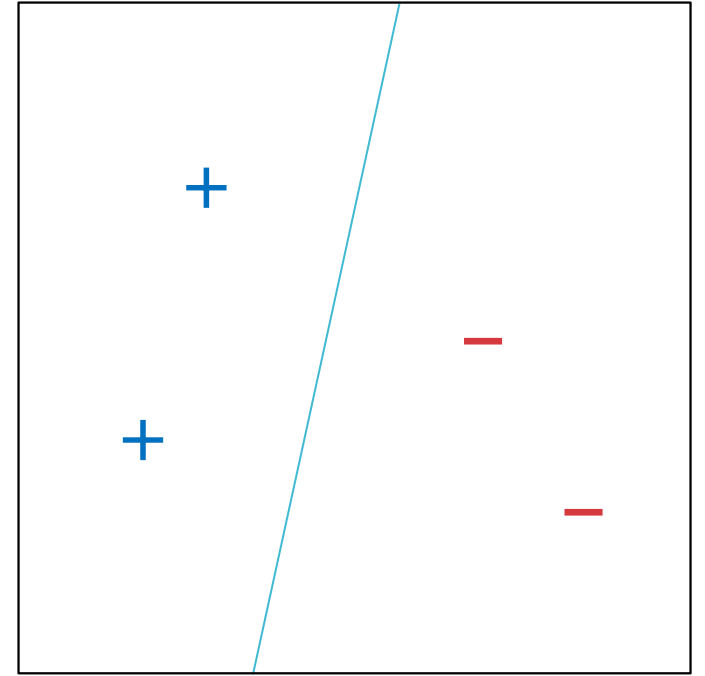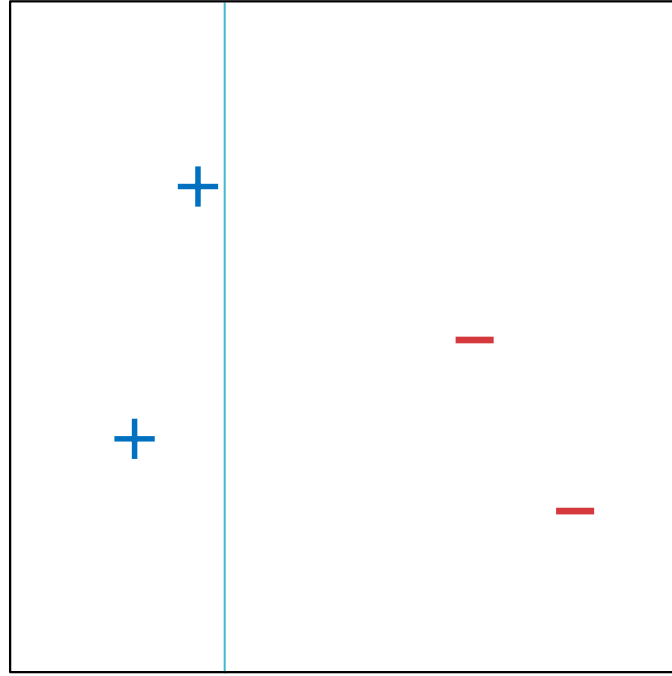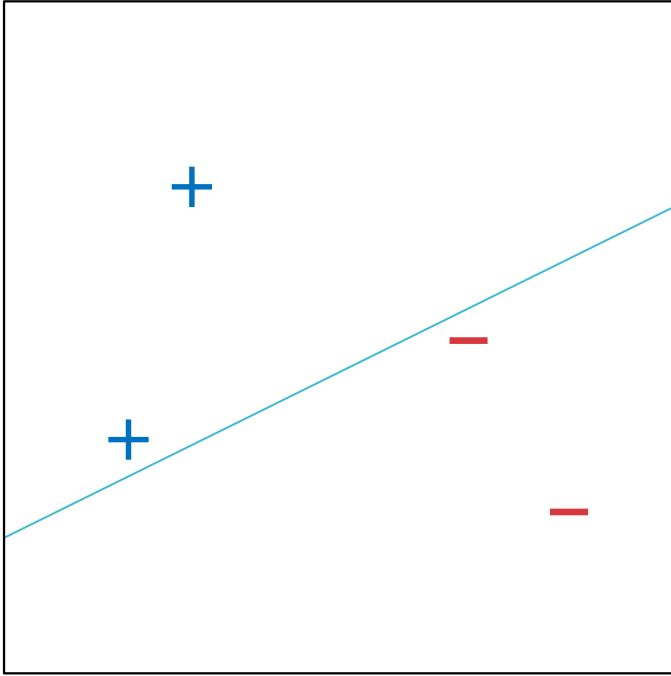# 10-315: Introduction to Machine Learning Lecture 15– Support Vector Machines
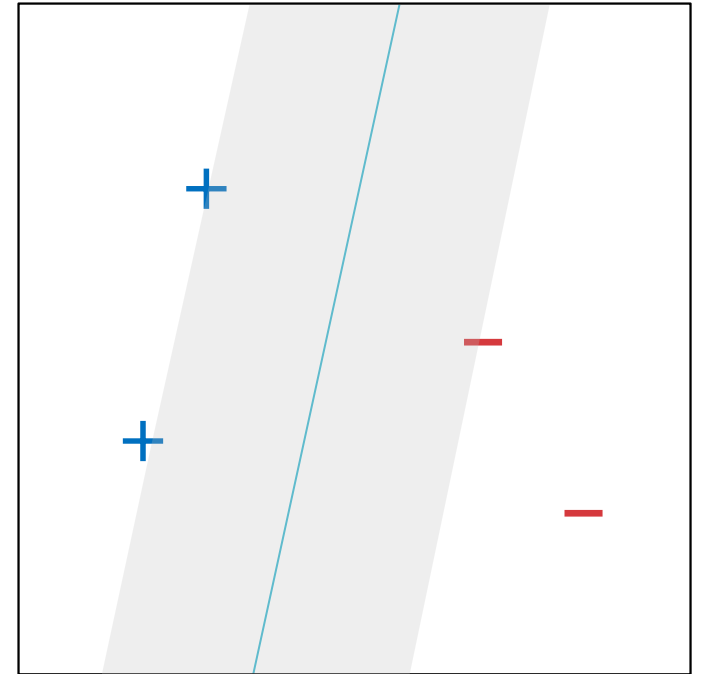
Henry Chai

