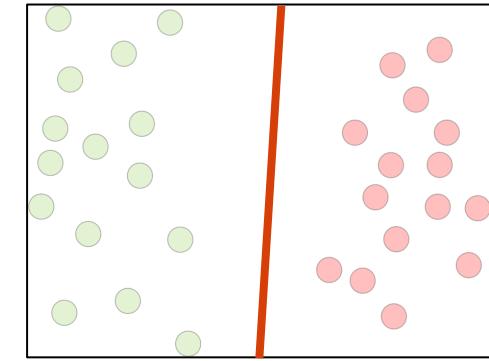
**A**



**B**



**C**



**D**

**A** I prefer A

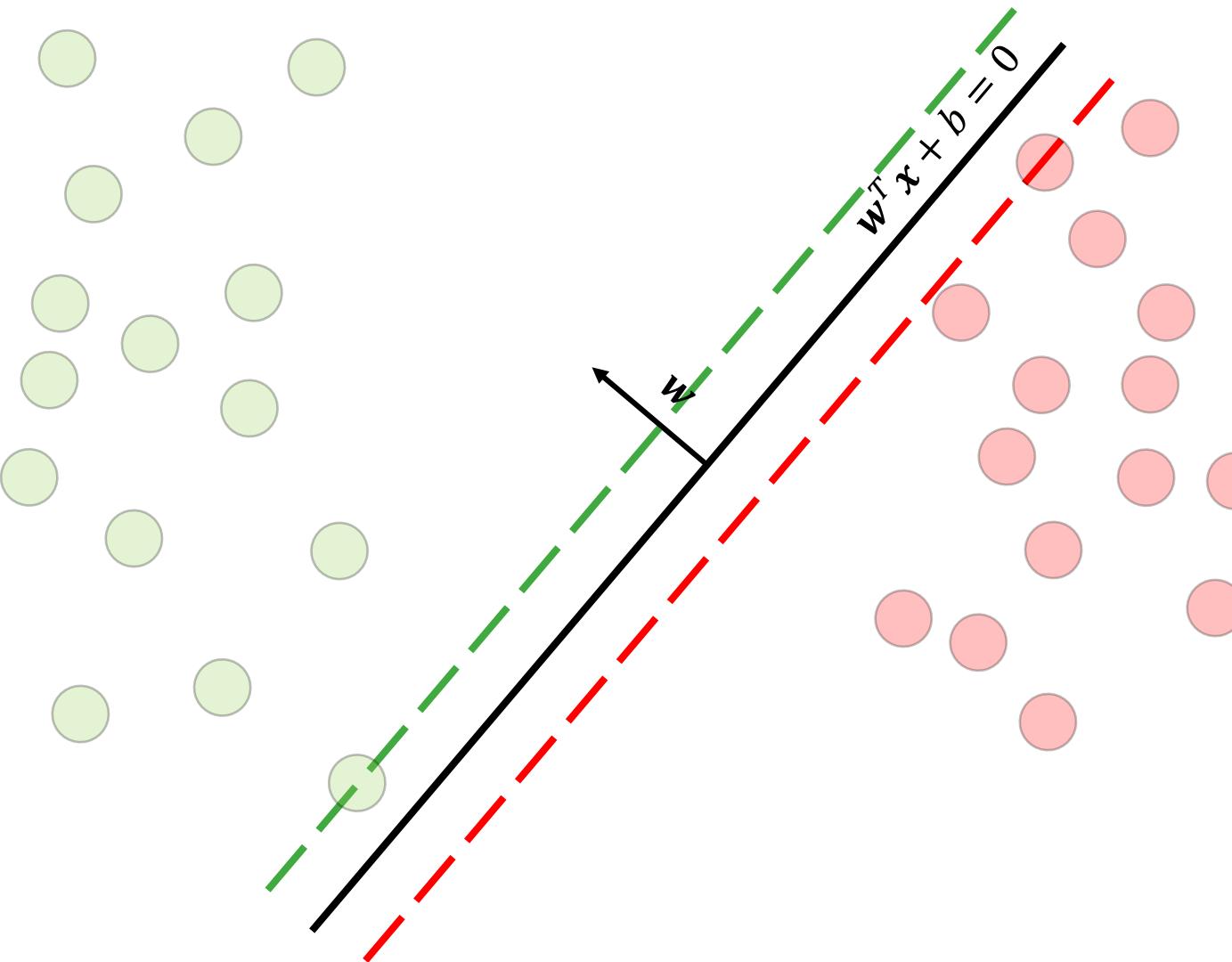
**B** I prefer B

**C** I prefer C

**D** I prefer D

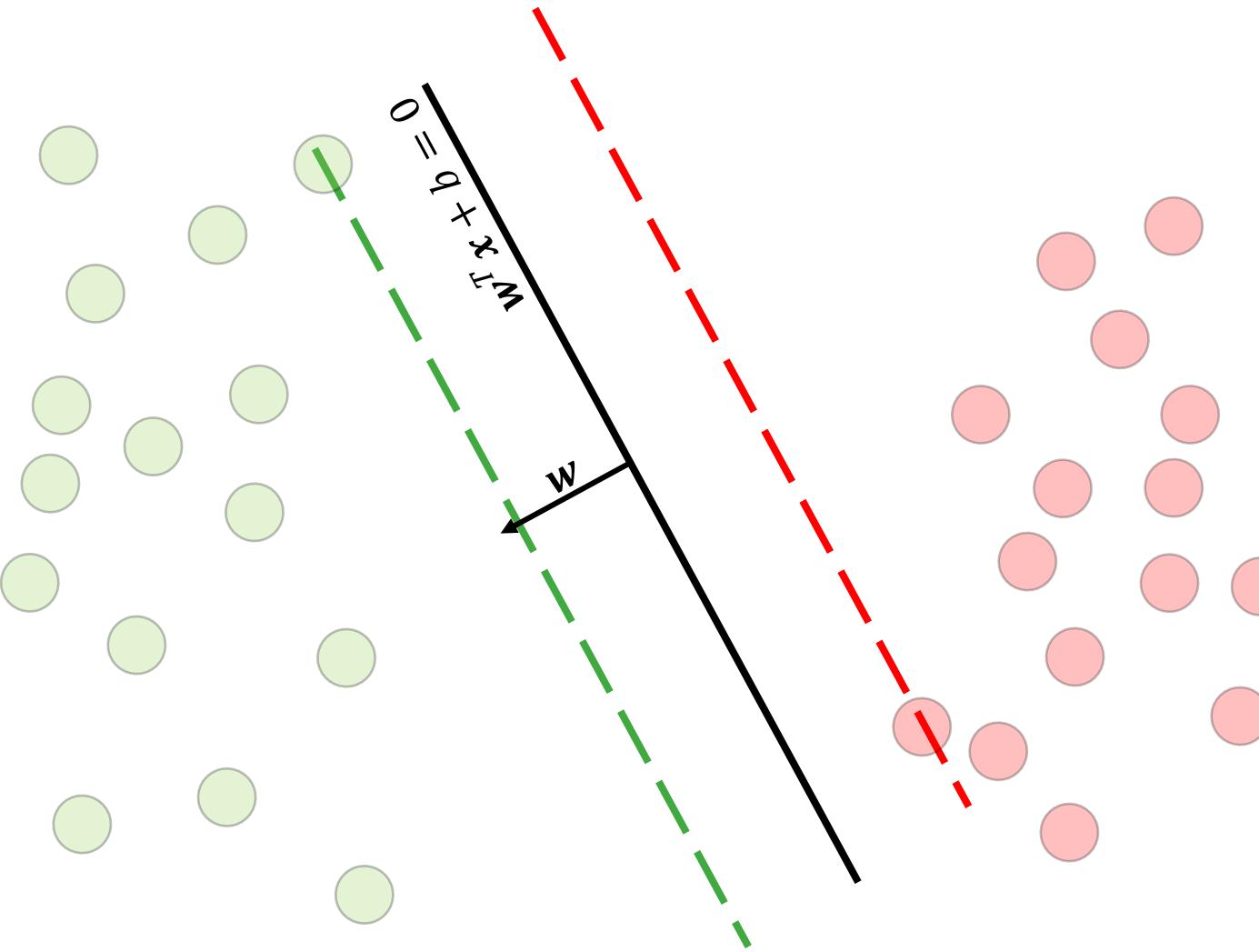


# A Motivating Discussion



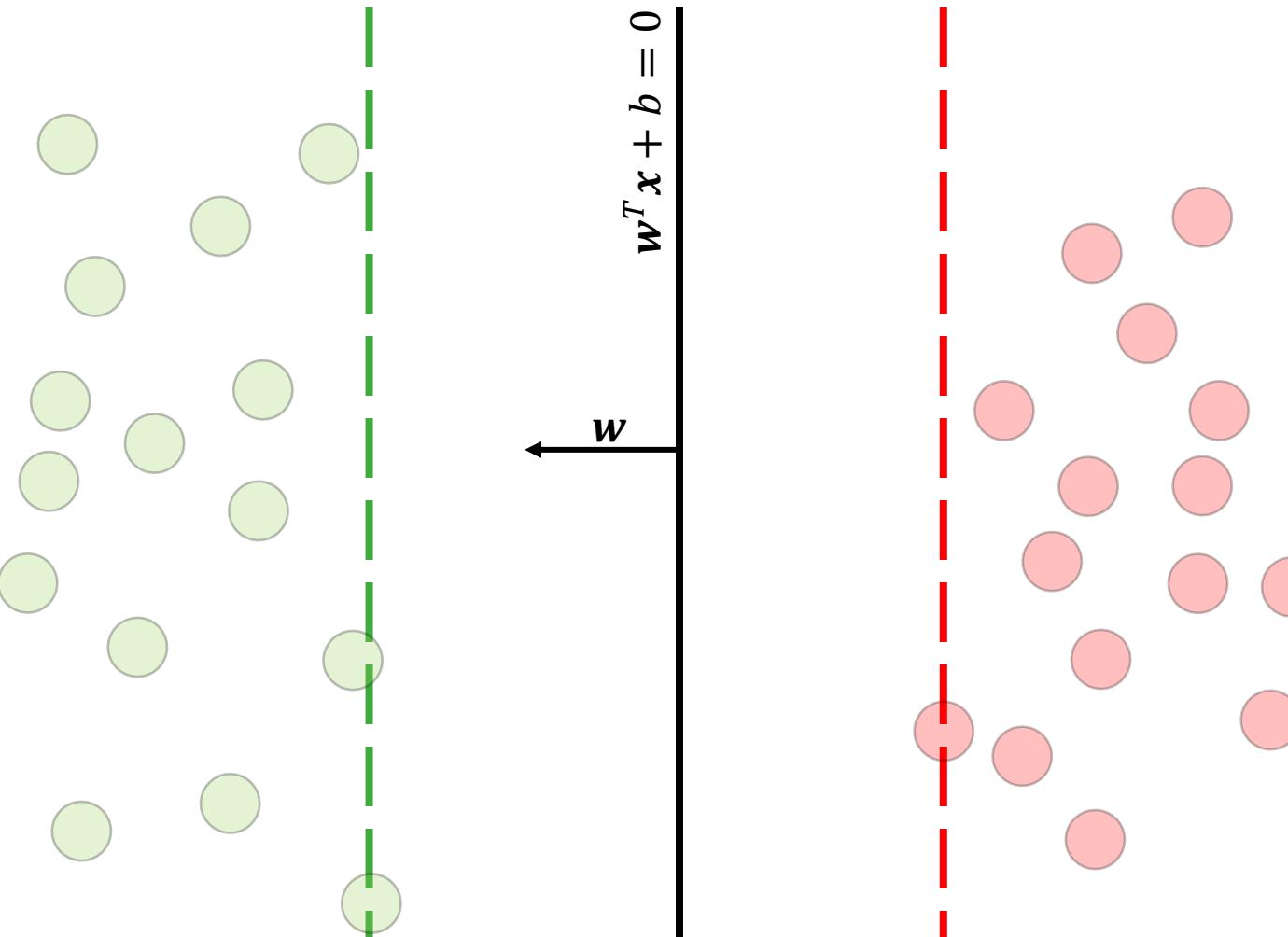


# A Motivating Discussion





# A Motivating Discussion





# New Topic: Support Vector Machines

**Algorithm Name:** Support Vector Machines (SVM)

**Type:** Supervised Linear Classifier

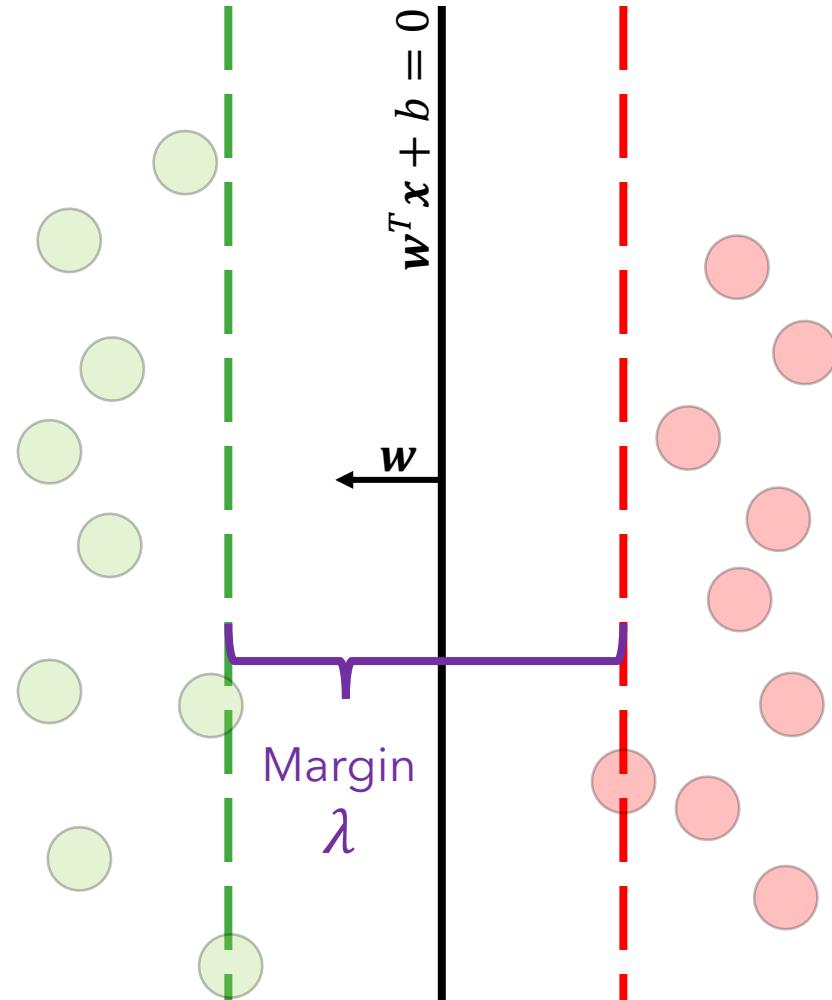
**Key Idea:** Find the linear decision boundary that maximizes the margin between positive and negative examples.



**What is a support vector?** We'll need to run through a lot of somewhat unintuitive math to really understand the structure of this problem. But the result will give us tremendous insight into this problem ... and it'll build character.



<https://jgreitemann.github.io/svm-demo>



## Our goal for today:

We want to find a classifier that separates the classes while maximizing the margin between positive and negative examples. *For now, we'll assume linear separability and fix that later.*

$$\max_{w,b} \text{margin}$$

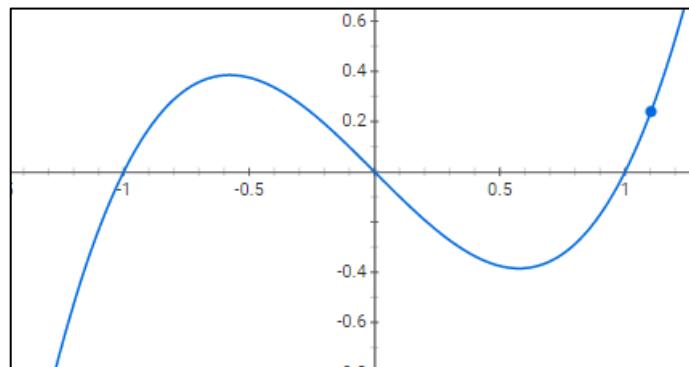
s.t. all examples are correct



