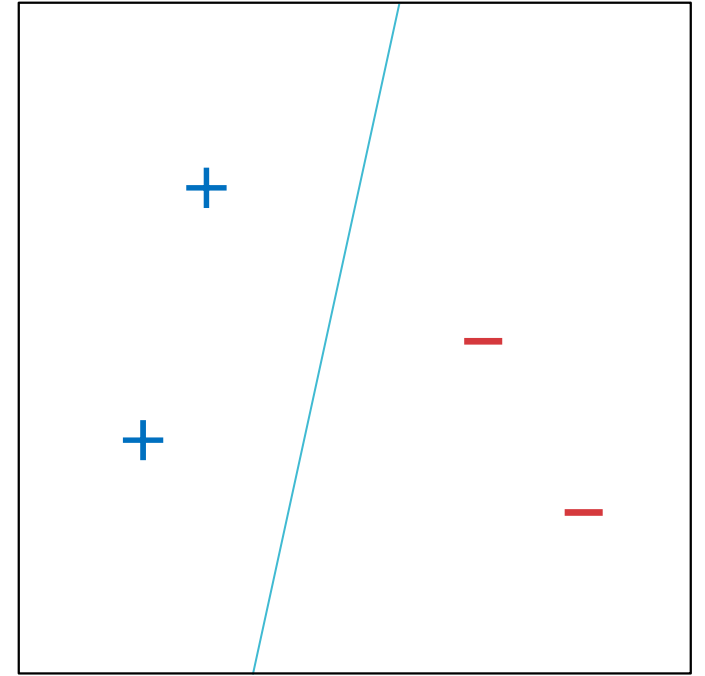
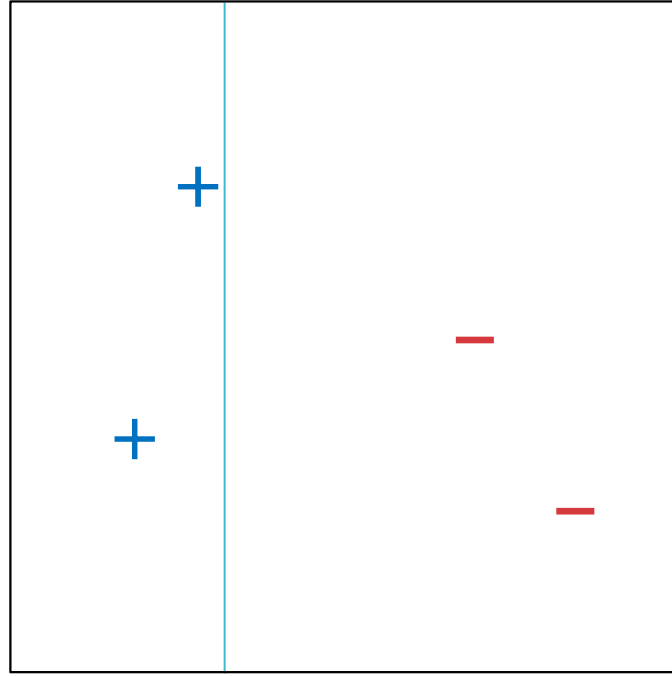
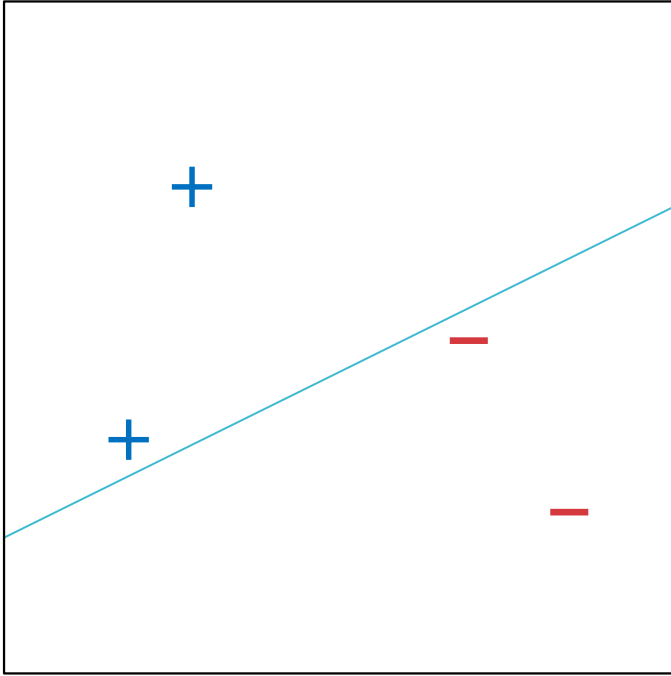


10-315: Introduction to Machine Learning

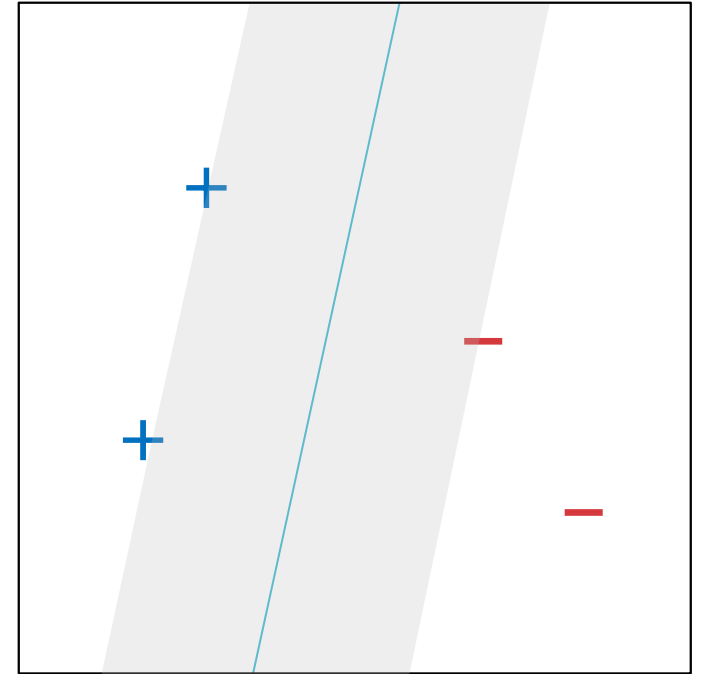
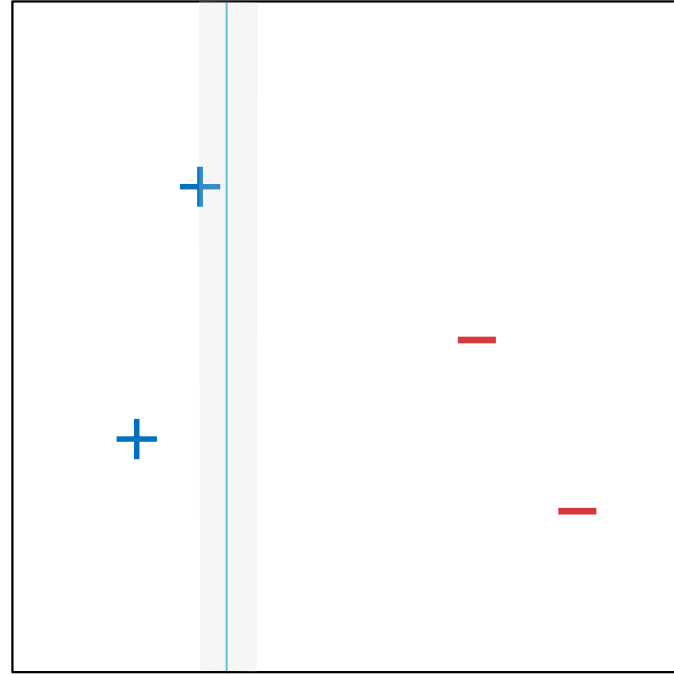
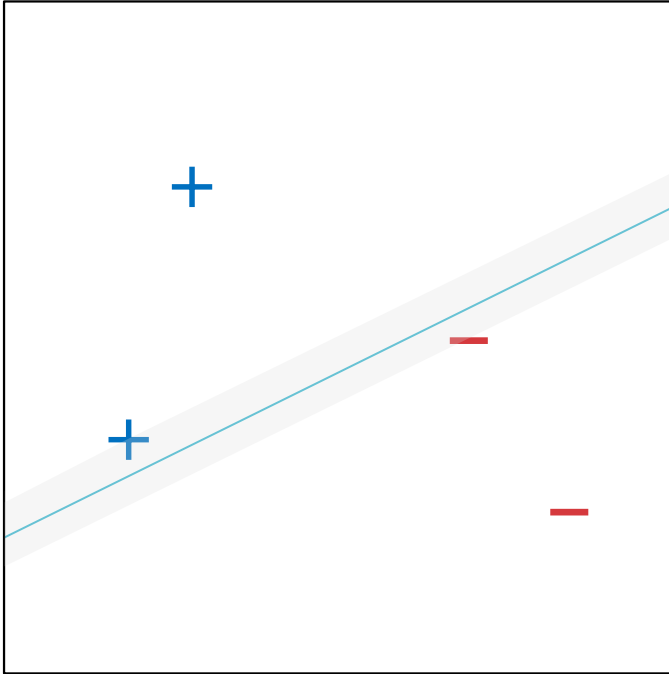
Lecture 15— Support Vector Machines

