



Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**



Machine Learning: Chenhao Tan

University of Colorado Boulder

LECTURE 17

Slides adapted from Jordan Boyd-Graber, Chris Ketelsen

Logistics

- Homework 3 is due on Sunday!
- Project team matches

Roadmap

- Last time: linear SVM formulation when data is linearly separable
- This time:
 - Make linear SVM work when data is not linearly separable
 - Introduce duality
- Next week: KKT conditions & Kernel tricks

Overview

Soft-margin SVM

Duality

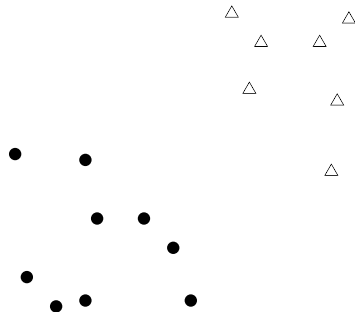
Outline

Soft-margin SVM

Duality

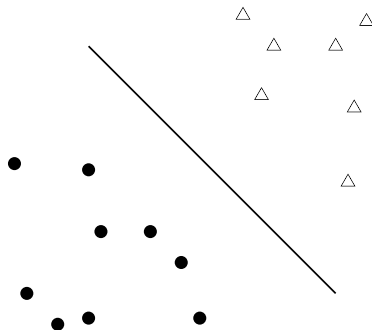
Support Vector Machines

- 2-class training data



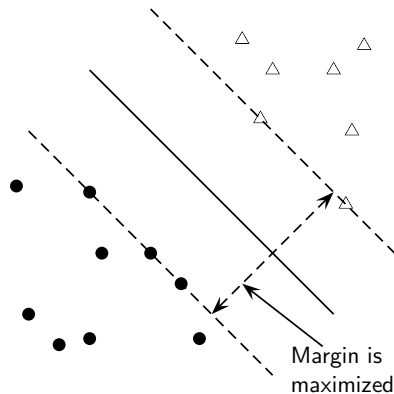
Support Vector Machines

- 2-class training data
- decision boundary \rightarrow **linear separator**



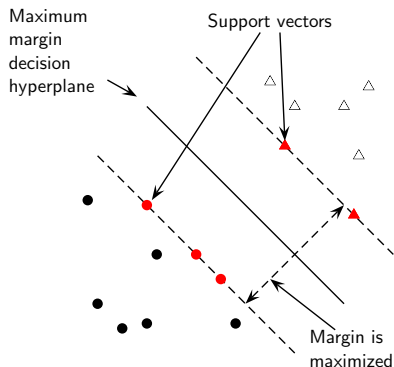
Support Vector Machines

- 2-class training data
- decision boundary \rightarrow **linear separator**
- criterion: being maximally far away from any data point \rightarrow determines classifier **margin**



Support Vector Machines

- 2-class training data
- decision boundary → **linear separator**
- criterion: being maximally far away from any data point → determines classifier **margin**
- linear separator position defined by **support vectors**
- other points have no impact on the decision boundary



Objective function for hard-margin SVM

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

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]$

Theoretical evidence that suggests SVMs will Work

- Leave-one-out error
- Margin analysis (omitted)
- VC Dimension (omitted for now)

Leave-one-out error (sketch)

Leave one out error is the error by using one point as your test set (averaged over all such points).

$$\hat{R}_{LOO} = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [h_{s-\{x_i\}} \neq y_i]$$

Leave-one-out error (sketch)

Leave one out error is the error by using one point as your test set (averaged over all such points).

$$\hat{R}_{LOO} = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [h_{s-\{x_i\}} \neq y_i]$$

This serves as an unbiased estimate of generalization error for samples of size $m - 1$.

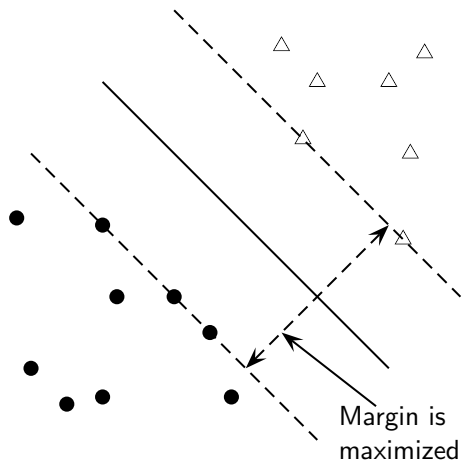
Leave-one-out error (sketch)

Leave-one-out error is bounded by the number of support vectors.

$$\mathbb{E}_{S \sim D^{m-1}} [R(h_S)] \leq \mathbb{E}_{S \sim D^m} \left[\frac{N_{SV}(S)}{m} \right]$$

Consider the held out error for x_i .

Pictorial proof



Leave-one-out error (sketch)

Leave-one-out error is bounded by the number of support vectors.

$$\mathbb{E}_{S \sim D^{m-1}} [R(h_S)] \leq \mathbb{E}_{S \sim D^m} \left[\frac{N_{SV}(S)}{m} \right]$$

Consider the held out error for x_i .