## Question Break!



Provide a solution to the following constrained optimization problem:



$$\begin{aligned} \min_x f(x) &= x^3 - x \\ \text{s.t. } &0 \leq x \leq 5 \end{aligned}$$

$$\begin{aligned} 3x^2 - 1 &= 0 \\ \Rightarrow x &= \pm\sqrt{0.33} \Rightarrow x = \sqrt{0.33} \end{aligned}$$

$$\begin{aligned} x = 0 &\Rightarrow f(0) = 0 \\ x = 5 &\Rightarrow f(5) = 120 \\ x = \sqrt{0.33} &\Rightarrow f(0.57) \approx -0.38 \end{aligned}$$

A

0

B

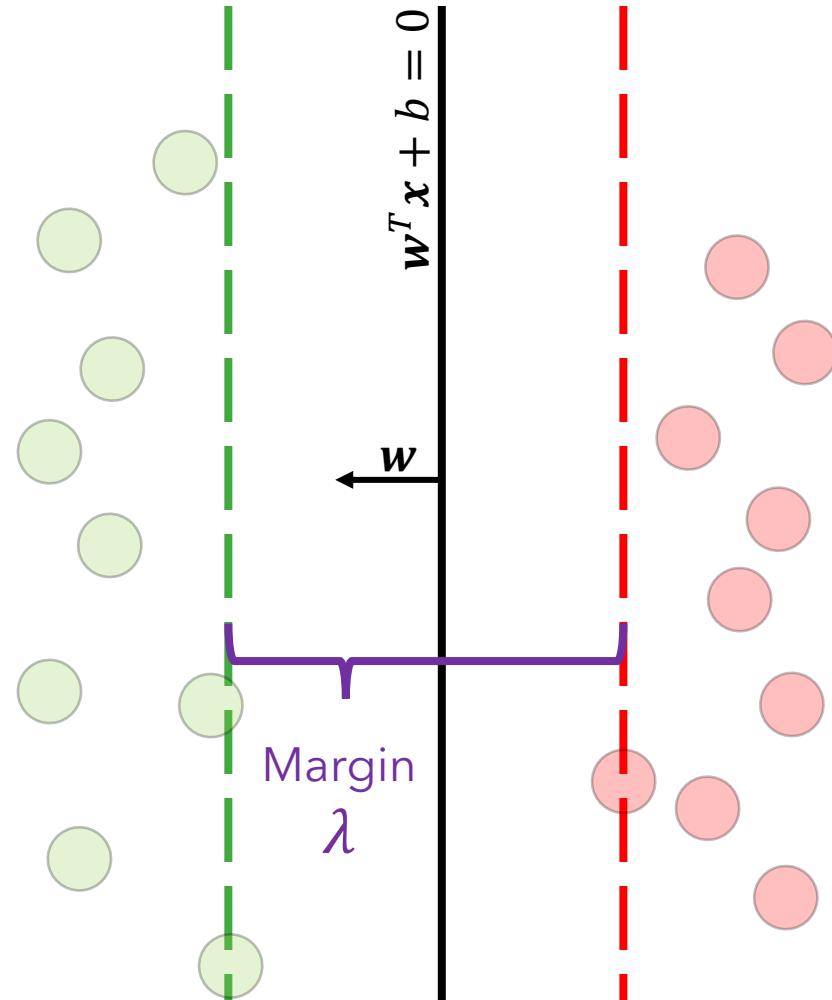
5

C

$\sqrt{0.33}$

D

$-\sqrt{0.33}$



## Our goal for today:

We want to find a classifier that separates the classes while maximizing the margin between positive and negative examples. *For now, we'll assume linear separability and fix that later.*

## Our plan of attack:

1. Set up some convenient notation
2. Derive an expression for our intuition
3. Set up a constrained optimization problem
4. Solve it to arrive at a Support Vector Machine  
(some new math tools to analyze this- Lagrange multipliers)



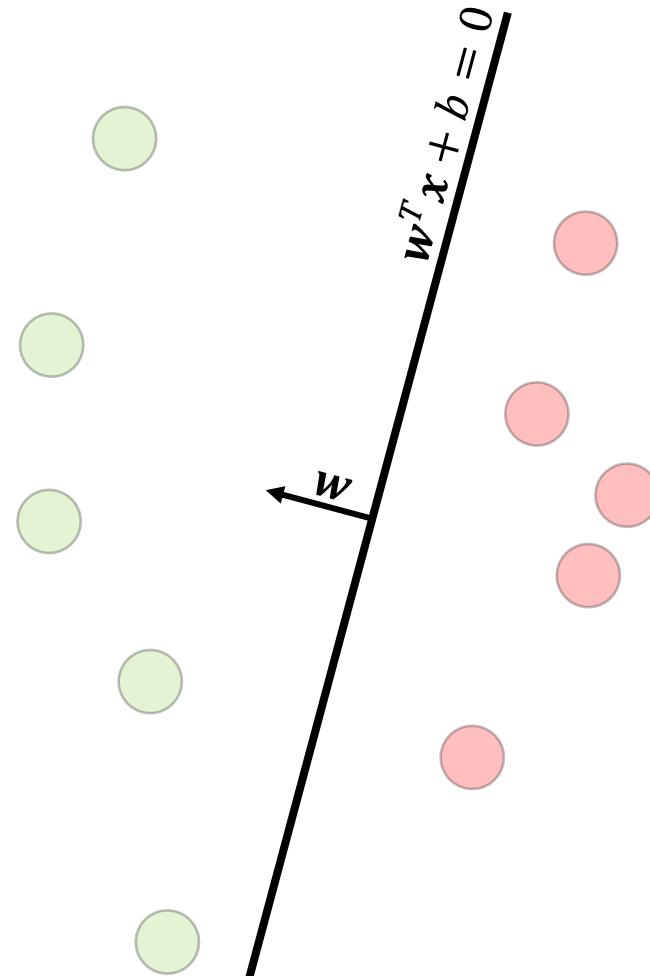


## Step 1) Establish Some Convenient Notation

### Problem Setup and Decision Rule

Consider a binary classification setting where we represent the classes as  $y \in \{-1,1\}$