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, $[b, \boldsymbol{w}]$

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

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

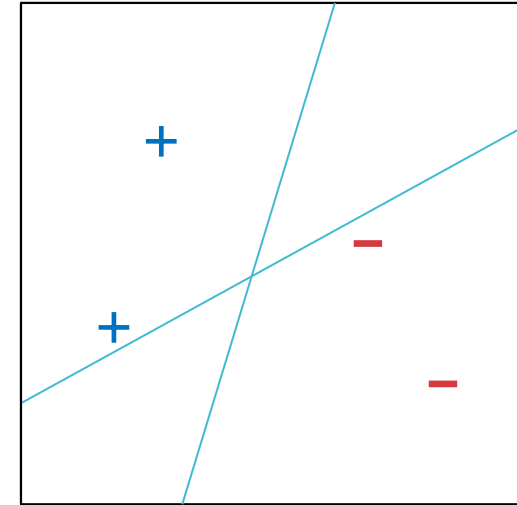- Solution: normalize weight vectors *w.r.t. the training data*

# Normalizing Hyperplanes

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

$$y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) > 0 \; \forall \left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$$

- For SVMs, we're going to consider linear separators in the set

$$\mathcal{H} = \left\{\hat{y} = \text{sign}(\boldsymbol{w}^T\boldsymbol{x} + b) : \min_{\left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}} y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) = 1\right\}$$

- If $\hat{y} = \text{sign}(\boldsymbol{w}^T\boldsymbol{x} + b)$ is a linear separator, then

$$\hat{y} = \text{sign}\left(\frac{\boldsymbol{w}^T}{\rho}\boldsymbol{x} + \frac{b}{\rho}\right) \in \mathcal{H} \text{ where}$$

$$\rho = \min_{\left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}} y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right)$$

# 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 | $\in \mathcal{H}$ |



| $x_1$ | $x_2$ | $y$ | $y(\boldsymbol{w}^T\boldsymbol{x} + b)$ |
|------|------|------|------|
| 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 |

# Computing the Margin

- Claim: $\boldsymbol{w}$ is orthogonal to the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ (the decision boundary)

- A vector is orthogonal to a hyperplane if it is orthogonal to every vector in that hyperplane

- Vectors $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are orthogonal if $\boldsymbol{\alpha}^T\boldsymbol{\beta} = 0$

- Proof:
  - Let $\boldsymbol{x}'$ and $\boldsymbol{x}''$ be two arbitrary points on $\boldsymbol{w}^T\boldsymbol{x} + b = 0$
    - $\boldsymbol{x}' - \boldsymbol{x}''$ is a vector on $\boldsymbol{w}^T\boldsymbol{x} + b = 0$
    - $\boldsymbol{w}^T\boldsymbol{x} + b = 0 \rightarrow \boldsymbol{w}^T\boldsymbol{x} = -b$
  - $\boldsymbol{w}^T(\boldsymbol{x}' - \boldsymbol{x}'') = \boldsymbol{w}^T\boldsymbol{x}' - \boldsymbol{w}^T\boldsymbol{x}'' = -b + b = 0$ ∎