- If x_i was not a support vector, the answer doesn't change.
- If x_i was a support vector, it could change the answer; this is when we can have an error.

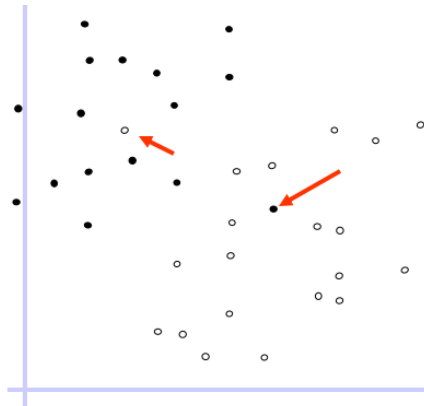
There are N_{SV} support vectors and thus N_{SV} possible errors.

Objective function for hard-margin SVM

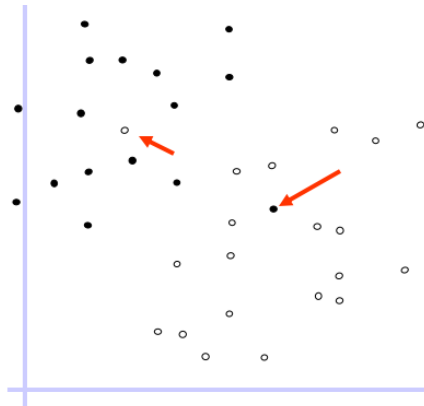
$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]$

Can SVMs Work Here?

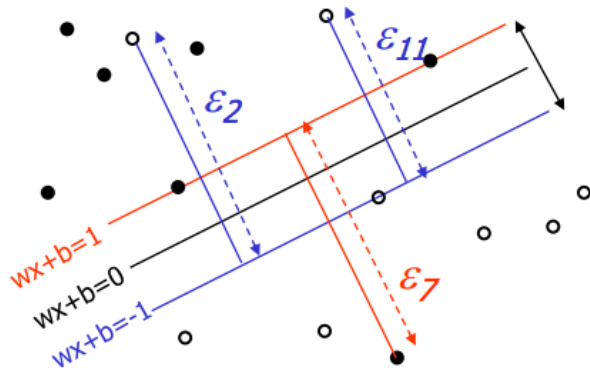


Can SVMs Work Here?



$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

Trick: Allow for a few bad apples



Hard-margin objective function

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

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]$

Relaxing the constraint

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

- $\xi_i = 0$ means at least one margin on correct side of decision boundary
- $\xi_i = 1/2$ means at least one-half margin on correct side of decision boundary
- $\xi_i = 2$ means at least one margin on wrong side of decision boundary

New objective function

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

subject to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, i \in [1, m]$$

$$\xi_i \geq 0, i \in [1, m]$$

New objective function

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

subject to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, i \in [1, m]$$

$$\xi_i \geq 0, i \in [1, m]$$

- Standard margin

New objective function

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

subject to

$$\begin{aligned} y_i(\mathbf{w} \cdot \mathbf{x}_i + b) &\geq 1 - \xi_i, i \in [1, m] \\ \xi_i &\geq 0, i \in [1, m] \end{aligned}$$

- Standard margin
- How wrong a point is (slack variables)

New objective function

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \textcolor{red}{C} \sum_i \xi_i$$

subject to

$$\begin{aligned} y_i(\mathbf{w} \cdot \mathbf{x}_i + b) &\geq 1 - \xi_i, i \in [1, m] \\ \xi_i &\geq 0, i \in [1, m] \end{aligned}$$

- Standard margin
- How wrong a point is (slack variables)
- **Tradeoff between margin and slack variables**

What is the role of C ?

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

subject to

$$\begin{aligned} y_i(\mathbf{w} \cdot \mathbf{x}_i + b) &\geq 1 - \xi_i, i \in [1, m] \\ \xi_i &\geq 0, i \in [1, m] \end{aligned}$$

- A. $C \uparrow \Rightarrow$ decrease bias, decrease variance
- B. $C \uparrow \Rightarrow$ decrease bias, increase variance
- C. $C \uparrow \Rightarrow$ increase bias, decrease variance
- D. $C \uparrow \Rightarrow$ increase bias, increase variance

Outline

Soft-margin SVM

Duality

Binary classification

Given: $S_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ training examples, $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$

Goal: Find hypothesis function $h : X \rightarrow Y$

Linear SVM: learn a linear decision rule of the form $\mathbf{w} \cdot \mathbf{x} + b$

Optimizing the objective function

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

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]$

This is a quadratic objective function with linear inequality constraints. Many off-the-shelf optimization methods are available.