Let our linear decision rule be:



$$\mathbf{w}^T \mathbf{x} + b > 0 \Rightarrow \text{predict } 1$$

$$\mathbf{w}^T \mathbf{x} + b < 0 \Rightarrow \text{predict } -1$$



## Question Break!



I can tell if a point  $x_i$  with label  $y_i$  is correctly or incorrectly classified by checking if:

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) > 0$$

**A** True

**B** False

**C** I don't know

**D**



## Step 1) Establish Some Convenient Notation

### Problem Setup and Decision Rule

Consider a binary classification setting where we represent the classes as  $y \in \{-1,1\}$

Let our linear decision rule be:

$$\mathbf{w}^T \mathbf{x} + b > 0 \Rightarrow \text{predict } 1$$

$$\mathbf{w}^T \mathbf{x} + b < 0 \Rightarrow \text{predict } -1$$

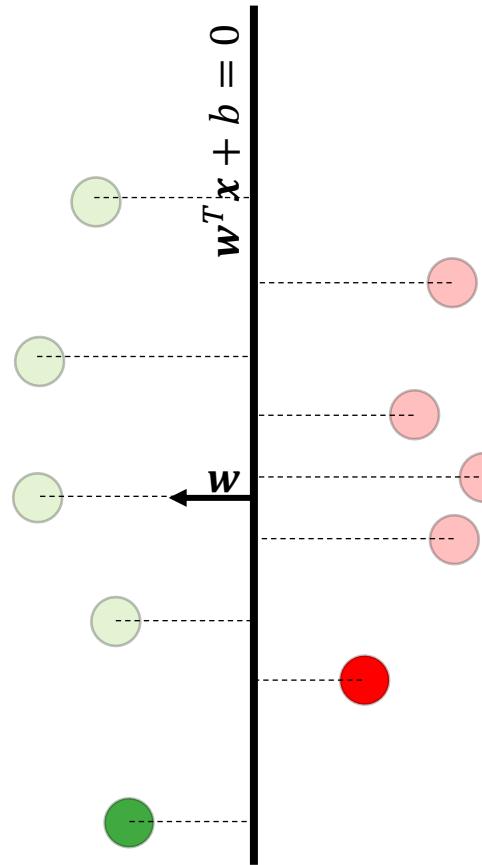
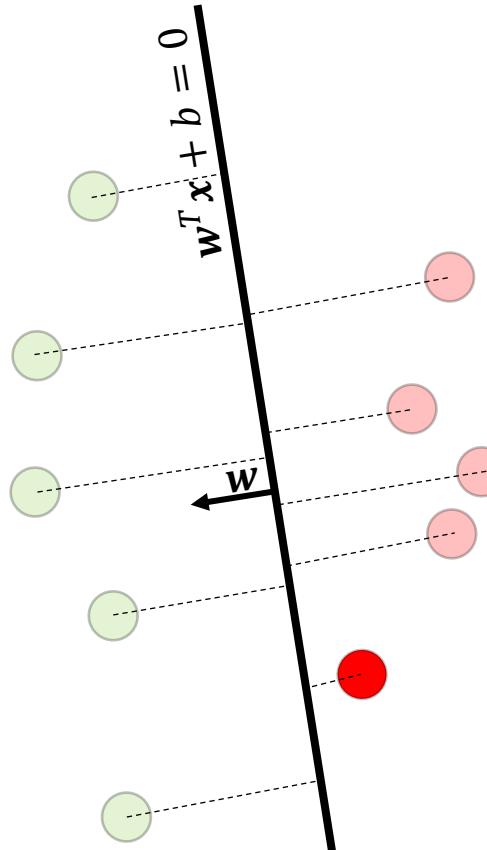
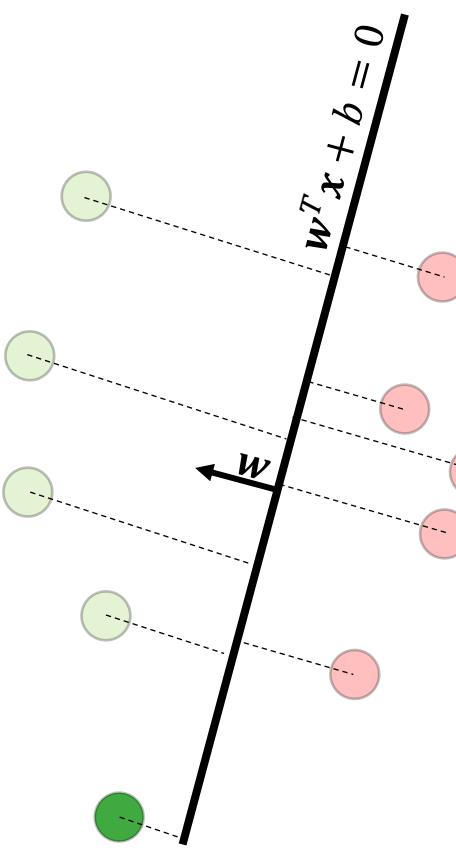
Means a correct classification has:

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) > 0$$



## Step 2) Derive an Expression for the Margin

**Intuition:** One way to talk about the margin of our classifier would be to measure how close the nearest point is to the boundary.

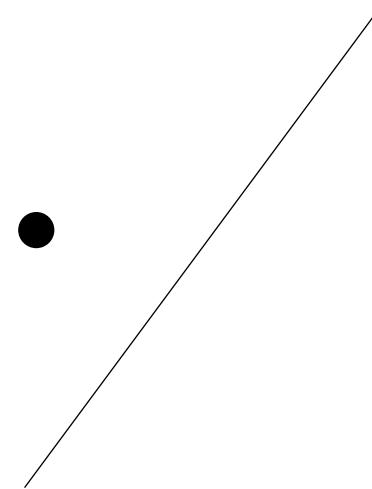




# Question Break!



The shortest path from a point to a line is perpendicular to the original line and parallel to the vector of weights defining it.



**A** True

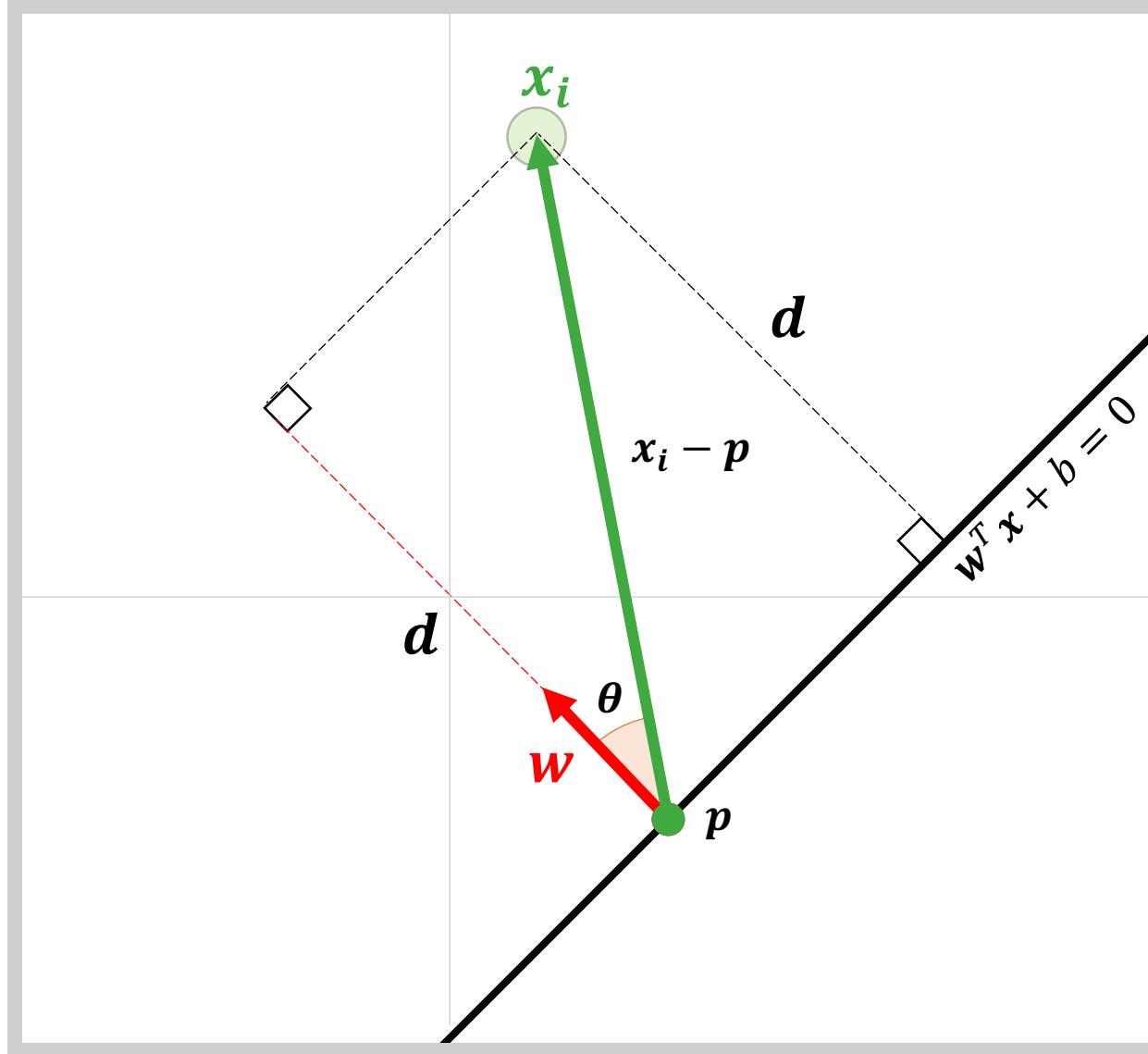
**B** False

**C** I don't know

**D** I don't understand the question.



## Step 2) Derive an Expression for the Margin



**Goal: Find distance  $d$ .**

Cos → adjacent over hypotenuse:

$$\cos(\theta) = \frac{d}{\|x_i - p\|}$$

Definition of dot product:

$$\frac{\mathbf{w}^T(x_i - p)}{\|\mathbf{w}\| \|x_i - p\|} = \frac{d}{\|x_i - p\|}$$