# Computing the Margin

- Claim: $\boldsymbol{w}$ is orthogonal to the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ (the decision boundary)

- A vector is orthogonal to a hyperplane if it is orthogonal to every vector in that hyperplane

- Vectors $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are orthogonal if $\boldsymbol{\alpha}^T\boldsymbol{\beta} = 0$



$$\boldsymbol{w}$$

$$\boldsymbol{x}''$$

$$\boldsymbol{x}'$$

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

# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane



$$\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$$

$\boldsymbol{x}''$

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

$\boldsymbol{x}'$

# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane

$$\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$$

$\boldsymbol{x}''$

$\boldsymbol{x}'$

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

# Computing the Margin

- Let $\boldsymbol{x}'$ be an arbitrary point on the hyperplane $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ and let $\boldsymbol{x}''$ be an arbitrary point

- The distance between $\boldsymbol{x}''$ and $\boldsymbol{w}^T\boldsymbol{x} + b = 0$ is equal to the magnitude of the projection of $\boldsymbol{x}'' - \boldsymbol{x}'$ onto $\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$, the unit vector orthogonal to the hyperplane



$$\boldsymbol{x}''$$

$$\dfrac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$$

$$\boldsymbol{x}'$$

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

# Computing the Margin

- Let $x'$ be an arbitrary point on the hyperplane $h(x) = w^T x + b = 0$ and let $x''$ be an arbitrary point

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

$$d(x'', h) = \left| \frac{w^T(x'' - x')}{\|w\|_2} \right| = \frac{|w^T x'' - w^T x'|}{\|w\|_2}$$

$$= \frac{|w^T x'' + b|}{\|w\|_2}$$

# Computing the Margin

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

$$\min_{\left(\boldsymbol{x}^{(i)},\boldsymbol{y}^{(i)}\right)\in\mathcal{D}} d\left(\boldsymbol{x}^{(i)},h\right) = \min_{\left(\boldsymbol{x}^{(i)},\boldsymbol{y}^{(i)}\right)\in\mathcal{D}} \frac{\left|\boldsymbol{w}^T\boldsymbol{x}^{(i)}+b\right|}{\|\boldsymbol{w}\|_2}$$

$$= \frac{1}{\|\boldsymbol{w}\|_2} \min_{\left(\boldsymbol{x}^{(i)},\boldsymbol{y}^{(i)}\right)\in\mathcal{D}} \left|\boldsymbol{w}^T\boldsymbol{x}^{(i)}+b\right|$$

$$= \frac{1}{\|\boldsymbol{w}\|_2} \min_{\left(\boldsymbol{x}^{(i)},\boldsymbol{y}^{(i)}\right)\in\mathcal{D}} \boldsymbol{y}^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)}+b\right)$$

$$= \frac{1}{\|\boldsymbol{w}\|_2}$$

# Maximizing the Margin

$$\text{maximize } \frac{1}{\|\boldsymbol{w}\|_2}$$

$$\text{subject to } \min_{\left(\boldsymbol{x}^{(i)},y^{(i)}\right)\in\mathcal{D}} y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)}+b\right)=1$$

$$\updownarrow$$

$$\text{minimize } \|\boldsymbol{w}\|_2$$

$$\text{subject to } \min_{\left(\boldsymbol{x}^{(i)},y^{(i)}\right)\in\mathcal{D}} y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)}+b\right)=1$$

$$\updownarrow$$

$$\text{minimize } \frac{1}{2}\|\boldsymbol{w}\|_2^2$$

$$\text{subject to } \min_{\left(\boldsymbol{x}^{(i)},y^{(i)}\right)\in\mathcal{D}} y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)}+b\right)=1$$

$$\updownarrow$$

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

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

# Maximizing the Margin

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

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

- If $\left[\hat{b}, \widehat{\boldsymbol{w}}\right]$ is the optimal solution, then $\exists$ at least one training data point $\left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$ s.t $y^{(i)}\left(\widehat{\boldsymbol{w}}^T\boldsymbol{x}^{(i)} + \hat{b}\right) = 1$