Henry Chai

3/21/22



Which linear separator is best?



Which linear separator is best?

Maximal Margin Linear Separators

- The margin of a linear separator is the distance between it and the nearest training data point
- Questions:
 1. How can we efficiently find a maximal-margin linear separator?
 2. Why are linear separators with larger margins better?
 3. What can we do if the data is not linearly separable?

Hyperplanes

- For linear models, decision boundaries are D -dimensional *hyperplanes* defined by a weight vector, $[b, \mathbf{w}]$

$$\mathbf{w}^T \mathbf{x} + b = 0$$

- Problem: there are infinitely many weight vectors that describe the same hyperplane

- $x_1 + 2x_2 + 2 = 0$ is the same line as

$2x_1 + 4x_2 + 4 = 0$, which is the same line as

$$1000000x_1 + 2000000x_2 + 2000000 = 0$$

Solution: normalize the weight vectors
w.r.t. training data set

Normalizing Hyperplanes

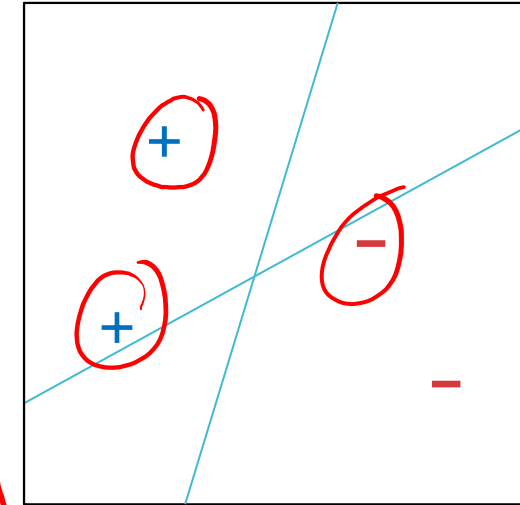
- Given a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ where $y \in \{-1, +1\}$,
 $\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$ is a valid linear separator if
$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) > 0 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$$

For SVMs, only consider
$$\mathcal{H} = \{ \hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x}^{(i)} + b) : \min_{(\mathbf{x}^{(i)}, y^{(i)})} y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1 \}$$

if $\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$, then $\exists \rho$ s.t.
$$\hat{y} = \text{sign}\left(\frac{\mathbf{w}^T}{\rho} \mathbf{x}^{(i)} + \frac{b}{\rho}\right) \in \mathcal{H}$$

Normalizing Hyperplanes: Example

b	w_1	w_2	
-0.2	-0.6	1	$\notin \mathcal{H}$
-0.4	-1.2	2	$\notin \mathcal{H}$
-2	-6	10	$\notin \mathcal{H}$
-10	-30	50	$\in \mathcal{H}$
0.2	-0.6	0.2	$\notin \mathcal{H}$
0.1	-0.3	0.1	$\notin \mathcal{H}$
1	-3	1	$\notin \mathcal{H}$
2	-6	2	$\in \mathcal{H}$



x_1	x_2	y	$y(w^T x + b)$
0.2	0.4	+1	1.6
0.3	0.8	+1	1.8
0.7	0.6	-1	1
0.8	0.3	-1	2.2

$$\begin{aligned}
 & -1(0.7(-6) + 0.6(10) - 2) \\
 & -1(-4.2 + 6 - 2) = 0.2
 \end{aligned}$$

Computing the Margin

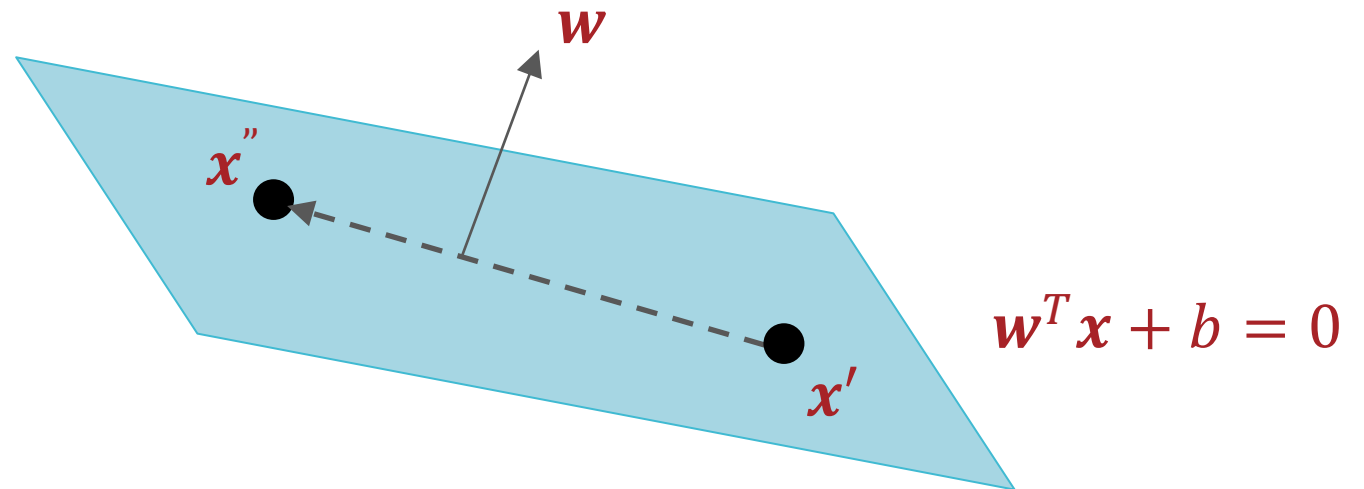
- Claim: \mathbf{w} is orthogonal to the hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ (the decision boundary)
- A vector is orthogonal to a hyperplane if it is orthogonal to every vector in that hyperplane
- Vectors α and β are orthogonal if $\alpha^T \beta = 0$
- Proof:

Let x' and x'' be two arbitrary points on $\mathbf{w}^T \mathbf{x} = -b \leftarrow \mathbf{w}^T \mathbf{x} + b = 0 \Rightarrow x' - x''$ is a vector on $\mathbf{w}^T \mathbf{x} + b = 0$

$$\begin{aligned}\mathbf{w}^T (x' - x'') &= \mathbf{w}^T x' - \mathbf{w}^T x'' \\ &= (-b) - (-b) = 0\end{aligned}$$

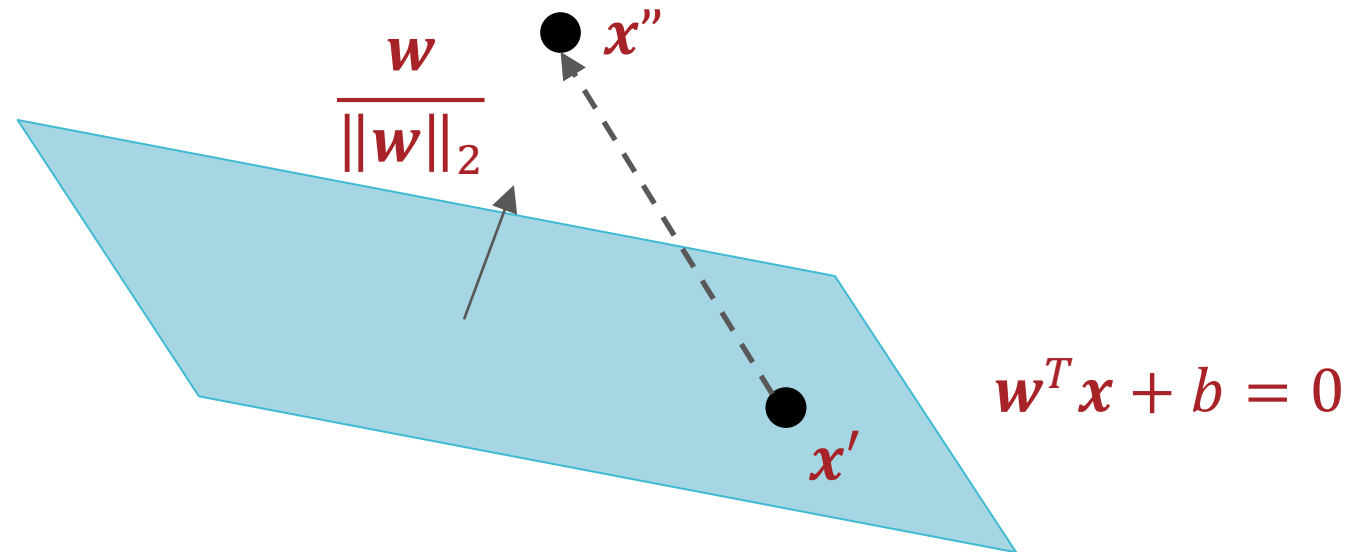
Computing the Margin

- Claim: \mathbf{w} is orthogonal to the hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ (the decision boundary)
- A vector is orthogonal to a hyperplane if it is orthogonal to every vector in that hyperplane
- Vectors α and β are orthogonal if $\alpha^T \beta = 0$



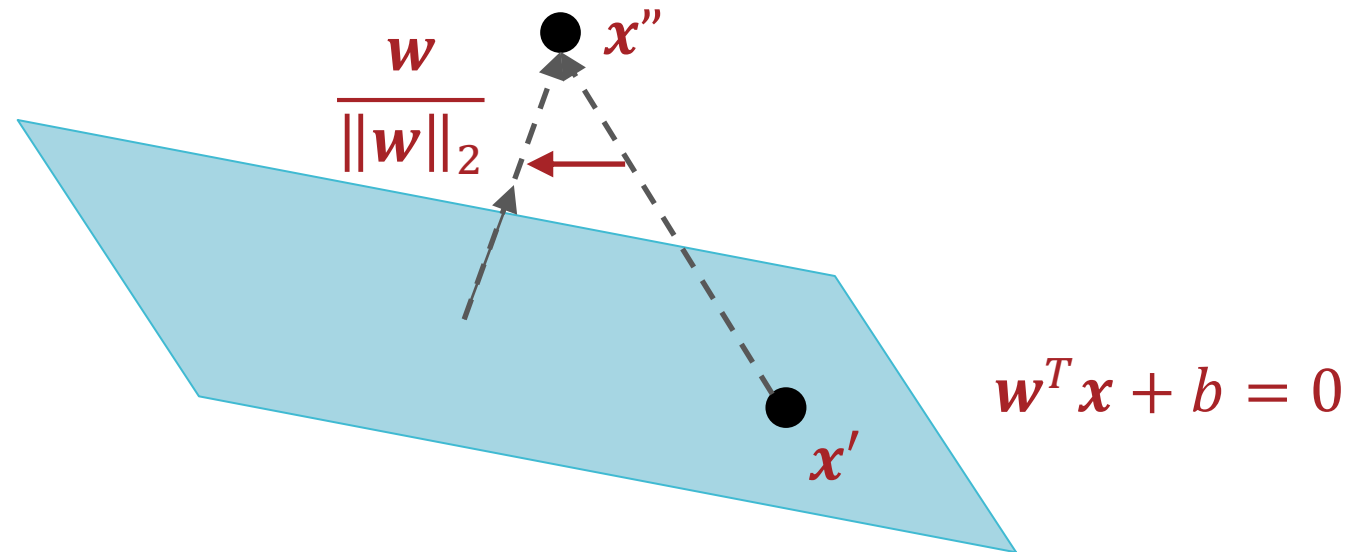
Computing the Margin

- Let \mathbf{x}' be an arbitrary point on the hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ and let \mathbf{x}'' be an arbitrary point
- The distance between \mathbf{x}'' and $\mathbf{w}^T \mathbf{x} + b = 0$ is equal to the magnitude of the projection of $\mathbf{x}'' - \mathbf{x}'$ onto $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, the unit vector orthogonal to the hyperplane