Point  $p$  is on the line so  $\mathbf{w}^T p = -b$

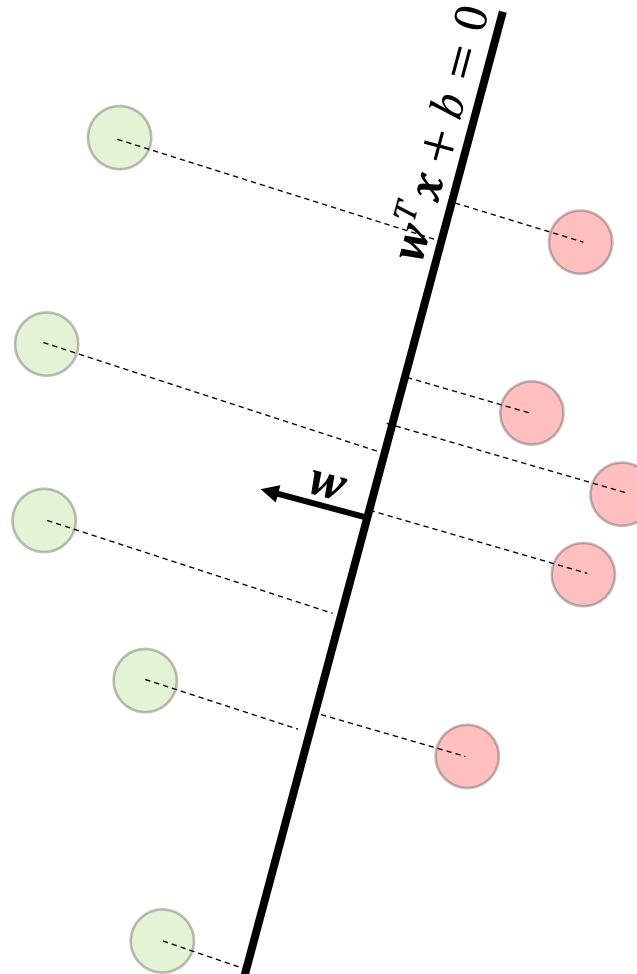
$$d = \frac{\mathbf{w}^T(x_i - p)}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T x_i + b}{\|\mathbf{w}\|}$$

Absolute value to remove sign:

$$d = \frac{|\mathbf{w}^T x_i + b|}{\|\mathbf{w}\|}$$



## Step 3) Set Up A Constrained Optimization Problem



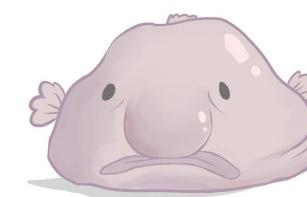
We can then write the distance to the decision boundary from the nearest point as:

$$d_{min} = \min_i \frac{|w^T x_i + b|}{\|w\|}$$

One thing we could do is say we want to maximize this minimum distance without making any errors:

Maximize distance to nearest point

$$\max_{w,b} \min_i \frac{|w^T x_i + b|}{\|w\|}$$



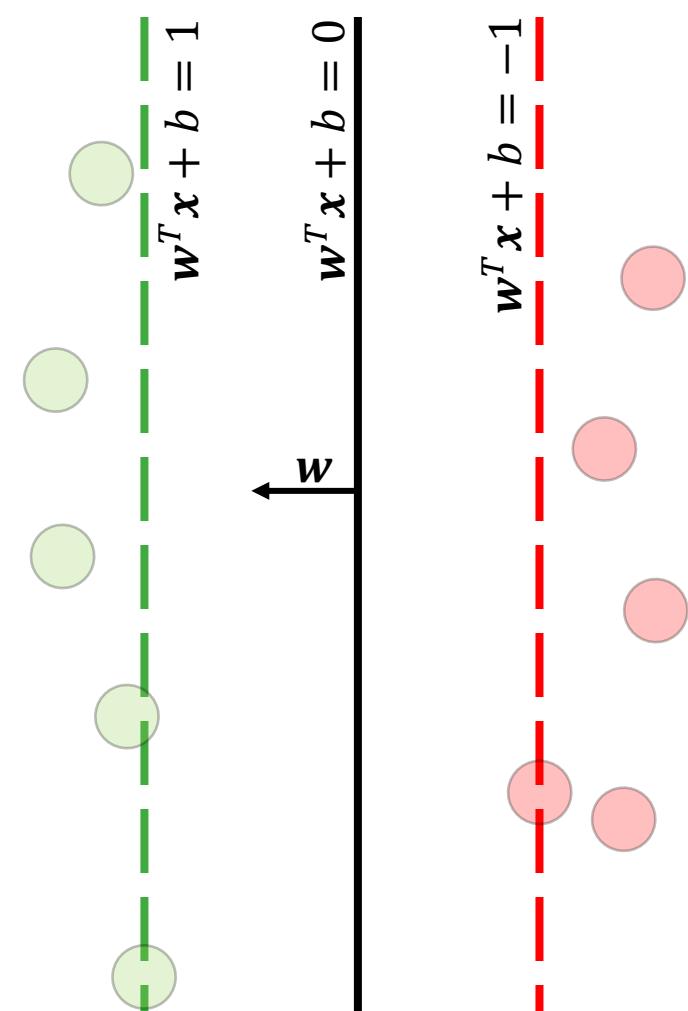
$$s.t. y_i(w^T x_i + b) > 0 \quad \forall i$$

While not making any errors

**This is an ugly problem! Max of min? w's everywhere?**



## Step 3) Set Up A Constrained Optimization Problem



**To avoid dealing with this, let's constrain the minimum distance to some constant w.r.t to  $x$  - say  $1/\|w\|$ :**

$$d_{min} = \min_i \frac{|w^T x_i + b|}{\|w\|} = \frac{1}{\|w\|} \rightarrow |w^T x_i + b| \geq 1 \quad \forall i$$

Can combine with "correctness" constraint  $y_i(w^T x_i + b) > 0$ :