  - All training data points $\left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$ where $y^{(i)}\left(\widehat{\boldsymbol{w}}^T\boldsymbol{x}^{(i)} + \hat{b}\right) = 1$ 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
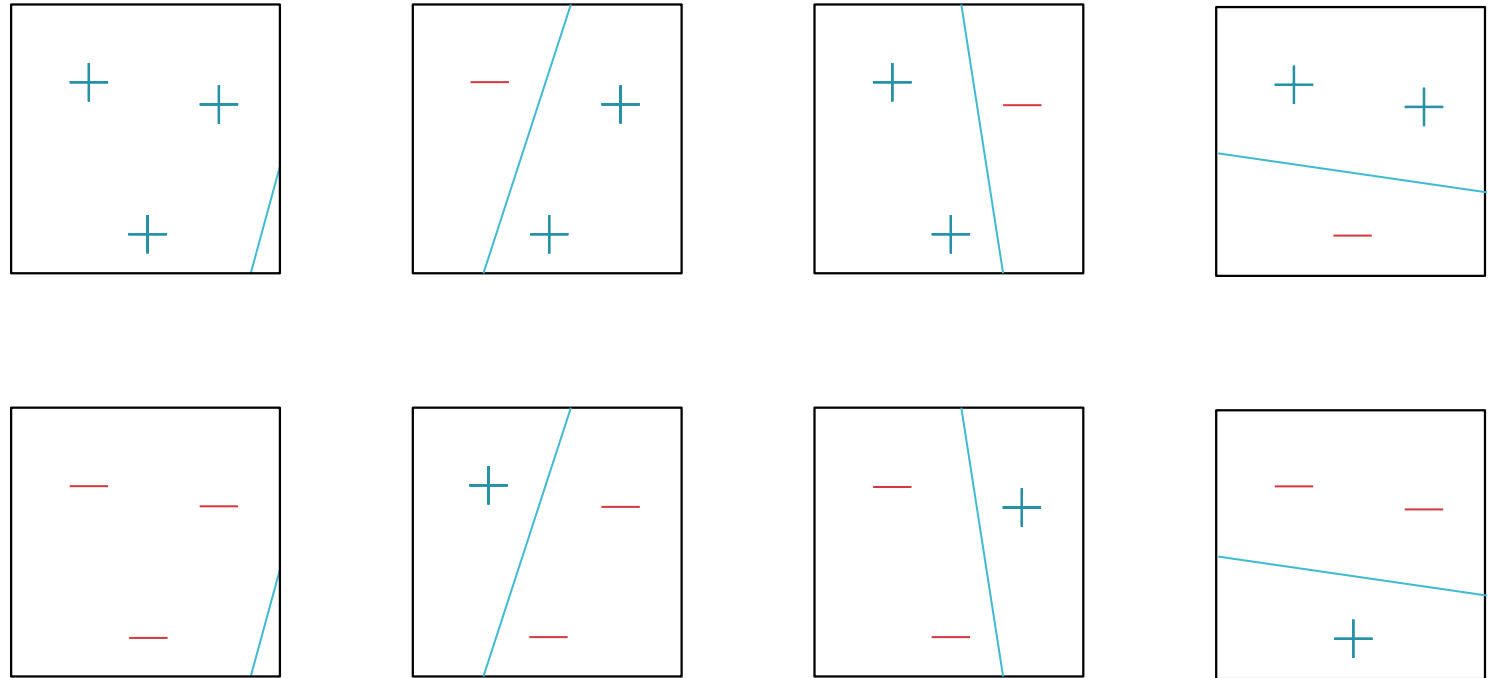
# Recipe for SVMs

- Define a model and model parameters

  - Assume a linear decision boundary (with normalized weights)

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

  - Parameters: $\boldsymbol{w} = [w_1, \dots, w_D]$ and $b$

- Write down an objective function (with constraints)

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

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

- 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 $\mathcal{H}$ = {all linear separators} and