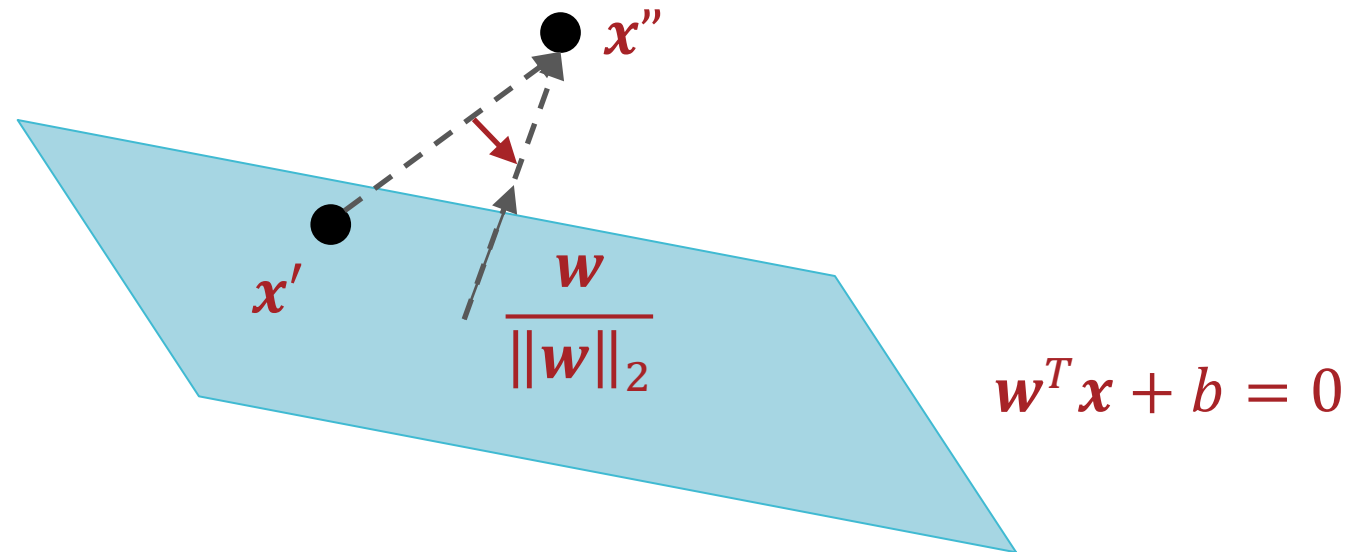
Computing the Margin

- Let \mathbf{x}' be an arbitrary point on the hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ and let \mathbf{x}'' be an arbitrary point
- The distance between \mathbf{x}'' and $\mathbf{w}^T \mathbf{x} + b = 0$ is equal to the magnitude of the projection of $\mathbf{x}'' - \mathbf{x}'$ onto $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, the unit vector orthogonal to the hyperplane



Computing the Margin

- Let \mathbf{x}' be an arbitrary point on the hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ and let \mathbf{x}'' be an arbitrary point
- The distance between \mathbf{x}'' and $\mathbf{w}^T \mathbf{x} + b = 0$ is equal to the magnitude of the projection of $\mathbf{x}'' - \mathbf{x}'$ onto $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, the unit vector orthogonal to the hyperplane



Computing the Margin

- Let \mathbf{x}' be an arbitrary point on the hyperplane $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$ and let \mathbf{x}'' be an arbitrary point
- The distance between \mathbf{x}'' and $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$ is equal to the magnitude of the projection of $\mathbf{x}'' - \mathbf{x}'$ onto $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, the unit vector orthogonal to the hyperplane

$$\begin{aligned} d(\mathbf{x}'', h) &= \left| \frac{\mathbf{w}^T (\mathbf{x}'' - \mathbf{x}')}{\|\mathbf{w}\|_2} \right| = \left| \frac{\mathbf{w}^T \mathbf{x}'' - \mathbf{w}^T \mathbf{x}'}{\|\mathbf{w}\|_2} \right| \\ &= \left| \frac{\mathbf{w}^T \mathbf{x}'' - (-b)}{\|\mathbf{w}\|_2} \right| = \left| \frac{\mathbf{w}^T \mathbf{x}'' + b}{\|\mathbf{w}\|_2} \right| \end{aligned}$$

Computing the Margin

- The margin of a linear separator is the distance between it and the nearest training data point

$$\begin{aligned} \min_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} d(x^{(i)}, h) &= \min_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left| \frac{w^T x^{(i)} + b}{\|w\|_2} \right| \\ &= \frac{1}{\|w\|_2} \min_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \underset{\substack{\uparrow \\ y^{(i)}}}{|w^T x^{(i)} + b|} \\ &= \frac{1}{\|w\|_2} \min_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \underbrace{y^{(i)}(w^T x^{(i)} + b)}_{(1)} \\ &= \frac{1}{\|w\|_2} \end{aligned}$$

Maximizing the Margin

$$\begin{aligned} &\text{maximize} \quad \frac{1}{\|w\|_2} = \frac{1}{\sqrt{w^T w}} \\ &\text{s.t.} \quad \min_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)}(w^T x^{(i)} + b) = 1 \end{aligned}$$

$$\begin{aligned} &\text{minimize} \quad w^T w \\ &\text{s.t.} \quad \min_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)}(w^T x^{(i)} + b) = 1 \end{aligned}$$

$$\begin{aligned} &\text{minimize} \quad w^T w \\ &\text{s.t.} \quad y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad \forall (x^{(i)}, y^{(i)}) \in \mathcal{D} \end{aligned}$$

Maximizing the Margin



$$\text{minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{subject to } y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$$

- If $[\hat{b}, \hat{\mathbf{w}}]$ is the optimal solution, then \exists at least one training data point $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$ s.t. $y^{(i)} (\hat{\mathbf{w}}^T \mathbf{x}^{(i)} + \hat{b}) = 1$
 - All training data points $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$ where $y^{(i)} (\hat{\mathbf{w}}^T \mathbf{x}^{(i)} + \hat{b}) = 1$ are known as support vectors
- Converting the non-linear constraint (involving the **min**) to N linear constraints means we can use quadratic programming (QP) to solve this problem in $O(D^3)$ time

Recipe for SVMs

- Define a model and model parameters
 - Assume a linear decision boundary (with normalized weights)

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$$

- Parameters: $\mathbf{w} = [w_1, \dots, w_D]$ and b
- Write down an objective function (with constraints)
minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w}$
subject to $y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$
- Optimize the objective w.r.t. the model parameters
 - Solve using quadratic programming

Why Maximal Margins?

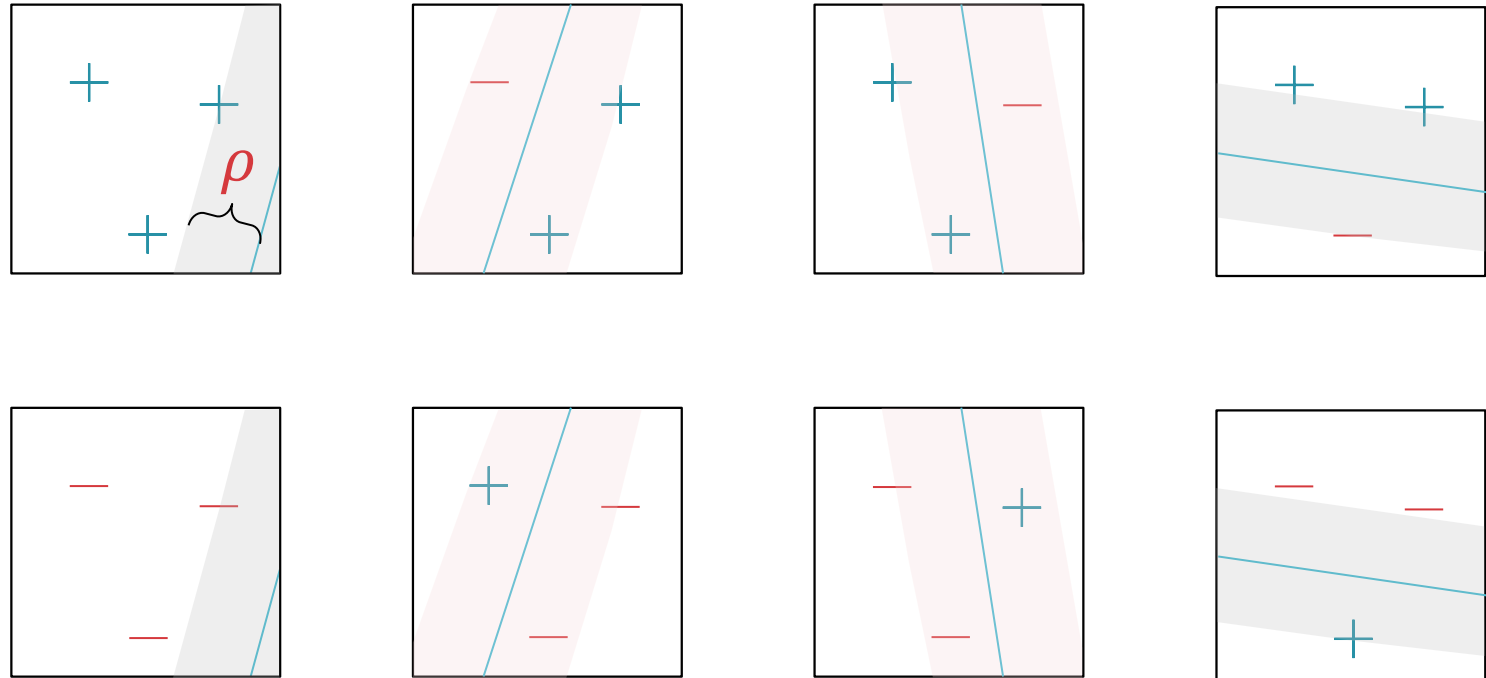
- Consider three binary data points in a **bounded** 2-D space
- Let $\mathcal{H} = \{\text{all linear separators}\}$ and
 $\mathcal{H}_\rho = \{\text{all linear separators with minimum margin } \rho\}$

Why Maximal Margins?