$$\rightarrow y_i(w^T x_i + b) \geq 1 \quad \forall i$$

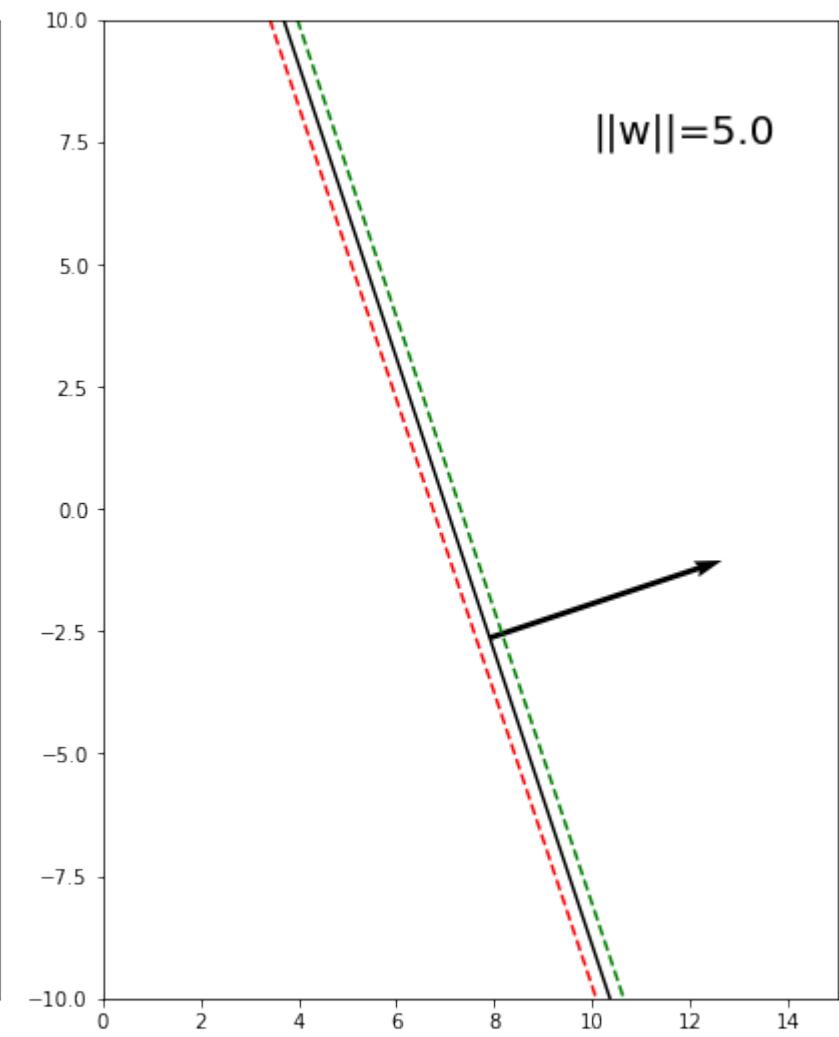
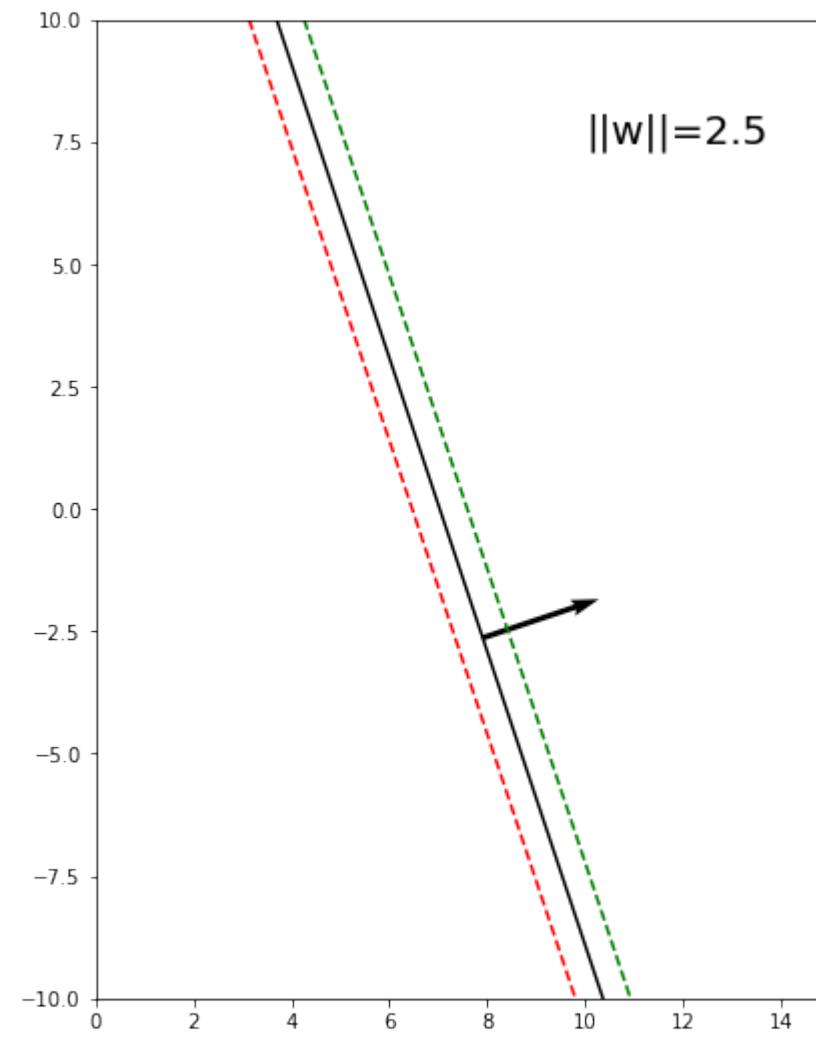
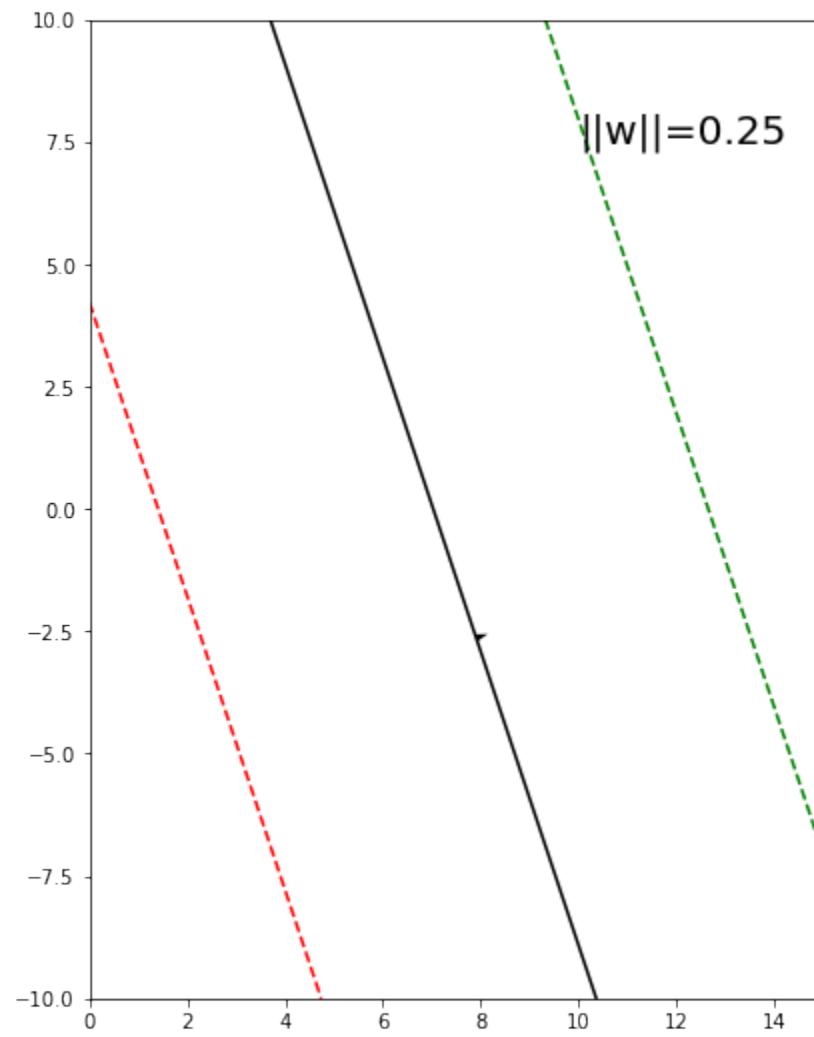
**What does this do?**

- Makes a space around the classifier where points aren't allowed to be - between the green and red lines.
- $w$  can be rescaled arbitrarily to make this true without changing the classifier's expressiveness.



## Step 3) Set Up A Constrained Optimization Problem

$\mathbf{w}$  can be rescaled to make this true without changing the classifier's expressiveness.

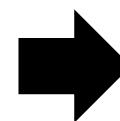




## Step 3) Set Up A Constrained Optimization Problem

**How does this change our “ugly” objective to something nicer?**

$$\max_{w,b} \min_i \frac{|w^T x_i + b|}{\|w\|}$$



$$\max_{w,b} \frac{1}{\|w\|}$$

$$s.t. y_i(w^T x_i + b) > 0 \quad \forall i$$

$$s.t. y_i(w^T x_i + b) \geq 1 \quad \forall i$$



**You:** Neat math. Cute fish. How does this relate to the margin?



## Question Break!



For a dataset  $\{x_i, y_i\}_{i=1}^n$ , how many constraints do we have?

$$\max_{w,b} \frac{1}{\|w\|}$$

$$s.t. y_i(w^T x_i + b) \geq 1 \quad \forall i$$

A

1

B

d

C

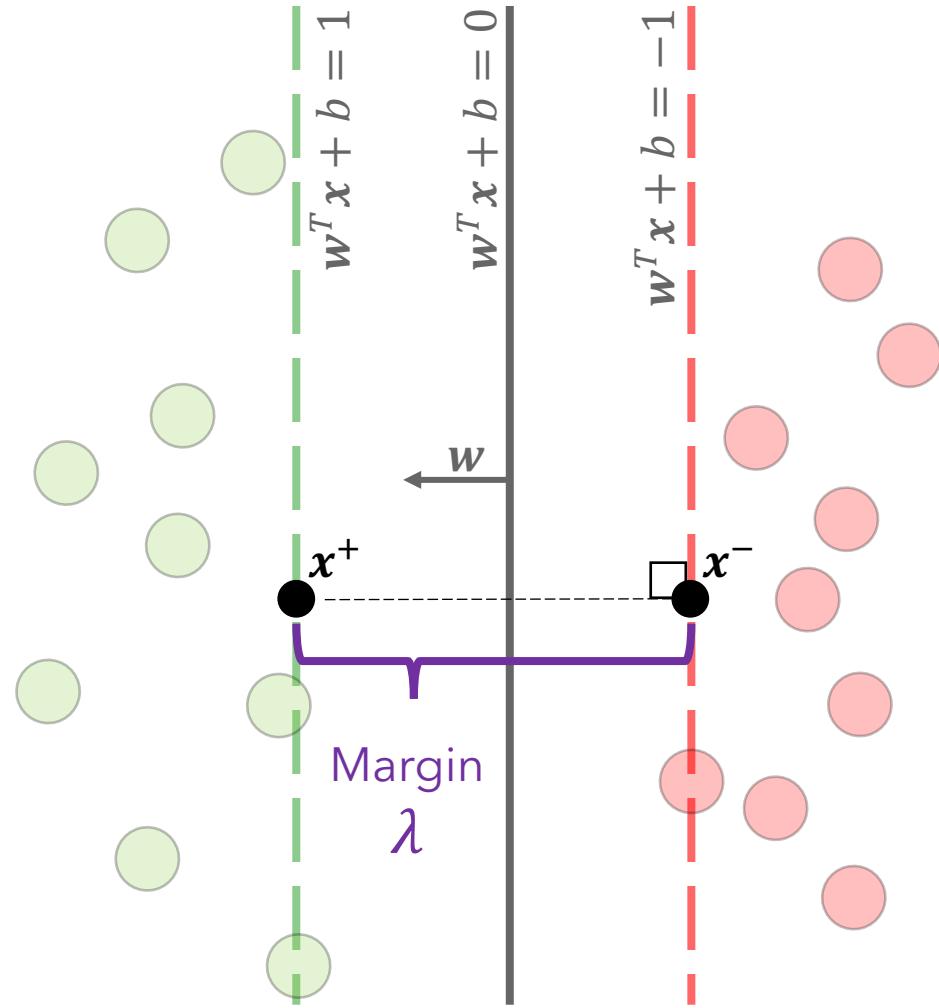
n

D

2



## Step 2) Derive an Expression for the Margin (Back Again)



Things we assumed plus some geometry:

$$\begin{aligned} 1) \quad & \mathbf{w}^T \mathbf{x}^- + b = -1 \\ 2) \quad & \mathbf{w}^T \mathbf{x}^+ + b = 1 \end{aligned}$$

**Goal:** Derive an expression for the margin  $\lambda$

$$3) \Rightarrow \mathbf{x}^+ = \lambda \mathbf{w} + \mathbf{x}^-$$

$$2) \Rightarrow \mathbf{w}^T(\lambda \mathbf{w} + \mathbf{x}^-) + b = 1$$

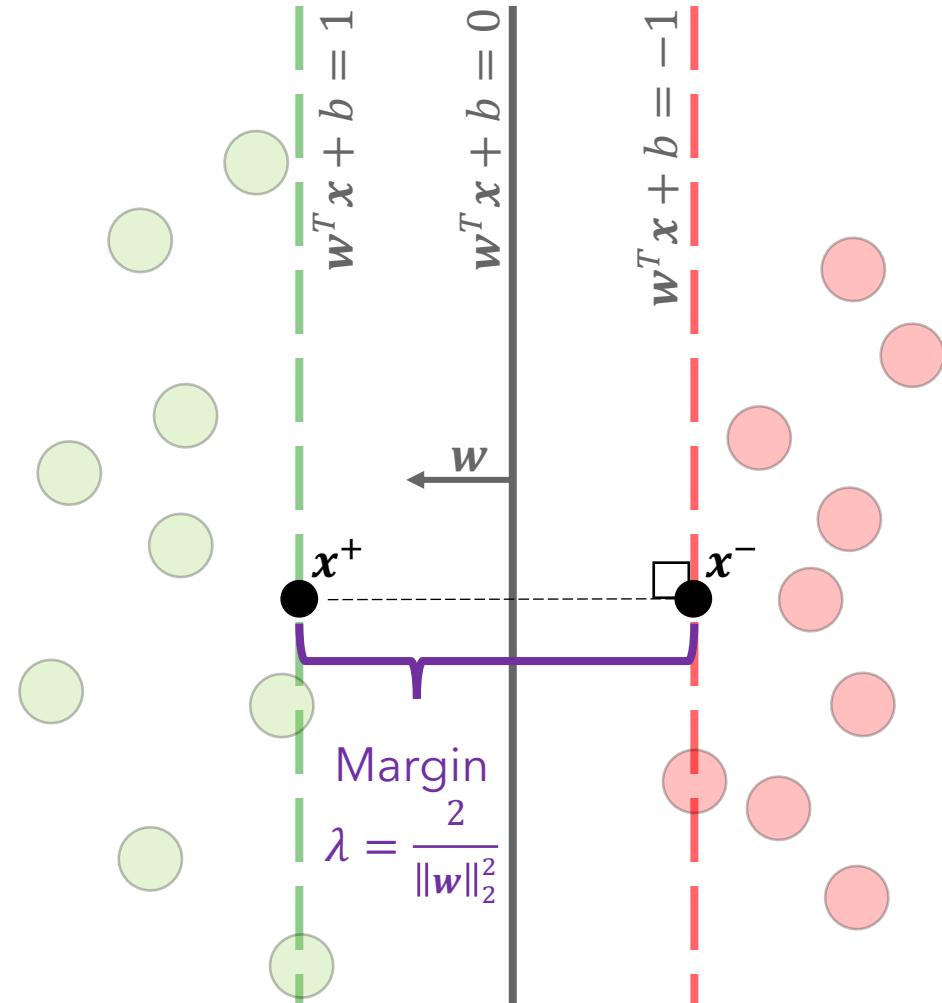
$$\lambda \mathbf{w}^T \mathbf{w} + \mathbf{w}^T \mathbf{x}^- + b = 1$$