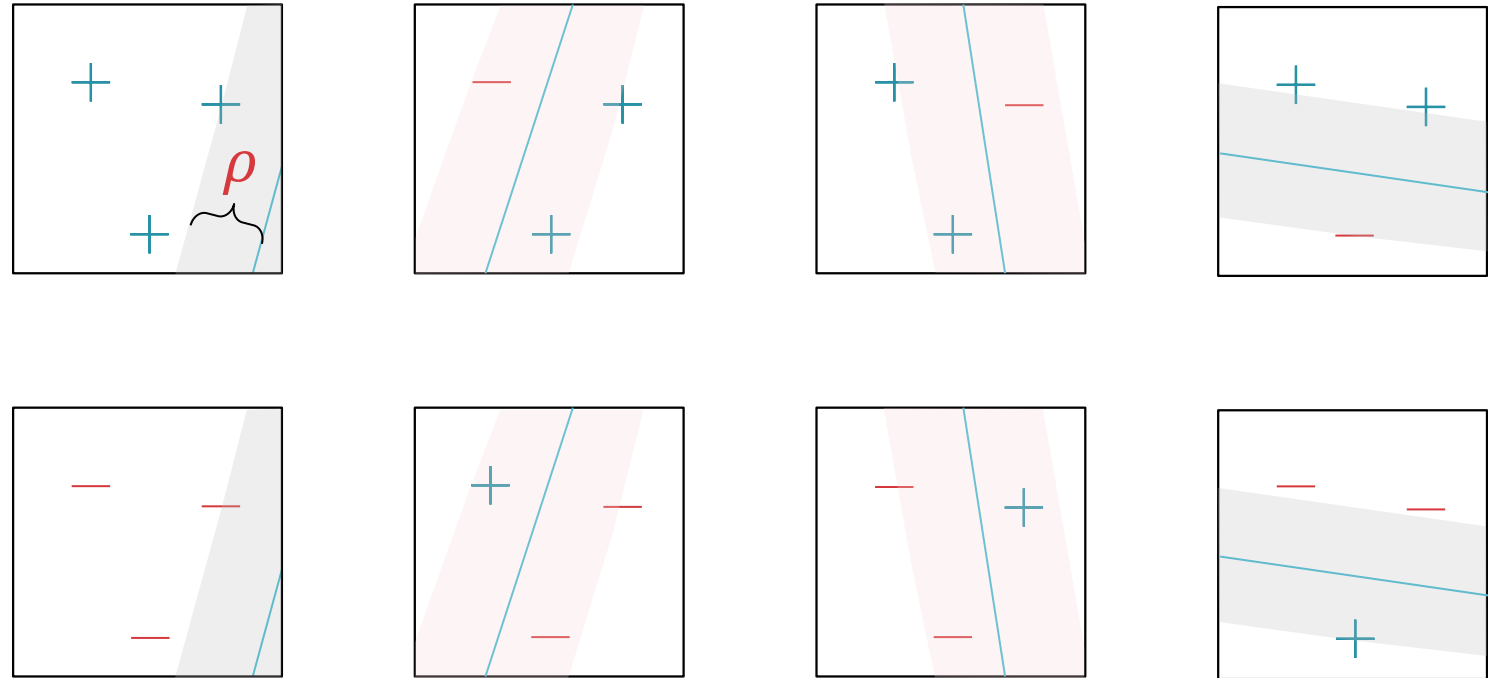  $\mathcal{H}_\rho$ = {all linear separators with minimum margin $\rho$}

# 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-colinear) data points in this space

# Why Maximal Margins?

- Consider three binary data points in a **bounded** 2-D space

- $\mathcal{H}_{\rho} = \{$all linear separators with minimum margin $\rho\}$ cannot always correctly classify three non-colinear 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

     - How much training error should we tolerate?

  2. Apply a non-linear transformation that shifts the data into a space where it is linearly separable

     - How can we pick a non-linear transformation?

# SVMs

minimize $\dfrac{1}{2}\boldsymbol{w}^T\boldsymbol{w}$

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

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

# Hard-margin SVMs

minimize $\dfrac{1}{2} \boldsymbol{w}^T \boldsymbol{w}$

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

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

# Soft-margin SVMs

minimize $\dfrac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C\displaystyle\sum_{i=1}^{N}\xi^{(i)}$

subject to $y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) \geq 1 - \xi^{(i)} \; \forall \left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$

$\xi^{(i)} \geq 0 \qquad\qquad\qquad\qquad\qquad \forall\, i \in \{1, \dots, N\}$

# Soft-margin SVMs

minimize $\dfrac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C\sum\limits_{i=1}^{N}\xi^{(i)}$

subject to $y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) \geq 1 - \xi^{(i)} \ \forall \left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$