- Consider three binary data points in a **bounded** 2-D space
- \mathcal{H} = {all linear separators} can always correctly classify any three (non-collinear) data points in this space

Why Maximal Margins?

- Consider three binary data points in a **bounded** 2-D space
- $\mathcal{H}_\rho = \{\text{all linear separators with minimum margin } \rho\}$ cannot always correctly classify three non-collinear data points



Summary Thus Far

- The margin of a linear separator is the distance between it and the nearest training data point
- Questions:
 1. How can we efficiently find a maximal-margin linear separator? By solving a constrained quadratic optimization problem using quadratic programming
 2. Why are linear separators with larger margins better? They're simpler *waves hands*
 3. What can we do if the data is not linearly separable? Next!

Linearly Inseparable Data

- What can we do if the data is not linearly separable?

→ 1. Accept some non-zero training error
today

→ 2. Apply some feature transformation
later that makes your data linearly
separable

SVMs

minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

subject to $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$

- When \mathcal{D} is not linearly separable, there are no feasible solutions to this optimization problem

Hard-margin SVMs

minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

subject to $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$

- When \mathcal{D} is not linearly separable, there are no feasible solutions to this optimization problem

Soft-margin SVMs

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi^{(i)} \\ &\text{subject to} \quad y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)} \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \\ &\quad \quad \quad \xi^{(i)} \geq 0 \quad \quad \quad \forall i \in \{1, \dots, N\} \end{aligned}$$

Soft-margin SVMs

- minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w} + c \sum_{i=1}^N \xi^{(i)}$ $\downarrow \geq 0$
- subject to $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)} \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$
- $\xi^{(i)} \geq 0 \quad \forall i \in \{1, \dots, N\}$
- $\xi^{(i)}$ is the “soft” error on the i^{th} training data point
 - If $\xi^{(i)} > 1$, then $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) < 0 \Rightarrow (\mathbf{x}^{(i)}, y^{(i)})$ is incorrectly classified
 - If $0 < \xi^{(i)} < 1$, then $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) > 0 \Rightarrow (\mathbf{x}^{(i)}, y^{(i)})$ is correctly classified but inside the margin
 - $\sum_{i=1}^N \xi^{(i)}$ is the “soft” training error

Soft-margin SVMs

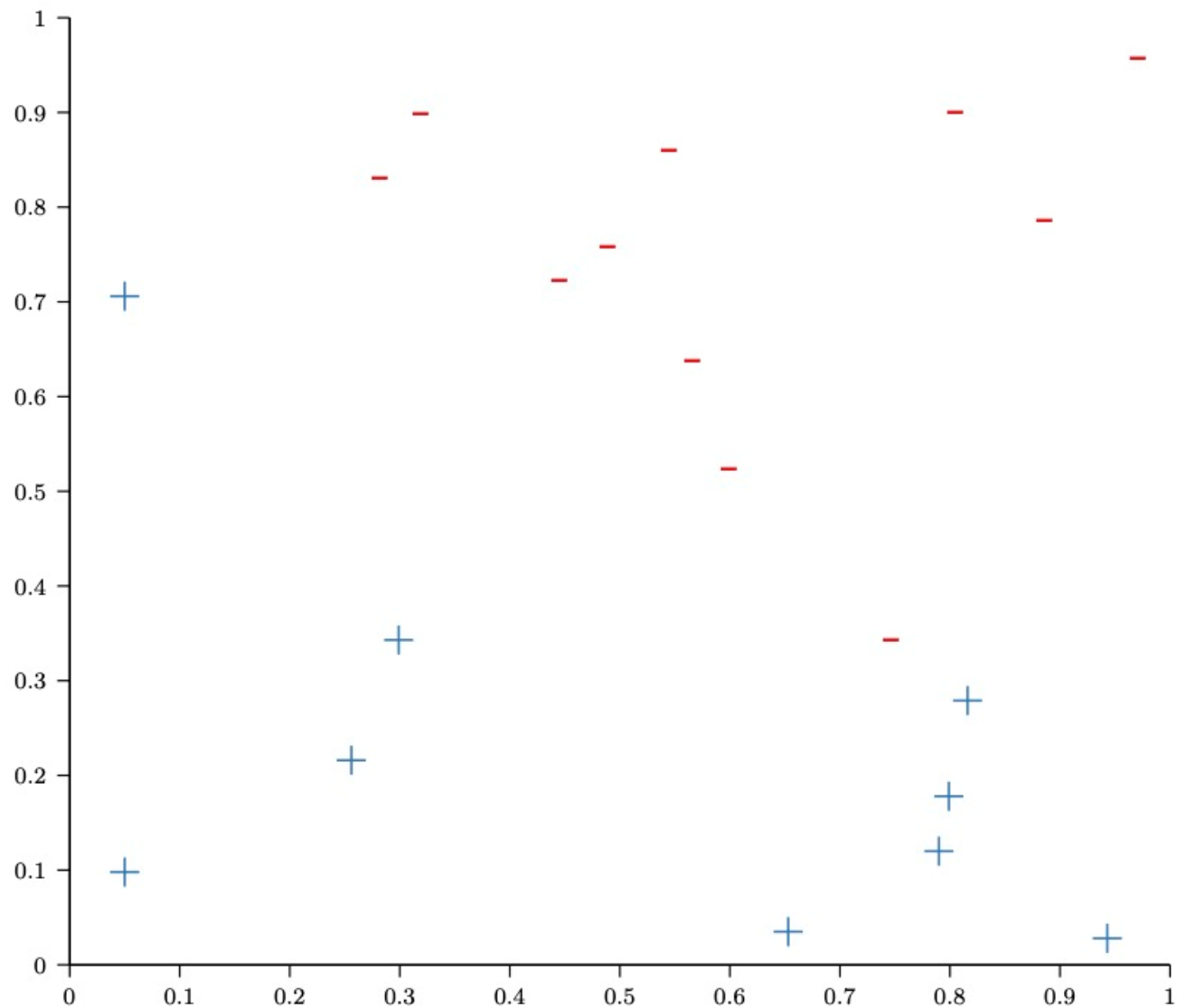
$$\begin{aligned} &\text{minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi^{(i)} \\ &\text{subject to } \underbrace{y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)}_{\xi^{(i)} \geq 0} \geq 1 - \xi^{(i)} \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \end{aligned}$$