$$1) \Rightarrow \lambda \mathbf{w}^T \mathbf{w} - 1 = 1$$

$$\Rightarrow \lambda = \frac{2}{\mathbf{w}^T \mathbf{w}} = \frac{2}{\|\mathbf{w}\|_2^2}$$



## Step 3) Set Up A Constrained Optimization Problem (Again)



**Maximize margin without making errors:**

$$\max_{w,b} \frac{2}{\|w\|_2^2}$$

s.t.  $y_i(w^T x_i + b) \geq 1 \quad \forall i$

Maximize  
margin

Such that no points  
are misclassified



## Recap Steps 1-3: We derived two objectives. How do they relate?

**Version 1) Maximize minimum distance to boundary without making errors:**

$$\max_{w,b} \frac{1}{\|w\|_2}$$

$$s.t. y_i(w^T x_i + b) \geq 1 \quad \forall i$$

**Version 2) Maximize margin width without making errors**

$$\max_{w,b} \frac{2}{\|w\|_2^2}$$

$$s.t. y_i(w^T x_i + b) \geq 1 \quad \forall i$$

These are proportional! Both are maximized by the same  $w^*$ .

Could also write this problem as:

$$\min_{w,b} \frac{1}{2} w^T w$$

$$s.t. y_i(w^T x_i + b) \geq 1 \quad \forall i$$



# Hard-Margin Support Vector Machine Objective

How do we actually solve this? And by "solve" I mean optimal find  $\mathbf{w}$  and  $b$

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$s.t. \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

This is a quadratic function with linear constraints. Many (non-trivial) algorithms to solve these sorts of things - called Quadratic Program Solvers (QP-Solvers).



# Hard-Margin Support Vector Machine Objective

**What is a Quadratic Program Solver?** Given equations of the form:

$$\min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T P \mathbf{z} + \mathbf{q}^T \mathbf{z}$$

Quadratic equation in matrix notation

$$s.t. \quad G\mathbf{z} \leq \mathbf{h}$$

Linear inequality constraints

$$A\mathbf{z} = \mathbf{b}$$

Linear equality constraints

$$\mathbf{b}_{lower} \leq \mathbf{z} \leq \mathbf{b}_{upper}$$

Limits on the values of  $\mathbf{x}$

Runs algorithms from optimization to provide the vector  $\mathbf{z}^*$  that minimizes the function without violating the constraints. (Or returns infeasible if no such solution exists)

- One example from Python (<https://pypi.org/project/qpsolvers/>)

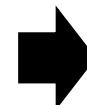


# Hard-Margin Support Vector Machine Objective

How can we put this into a quadratic form?

$$\min_{w,b} \frac{1}{2} w^T w$$

$$s.t. \quad y_i(w^T x_i + b) \geq 1 \quad \forall i$$



$$\min_z \frac{1}{2} z^T P z + q^T z$$

$$s.t. \quad Gz \leq h \\ Az = b$$

$$b_{lower} \leq z \leq b_{upper}$$

$$z = \begin{bmatrix} w \\ b \end{bmatrix}_{d+1 \times 1}$$

$$P = \begin{bmatrix} I_d & 0 \\ 0 & 0 \end{bmatrix}_{d+1 \times d+1}$$

$$q = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}_{d+1}$$

$$h = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{d+1}$$

$$G = \begin{bmatrix} y_1 x_1^T, y_1 \\ y_2 x_2^T, y_2 \\ y_3 x_3^T, y_3 \\ \vdots \\ y_n x_n^T, y_n \end{bmatrix}_{n \times d+1}$$

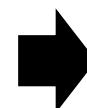


# Hard-Margin Support Vector Machine Objective

How can we put this into a quadratic form?

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

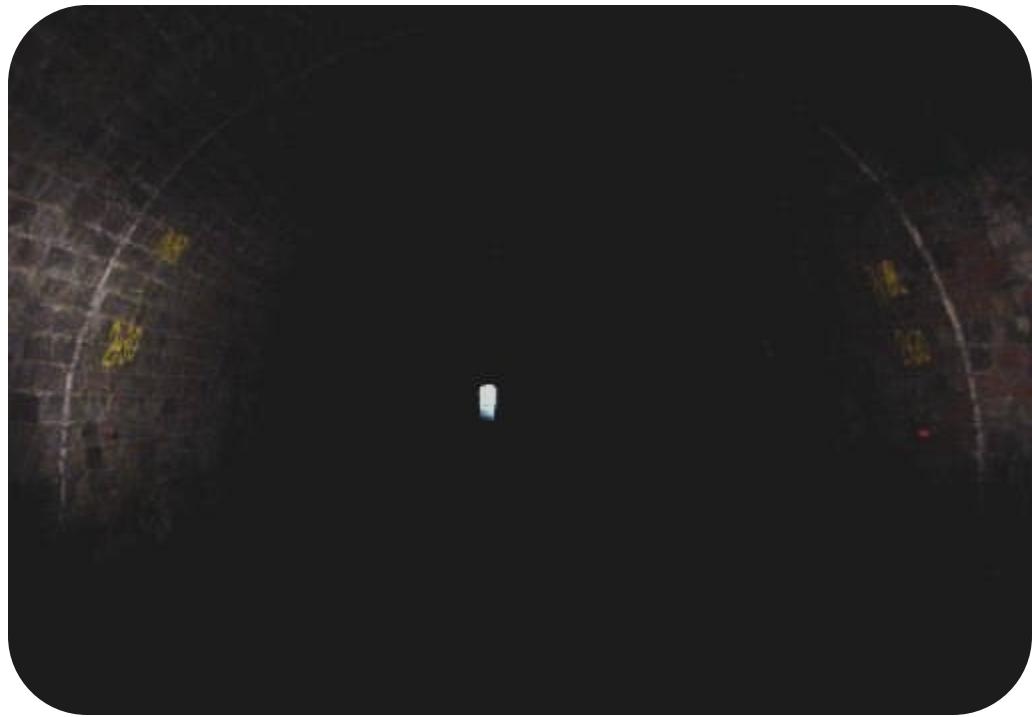
$$s.t. \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$



$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} [\mathbf{w} \ b] \begin{bmatrix} I_d & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} + \vec{0}^T \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \\ & s.t. \begin{bmatrix} y_1 \mathbf{x}_1^T, y_1 \\ y_2 \mathbf{x}_2^T, y_2 \\ y_3 \mathbf{x}_3^T, y_3 \\ \vdots \\ y_n \mathbf{x}_n^T, y_n \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \geq \vec{1} \end{aligned}$$

Can pass this to a QP solver and let some optimization people deal with it.

For our purposes however, there is some deeper bit of understanding to be gained from pushing this further with another method - Lagrange multipliers.





To be *very* clear, the following is beyond the scope of this class and would be content in an optimization course. But we need the tool to move forward with our work here.

I'm going through this for completeness and to give you intuition for when I apply it later.

If you understand it, great. If not, don't worry and focus on the outcome and conclusions.



# Constrained Optimization with Lagrange Multipliers

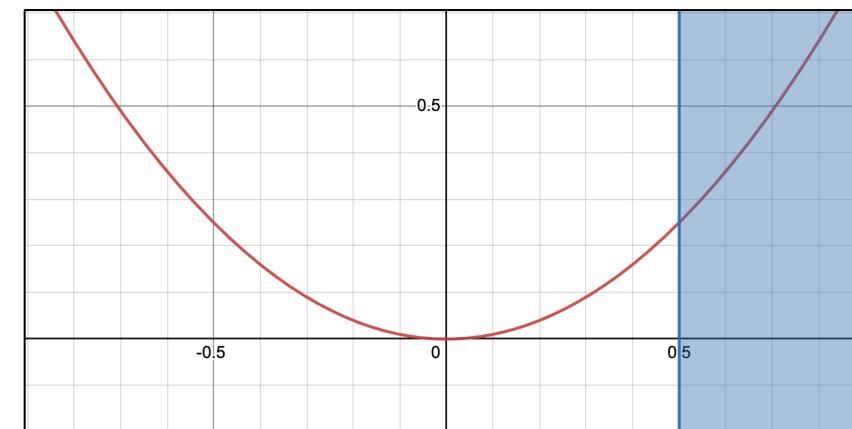
Suppose we need to minimize some function given a set of constraints:

$$\begin{aligned} & \min_x f(x) \\ & s.t. g_i(x) \leq 0 \text{ for } i = 1, 2, \dots, m \end{aligned}$$

---

For example, a quadratic function with one linear constraint:

$$\begin{aligned} & \min_x x^2 \\ & s.t. -10x + 5 \leq 0 \end{aligned}$$



In general,  $f$  and  $g_i$  can be arbitrary.