$\qquad\qquad \xi^{(i)} \geq 0 \qquad\qquad \forall\, i \in \{1, \dots, N\}$

- $\xi^{(i)}$ is the "soft" error on the $i^{th}$ training data point
  - If $\xi^{(i)} > 1$, then $y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) < 0 \Rightarrow$
    $\left(\boldsymbol{x}^{(i)}, y^{(i)}\right)$ is incorrectly classified
  - If $0 < \xi^{(i)} < 1$, then $y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) > 0 \Rightarrow$
    $\left(\boldsymbol{x}^{(i)}, y^{(i)}\right)$ is correctly classified but inside the margin
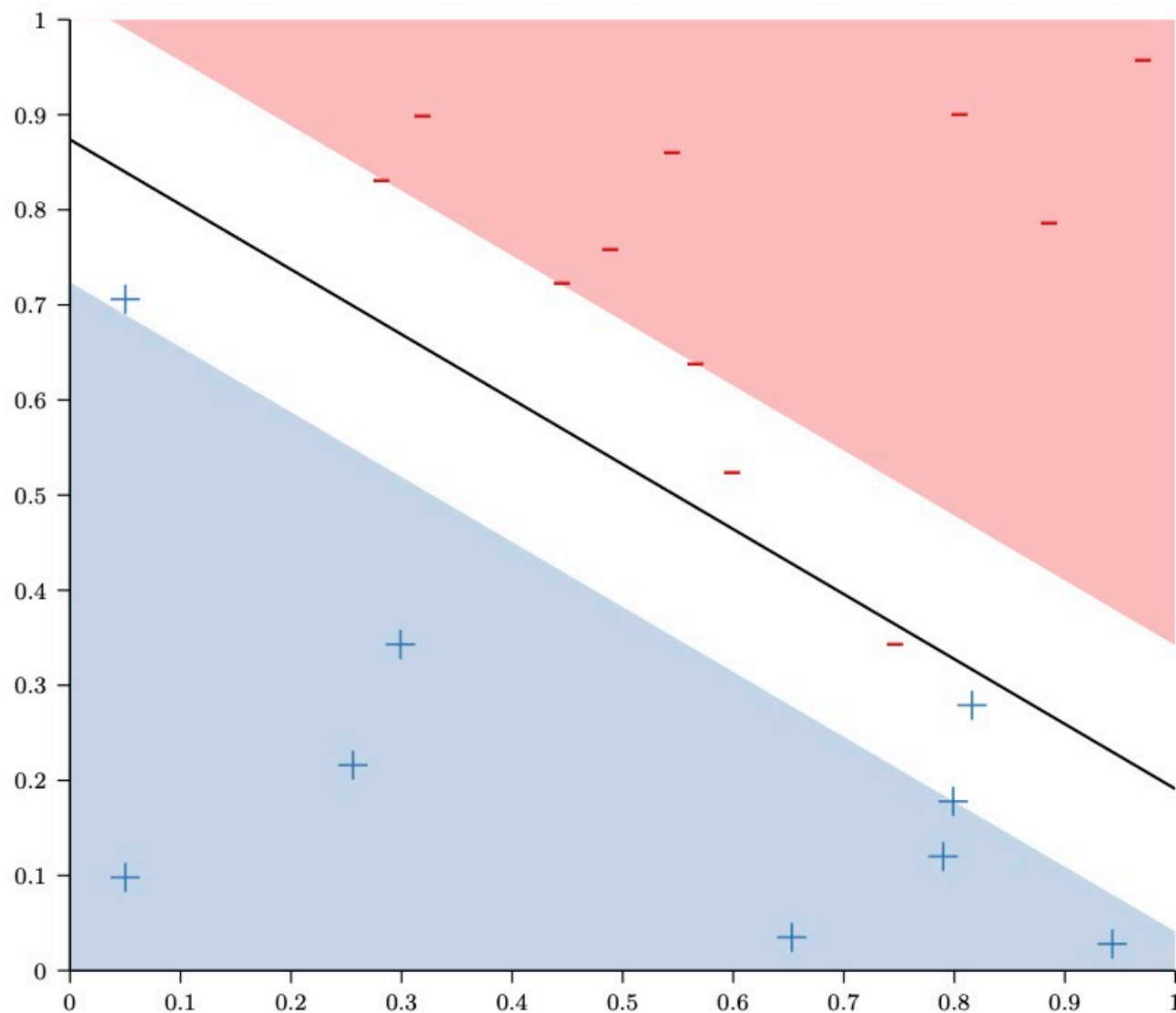- $\sum\limits_{i=1}^{N}\xi^{(i)}$ is the "soft" training error

# Soft-margin SVMs

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

subject to $y^{(i)}\left(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\right) \geq 1 - \xi^{(i)} \ \forall \left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$