$\checkmark \quad \neq 0$

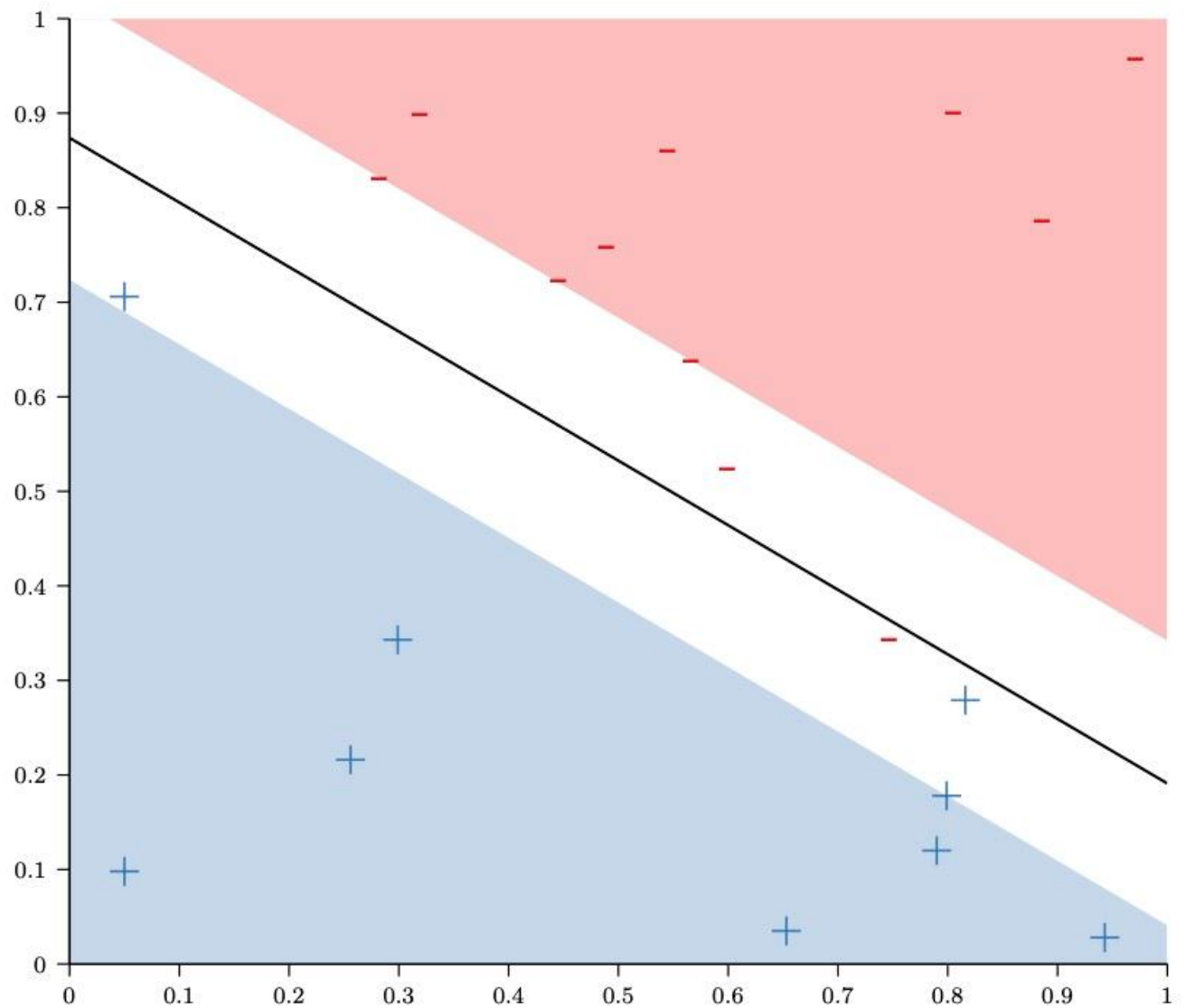
$\forall i \in \{1, \dots, N\}$

- Still solvable using quadratic programming
 - All training data points $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$ where $\underbrace{y^{(i)}(\hat{\mathbf{w}}^T \mathbf{x}^{(i)} + \hat{b}) \leq 1}$ are known as *support vectors*
- $=$ for hard-margin