# Constrained Optimization with Lagrange Multipliers

Write down an unconstrained function called the Lagrangian:

$$\mathcal{L}(x, \alpha_1, \dots, \alpha_m) = f(x) + \sum_i^m \alpha_i g_i(x)$$

The original optimization problem is equivalent to solving the following:

**Primal:**  $\min_x \max_{\alpha_1, \dots, \alpha_m} \mathcal{L}(x, \alpha_1, \dots, \alpha_m) \quad s.t. \alpha_i \geq 0 \forall i$

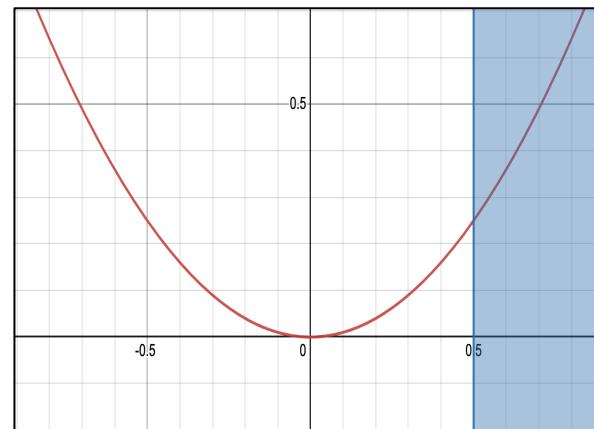
By exchanging the order of min and max, we get the dual problem:

**Dual:**  $\max_{\alpha_1, \dots, \alpha_m} \min_x \mathcal{L}(x, \alpha_1, \dots, \alpha_m) \quad s.t. \alpha_i \geq 0 \forall i$



# Constrained Optimization with Lagrange Multipliers

$$\begin{array}{ll} \min & \color{red} x^2 \\ \text{s.t.} & -10\color{blue}x + 5 \leq 0 \end{array} \quad \rightarrow \quad \mathcal{L}(x, \alpha) = x^2 + \alpha(-10x + 5)$$



Dual:  $\max_{\alpha} \min_x x^2 + \alpha(-10x + 5) \quad \text{s.t. } \alpha_i \geq 0 \ \forall i$

$$\frac{d\mathcal{L}(x, \alpha)}{dx} = 2x - 10\alpha = 0 \Rightarrow x = 5\alpha$$

$$\begin{aligned} \frac{d\mathcal{L}(x, \alpha)}{d\alpha} &= -10x + 5 = 0 \\ &\Rightarrow -50\alpha + 5 = 0 \end{aligned}$$

$$\begin{aligned} \Rightarrow \alpha &= \frac{1}{10} & \Rightarrow x &= \frac{1}{2} \end{aligned}$$



$$\mathcal{L}(x, \alpha_1, \dots, \alpha_m) = f(x) + \sum_i^m \alpha_i g_i(x)$$

Think of these like a 2-player adversarial games with different turn order.

Primal:

$$\min_x \max_{\alpha_1, \dots, \alpha_m} \mathcal{L}(x, \alpha_1, \dots, \alpha_m) \quad s.t. \alpha_i \geq 0 \forall i$$

Dual:

$$\max_{\alpha_1, \dots, \alpha_m} \min_x \mathcal{L}(x, \alpha_1, \dots, \alpha_m) \quad s.t. \alpha_i \geq 0 \forall i$$

In the **primal**, the  $x$ -player goes first and tries to set whatever value of  $x$  minimizes things (ignoring constraints). Then the  $\alpha$ -player goes and can set huge penalties anywhere a constraint is violated.

In the **dual**, the  $\alpha$ -player goes first and tries to maximize things. Then the  $x$ -player goes and can try to minimize from there.



$$\mathcal{L}(x, \alpha_1, \dots, \alpha_m) = f(x) + \sum_i^m \alpha_i g_i(x)$$

**Duality Gap:** max of min is often lower than min of max.

Primal:

$$\min_x \max_{\alpha_1, \dots, \alpha_m} \mathcal{L}(x, \alpha_1, \dots, \alpha_m) \quad s.t. \alpha_i \geq 0 \forall i$$

Dual:

$$\max_{\alpha_1, \dots, \alpha_m} \min_x \mathcal{L}(x, \alpha_1, \dots, \alpha_m) \quad s.t. \alpha_i \geq 0 \forall i$$

**Karush-Kuhn-Tucker (KKT)** Conditions tell us when the dual and primal have the same optimal solutions  $x^*, \alpha_1^*, \dots, \alpha_m^*$  for convex  $f(x)$ :

- $\nabla \mathcal{L}(x^*, \alpha_1^*, \dots, \alpha_m^*) = 0$  At a critical point of the Lagrangian
- $g(x^*) \leq 0$  The solution satisfies the primal constraints
- $\alpha_i^* \geq 0$  The solution satisfies the dual constraints
- $\alpha^* g(x^*) = 0$  Complementary slackness.



$$\begin{array}{ll} \min_x & \textcolor{red}{x^2} \\ \text{s.t.} & -10\textcolor{blue}{x} + 5 \leq 0 \end{array}$$

$$\mathcal{L}(x, \alpha) = x^2 + \alpha(-10x + 5)$$

$$x^* = \frac{1}{2} \quad \alpha^* = \frac{1}{5}$$

**Karush-Kuhn-Tucker (KKT)** Conditions tell us when the dual and primal have the same optimal solutions  $x^*, \alpha_1^*, \dots, \alpha_m^*$  for convex  $f(x)$ :

- $\frac{\delta \mathcal{L}(x^*, \alpha^*)}{dx} = 0, \frac{\delta \mathcal{L}(x^*, \alpha^*)}{d\alpha} = 0$  At a critical point of the Lagrangian
- $-10x^* + 5 \leq 0$  The solution satisfies the primal constraints
- $\alpha^* \geq 0$  The solution satisfies the dual constraints
- $\alpha^*(-10x^* + 5) = 0$  Complementary slackness.



# Constrained Optimization with Lagrange Multipliers





# Hard-Margin Support Vector Machine Objective

**Where were we?** We wanted to solve the constrained optimization problem below because it maximizes the margin while having no errors.

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$s.t. \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

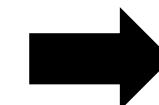
Let's start getting this into shape to apply the Langrange multipliers bit.



# Hard-Margin Support Vector Machine Objective

Need to get our inequality facing the right way. Multiply by -1 and move the right side over:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$



$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$s.t. \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

$$s.t. \quad 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0 \quad \forall i$$



# Hard-Margin Support Vector Machine Objective

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) \leq 0 \quad \forall i \end{aligned}$$

Let's write out the Lagrangian:

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \alpha_1, \dots, \alpha_n) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \alpha_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^n \alpha_i y_i b \end{aligned}$$

Want to solve the dual problem:

$$\max_{\alpha_i \geq 0} \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^n \alpha_i y_i b$$



# Hard-Margin Support Vector Machine Objective

$$\max_{\alpha_i \geq 0} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^n \alpha_i y_i b$$

Take derivatives and set equal to zero:

$$\frac{\delta \mathcal{L}}{\delta \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

$$\frac{\delta \mathcal{L}}{\delta b} = - \sum_{i=1}^n \alpha_i y_i = 0$$



# Hard-Margin Support Vector Machine Objective

$$\max_{\alpha_i \geq 0} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^n \alpha_i y_i b$$

Apply what we learned from the derivatives:

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

$$\sum_{i=1}^n \alpha_i y_i = 0$$

$$\max_{\alpha_i \geq 0} \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right)^T \left( \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right) + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \left( \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right)^T \mathbf{x}_i - b \sum_{i=1}^n \alpha_i y_i$$

$$\max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0$$



# Hard-Margin Support Vector Machine Objective

$$\max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i^n \sum_j^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0$$

This is a quadratic function with linear constraints. Many (non-trivial) algorithms to solve these sorts of things - called Quadratic Program Solvers (QP-Solvers). I can put this into a QP solver and get back out optimal  $\alpha_i^*$ .

Further, **the KKT conditions hold** (not shown here) such that solving this yields the exact same optimal solution as the original formalization!





Please accept this majestic photo of my friend's cat Simba for sticking with me





# Recap of what we've been doing so far

**Motivation:** We wanted to find a classifier that maximizes the margin between classes.

**Derivation:** We derived the following constrained optimization problem that does that:

Hard-Margin  
SVM Primal:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad s.t. \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

We said that is a Quadratic Programming problem and totally solvable. Then I ran us through a lot of math to arrive at the dual formalization below which solves the same problem (KKT's hold).

Hard-Margin  
SVM Dual:

$$\max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i^n \sum_j^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0$$

Which is also a Quadratic Programming problem (in  $\alpha$ 's this time rather than  $\mathbf{w}$ ) and is totally solvable to global optimality..... .... **so why did I do this to us all?**



# Why have we done all this?

Hard math builds character.





# Examining the Dual Solution

**Real Reason:** We can examine this dual formalization to observe some very, very interesting things that are not at all apparent in the primal formalization!

**Observation 1:** For optimal  $\alpha_i^*$  found by maximizing the dual problem, the optimal weight vector turns out to be a weighted sum of input vectors times their labels:

We found this back when  
we took derivatives.


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

Arises from KKT  
conditions, next slide.


$$b^* = \frac{1}{\sum_{i=1}^n \alpha_i^*} \sum_{i=1}^n \alpha_i^* (y_i - \mathbf{w}^{*T} \mathbf{x}_i)$$



# Examining the Dual Solution

**Observation 2:** Only certain training examples influence the result. Many of the  $\alpha_i$ 's are zero.

At optimality, the KKT conditions require that:

$$\alpha_i \geq 0 \quad \forall i \quad \leftarrow \text{Satisfies dual constraints}$$

$$y_i(w^{*T}x_i + b) \geq 1, \quad \forall i \quad \leftarrow \text{Satisfies primal constraints}$$

$$\rightarrow \alpha_i^*(y_i(w^{*T}x_i + b) - 1) = 0, \quad \forall i \quad \leftarrow \text{Complementary slackness}$$



What does the complementary slackness condition

$\alpha_i^*(y_i(w^{*T}x_i + b) - 1) = 0, \forall i$  imply about the values of the alpha variables?