$\xi^{(i)} \geq 0 \qquad\qquad\qquad \forall\, i \in \{1, \dots, N\}$

- Still solvable using quadratic programming

- All training data points $\left(\boldsymbol{x}^{(i)}, y^{(i)}\right) \in \mathcal{D}$ where $y^{(i)}\left(\widehat{\boldsymbol{w}}^T\boldsymbol{x}^{(i)} + \hat{b}\right) \leq 1$ are known as *support vectors*
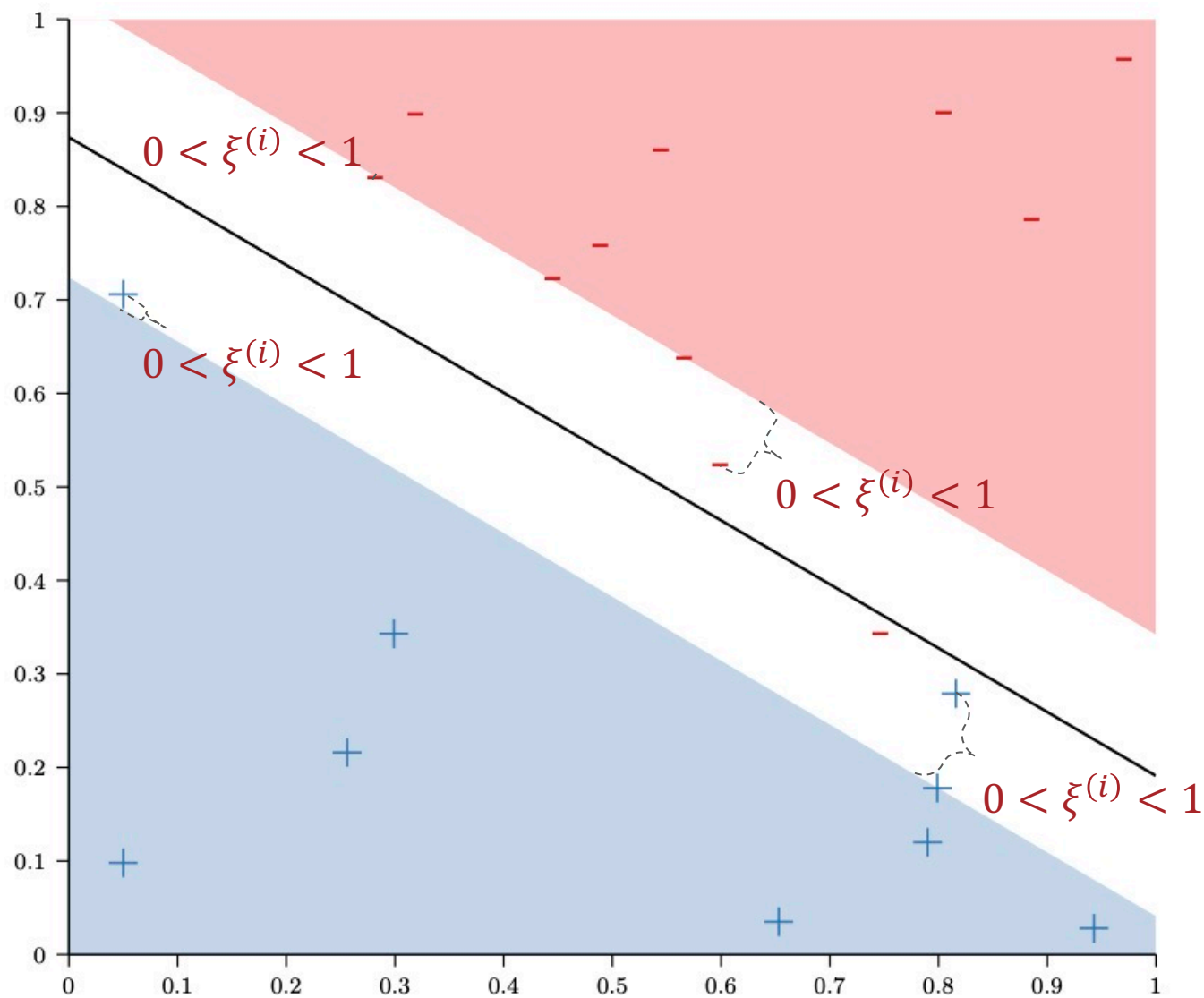
# Interpreting $\xi^{(i)}$