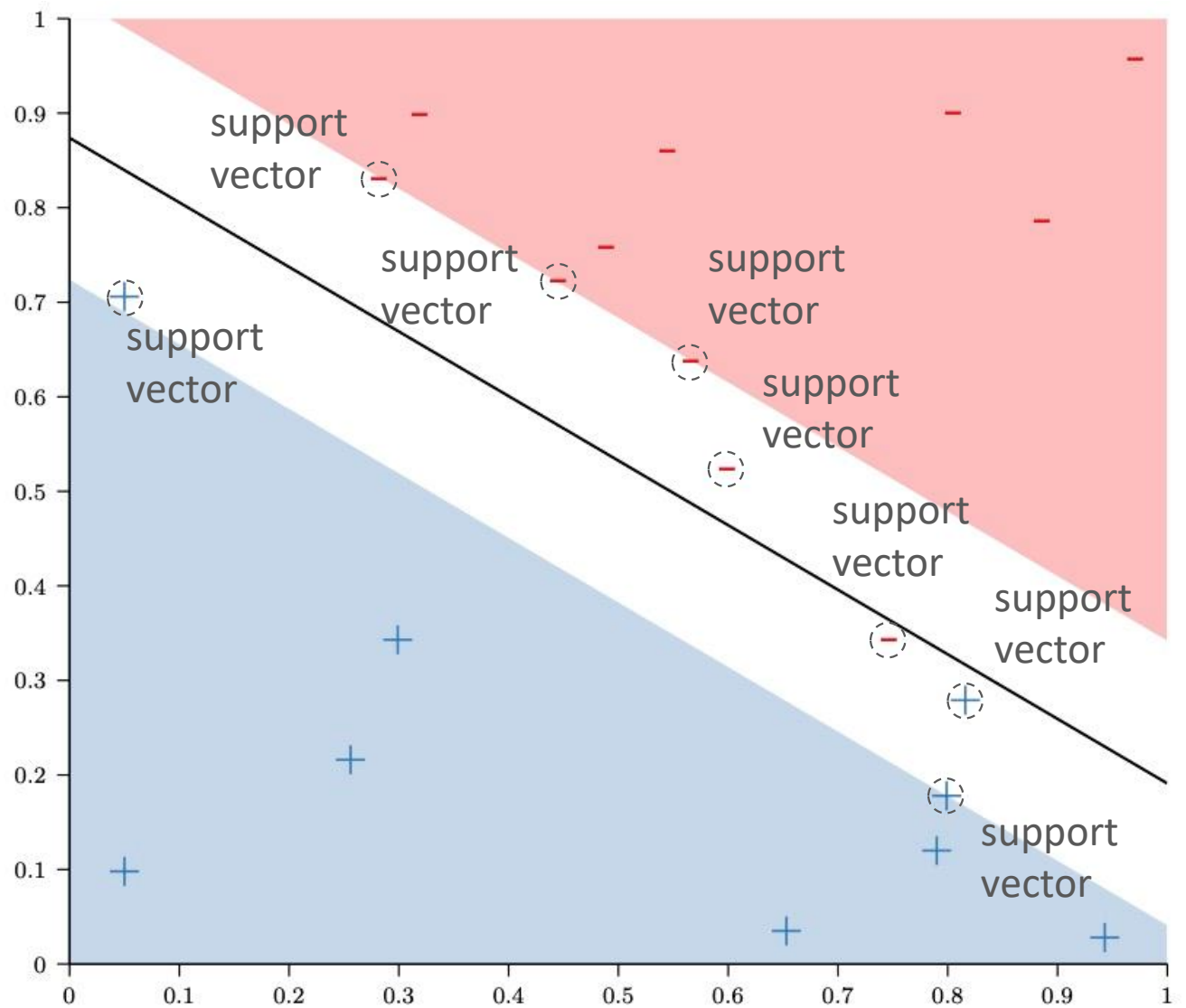
Interpreting $\xi^{(i)}$



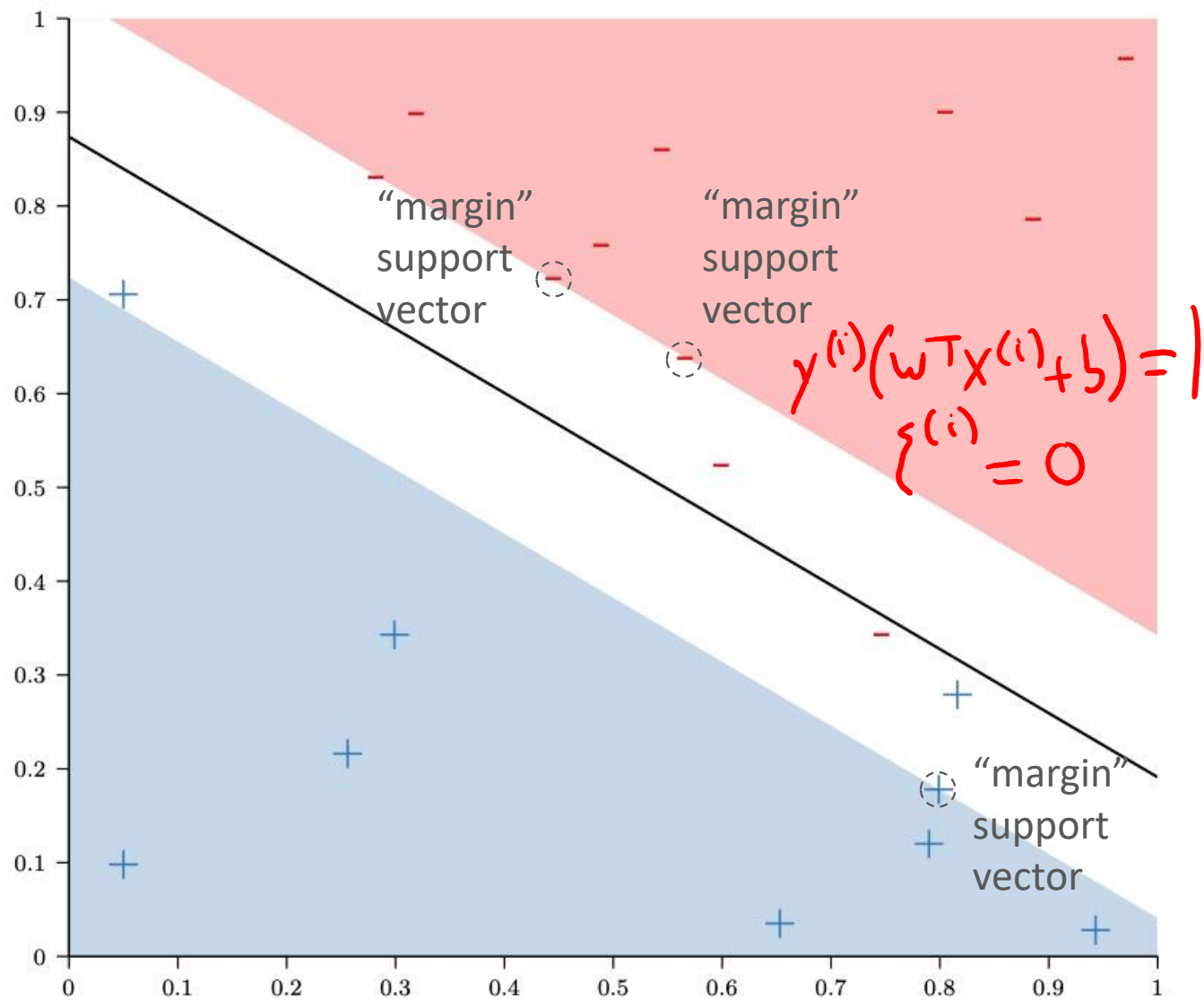
Interpreting $\xi^{(i)}$



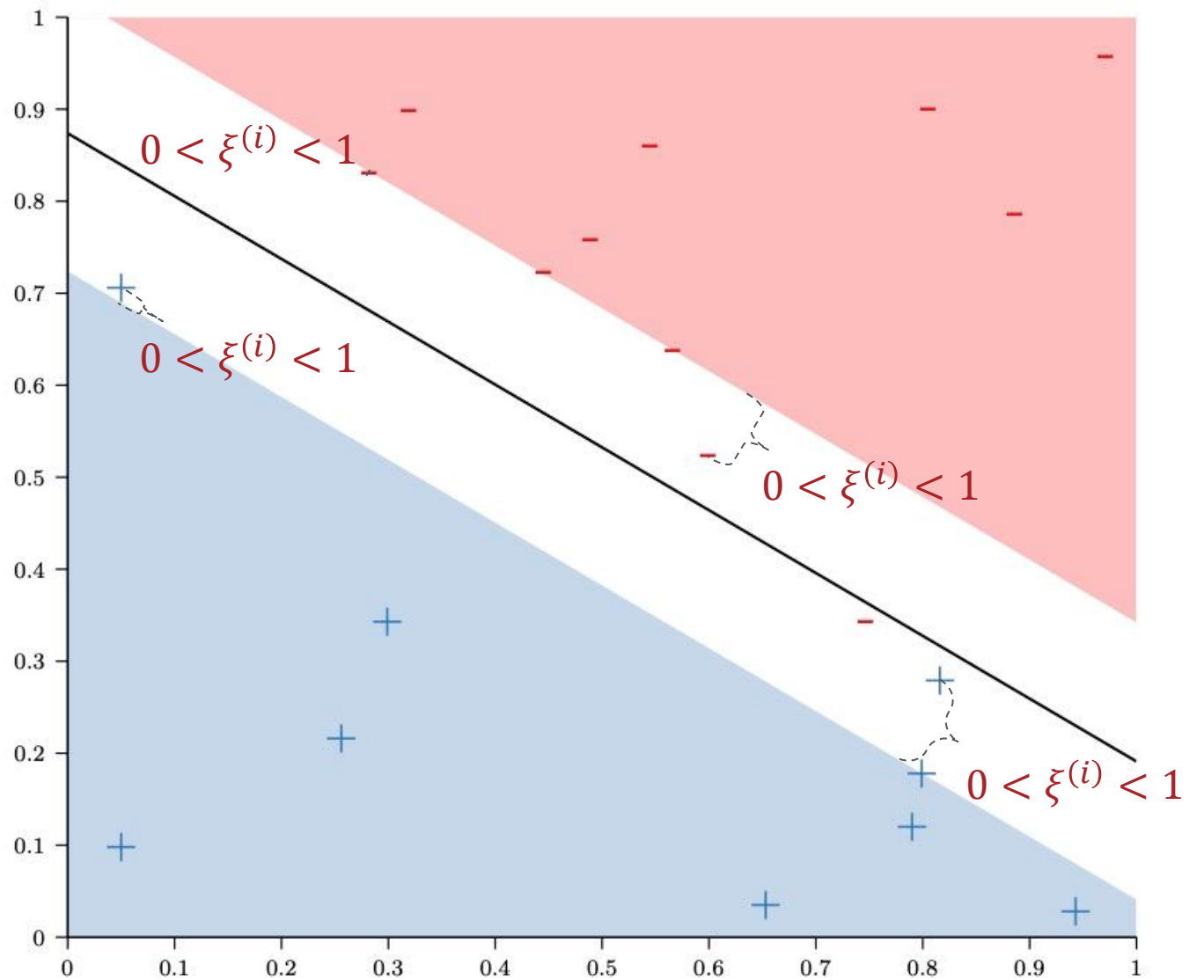
Interpreting $\xi^{(i)}$



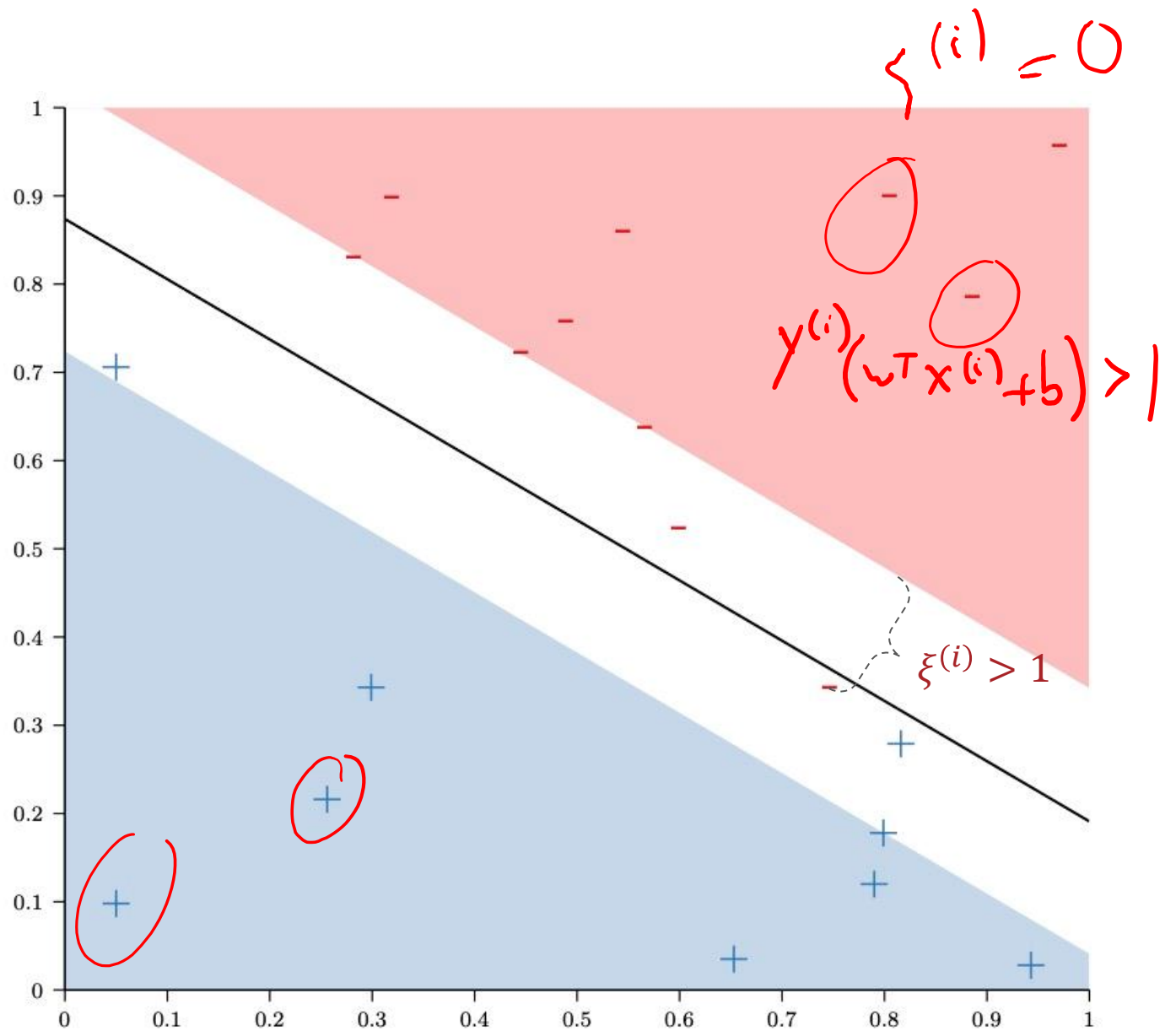
Interpreting $\xi^{(i)}$

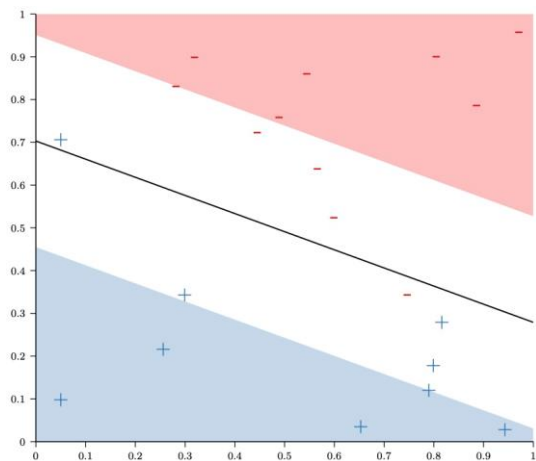


Interpreting $\xi^{(i)}$