(Think about when this equation is true for a single training point i)



# Examining the Dual Solution

**Observation 2:** Only certain training examples influence the result. Many of the  $\alpha_i$ 's are zero.

At optimality, the KKT conditions require that:

$$\alpha_i \geq 0 \quad \forall i \quad \leftarrow \text{Satisfies dual constraints}$$

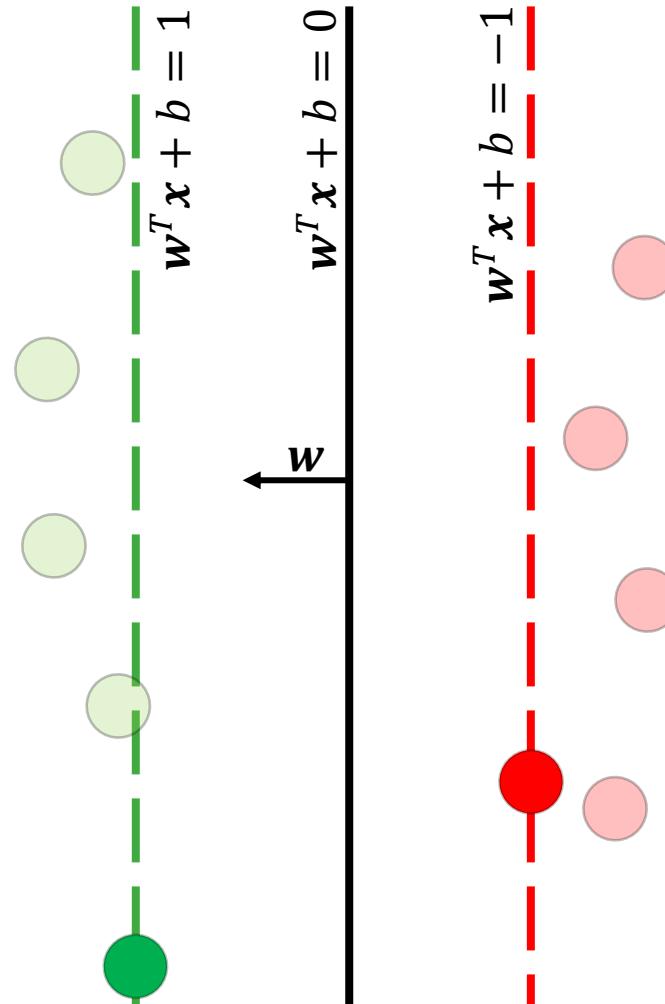
$$y_i(w^{*T}x_i + b) \geq 1, \quad \forall i \quad \leftarrow \text{Satisfies primal constraints}$$

$$\rightarrow \alpha_i^*(y_i(w^{*T}x_i + b) - 1) = 0, \quad \forall i \quad \leftarrow \text{Complementary slackness}$$

**$\alpha_i^*$  is only non-zero if  $y_i(w^T x_i + b) = 1!$**



# Examining the Dual Solution



**$\alpha_i^*$  is only non-zero if  $y_i(w^T x_i + b) = 1$ !**



**$\alpha_i^*$  is only non-zero if  $x_i$  is on the margin edges (green or red line)**

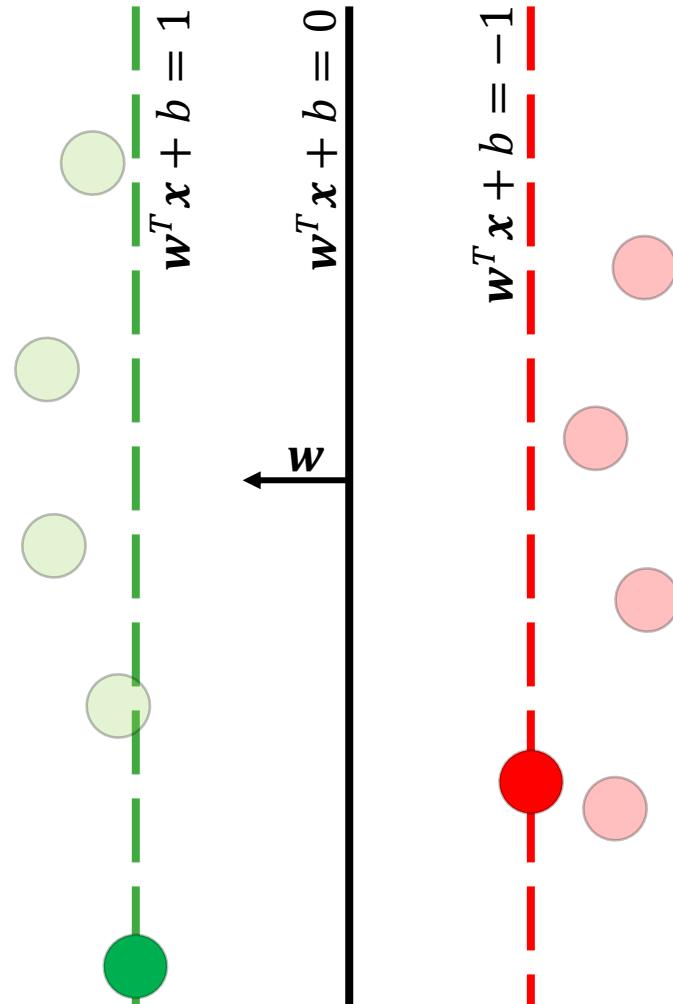


$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

**Only the points on the margin matter to the classifier weights.**



# Examining the Dual Solution

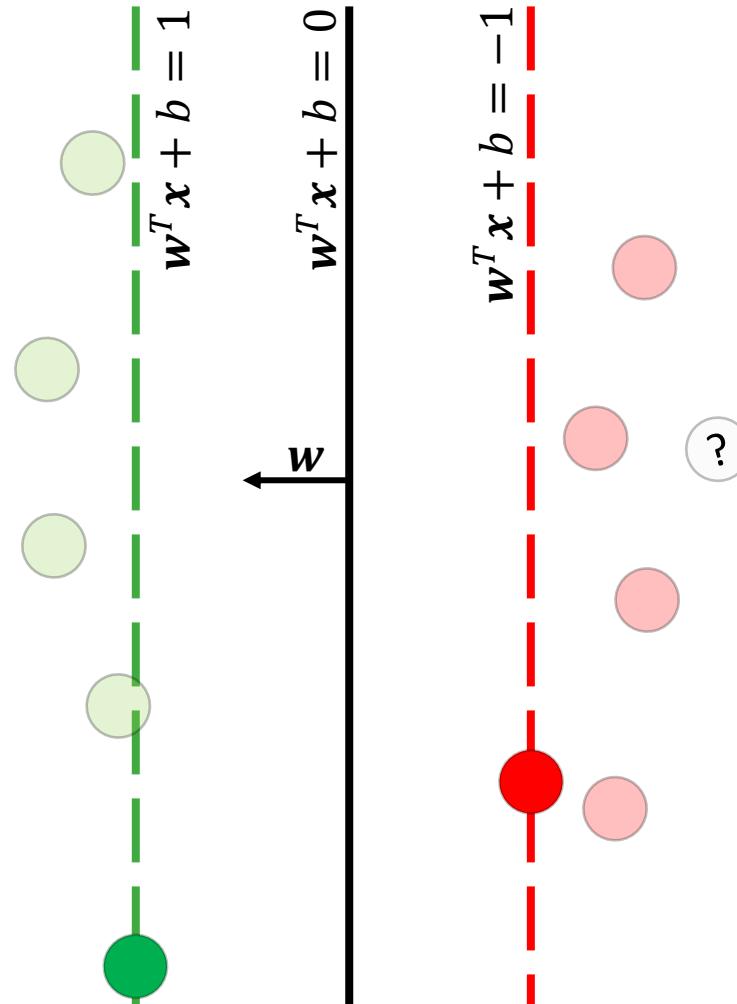


Call these points with non-zero  
alpha **support vectors**.





# Examining the Dual Solution



## How is a new point classified?

Our decision rule is:

$$\begin{aligned} w^T x + b > 0 &\Rightarrow \text{predict } 1 \\ w^T x + b < 0 &\Rightarrow \text{predict } -1 \end{aligned}$$

Plugging in our optimal weight:

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i \Rightarrow w^T x + b = b + \sum_{i=1}^n \alpha_i^* y_i x_i^T x$$

The prediction is based on dot product with a sparse set of training examples.



<https://jgreitemann.github.io/svm-demo>



## What you need to know from today's lecture:

**Support Vector Machines (SVMs)** find the decision boundary that maximizes the margin between the positive and negative points.

The weight vector corresponding to this max-margin classifier can be found by solving a constrained optimization problem with a QP solver.

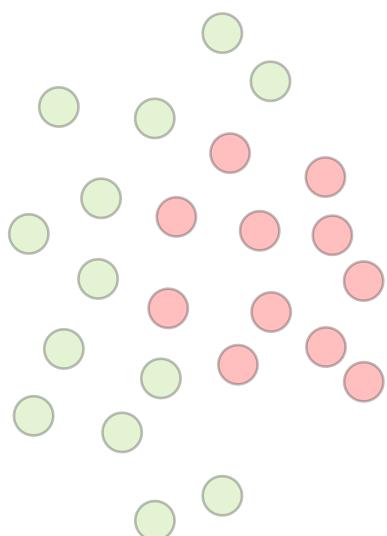
By going through the effort of examining the Lagrangian dual of the SVM objective, we found that the **decision boundary is a function of a small, sparse set of examples** that lie on the classifier's margin.

What we derived today is referred to as a **hard-margin SVM** because it does not tolerate any errors. If we give it a non-linearly separable dataset, the QP solver will report no feasible solution.



## Teaser for Next Time: What about the non-linearly separable case?

**Teaser 1:** We can modify our objective function to try to maximize margin while making as few mistakes as possible. Will need to introduce some notion of the tradeoff between the two.



Maximize margin



Minimize violations  
of the margin

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \epsilon_i$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i \quad \forall i \\ \epsilon_i \geq 0$$

Allow some  
violations

We'll talk about how this works out next time.



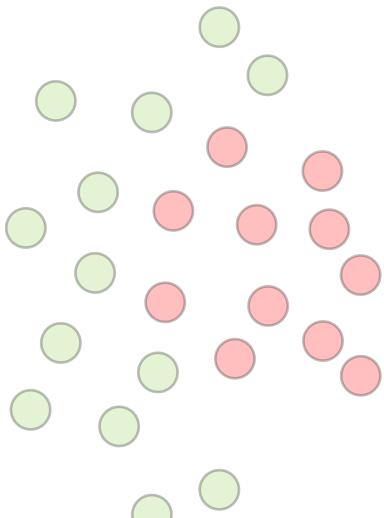
## Teaser for Next Time: What about the non-linearly separable case?

**Teaser 2:** Apply the non-linear basis function idea again and just replace each point with some  $\Phi(x)$  that is non-linear. Then fit a linear classifier in that higher-dimensional space.

Train and find  $\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \Phi(\mathbf{x}_i)$

Predict based on  $\mathbf{w}^T \Phi(\mathbf{x}) + b$

$$\longrightarrow b + \sum_{i=1}^n \alpha_i^* y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x})$$



This is expensive and we have to deal with those higher dimensions. What if we just defined a kernel that compute non-linear similarity without actually representing these high dimensional vectors.

$$b + \sum_{i=1}^n \alpha_i^* y_i \mathcal{K}(\mathbf{x}_i, \mathbf{x})$$

$$\max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i^n \sum_j^n \alpha_i \alpha_j y_i y_j \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$$

SVMs operate on dot products between examples. Can “pretend” we are in a higher dimension – call this the **kernel trick**. More next time.



Next time



**Next Time:** How do we deal with non-linearly separable data? Kernels + Soft-Margin SVM.