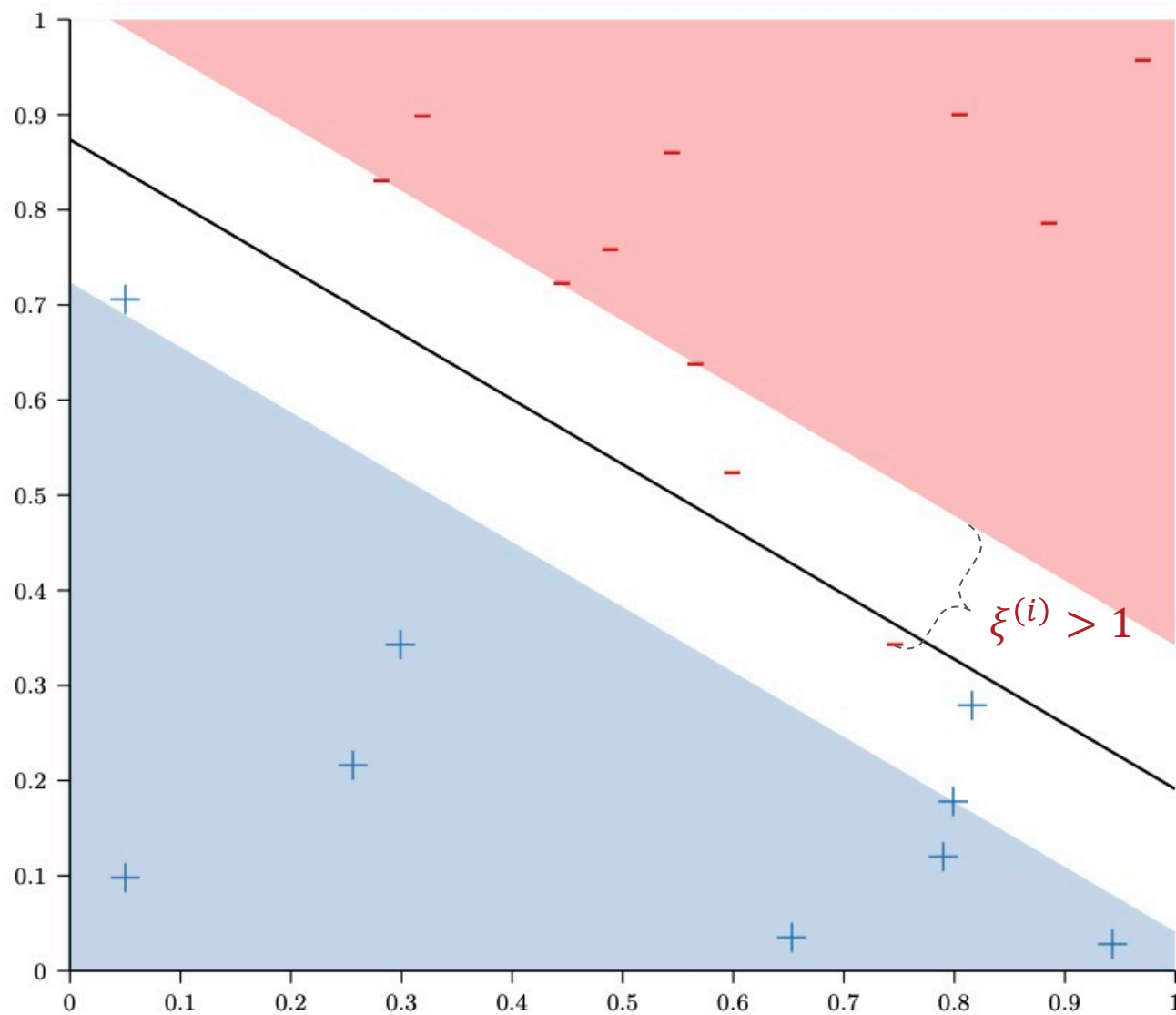# Interpreting $\xi^{(i)}$

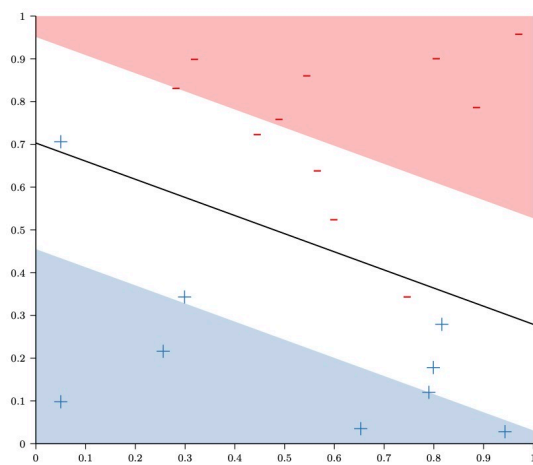# Interpreting $\xi^{(i)}$

# Interpreting $\xi^{(i)}$



"margin" support vector

"margin" support vector

"margin" support vector

# Interpreting $\xi^{(i)}$

# Interpreting $\xi^{(i)}$



$\xi^{(i)} > 1$

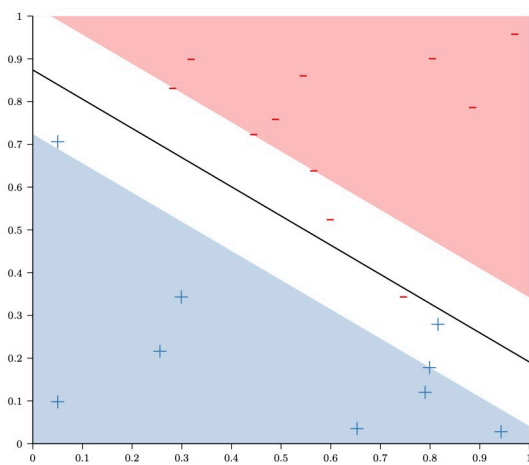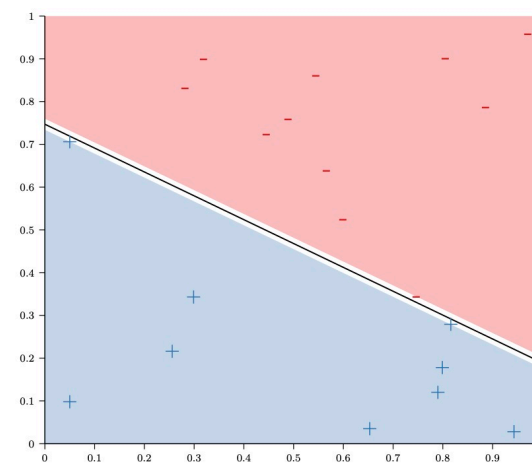Smaller $C$                                   Larger $C$          Hard Margin

# Setting $C$

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