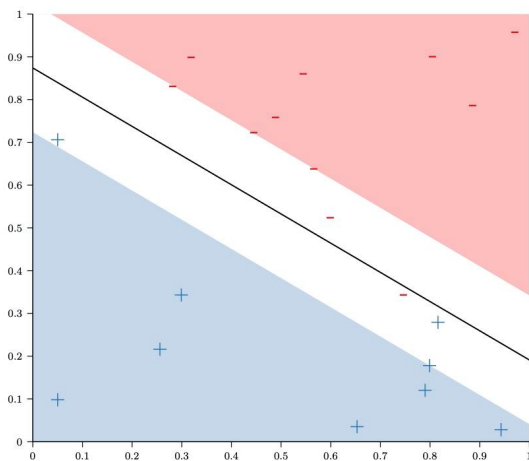


Interpreting $\xi^{(i)}$

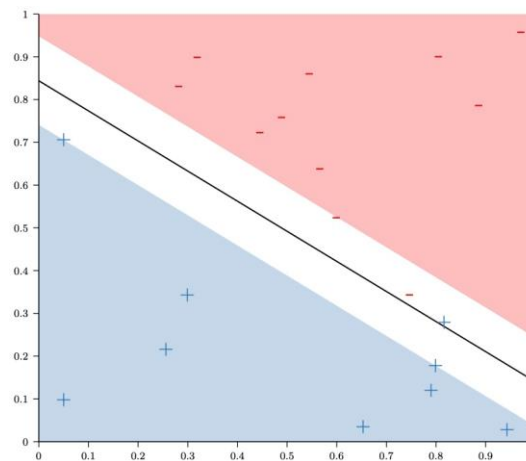




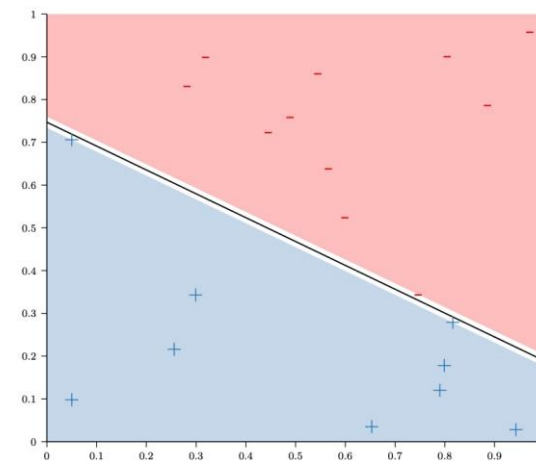
Smaller C



$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi^{(i)}$$



Larger C



Hard Margin

Setting C

C is a tradeoff parameter (much like the tradeoff parameter in regularization)