Optimizing Constrained Functions

The Method of Lagrange Multipliers

Constrained problem (Primal problem)

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t. } & g_i(\mathbf{x}) \geq 0, i \in [1, n] \end{aligned}$$

Lagrange Multiplier

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}) &= f(\mathbf{x}) - \sum_{i=1}^n \alpha_i g_i(\mathbf{x}), \\ \alpha_i &\geq 0, i \in [1, n] \end{aligned}$$

Lagrange Multiplier

p^* : the optimal value in the primal problem

We claim that

$$p^* = \min_{\mathbf{x}} \max_{\boldsymbol{\alpha}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}) = \min_{\mathbf{x}} \max_{\boldsymbol{\alpha}} f(\mathbf{x}) - \sum_{i=1}^n \alpha_i g_i(\mathbf{x})$$

Lagrange Multiplier

p^* : the optimal value in the primal problem

We claim that

$$p^* = \min_{\mathbf{x}} \max_{\boldsymbol{\alpha}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}) = \min_{\mathbf{x}} \max_{\boldsymbol{\alpha}} f(\mathbf{x}) - \sum_{i=1}^n \alpha_i g_i(\mathbf{x})$$

This is because

$$\max -\alpha y = \begin{cases} 0 & y \geq 0 \\ +\infty & \text{otherwise} \end{cases}$$

Lagrange Multiplier

What happens if we reverse min and max:

$$\max_{\alpha} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \alpha) \geq \text{or} \leq \min_{\mathbf{x}} \max_{\alpha} \mathcal{L}(\mathbf{x}, \alpha)$$

Lagrange Multiplier

What happens if we reverse min and max:

$$\max_{\alpha} \min_x \mathcal{L}(\mathbf{x}, \alpha) \leq \min_x \max_{\alpha} \mathcal{L}(\mathbf{x}, \alpha)$$

The left leads to the dual problem.

Primal vs. Dual

Primal problem

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

$$\text{s.t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]$$

Derive the function for dual problem.
Replace \mathbf{w}, b with stationarity conditions.
(There will be detailed derivations for the soft-margin case later.)

Primal vs. Dual

Primal problem

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

$$\text{s.t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]$$

Dual problem

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_j \cdot \mathbf{x}_i)$$

$$\text{s.t. } \alpha_i \geq 0, i \in [1, m]$$

$$\sum_i \alpha_i y_i = 0$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \alpha_i \geq 0$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \alpha_i \geq 0$$

Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^m \alpha_i y_i = 0$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \alpha_i \geq 0$$

Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \quad \sum_{i=1}^m \alpha_i y_i = 0$$

Remember that two properties about support vector machine directly follows from this:

- Only support vectors affect the weights ($\alpha_i > 0$).
- There must be both positive and negative support vectors.

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \alpha_i \geq 0$$

Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^m \alpha_i y_i = 0$$

Complementary slackness

$$\alpha_i = 0 \vee y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$$

What is the dual problem of soft-margin SVM?

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1} \xi_i$$

subject to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, i \in [1, m]$$

$$\xi_i \geq 0, i \in [1, m]$$

New Lagrangian

$$\begin{aligned}\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ & - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] \\ & - \sum_{i=1}^m \beta_i \xi_i\end{aligned}$$

New Lagrangian

$$\begin{aligned}\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ & - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] \\ & - \sum_{i=1}^m \beta_i \xi_i\end{aligned}$$

Taking the gradients $(\nabla_{\mathbf{w}} \mathcal{L}, \nabla_b \mathcal{L}, \nabla_{\xi_i} \mathcal{L})$ and solving for zero gives us

$$\begin{aligned}\mathbf{w} &= \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i & \sum_{i=1}^m \alpha_i y_i &= 0 & \alpha_i + \beta_i &= C\end{aligned}$$

New Lagrangian

$$\begin{aligned}\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ & - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] \\ & - \sum_{i=1}^m \beta_i \xi_i\end{aligned}$$

Taking the gradients $(\nabla_{\mathbf{w}} \mathcal{L}, \nabla_b \mathcal{L}, \nabla_{\xi_i} \mathcal{L})$ and solving for zero gives us

$$\begin{aligned}\mathbf{w} &= \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i & \sum_{i=1}^m \alpha_i y_i &= 0 & \alpha_i + \beta_i &= C\end{aligned}$$

New Lagrangian

$$\begin{aligned}\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ & - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] \\ & - \sum_{i=1}^m \beta_i \xi_i\end{aligned}$$

Taking the gradients ($\nabla_{\mathbf{w}} \mathcal{L}$, $\nabla_b \mathcal{L}$, $\nabla_{\xi_i} \mathcal{L}$) and solving for zero gives us

$$\begin{aligned}\mathbf{w} &= \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i & \sum_{i=1}^m \alpha_i y_i &= 0 & \alpha_i + \beta_i &= C\end{aligned}$$

New Lagrangian

$$\begin{aligned}\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ & - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] \\ & - \sum_{i=1}^m \beta_i \xi_i\end{aligned}$$

Taking the gradients ($\nabla_{\mathbf{w}} \mathcal{L}$, $\nabla_b \mathcal{L}$, $\nabla_{\xi_i} \mathcal{L}$) and solving for zero gives us

$$\begin{aligned}\mathbf{w} &= \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i & \sum_{i=1}^m \alpha_i y_i &= 0 & \alpha_i + \beta_i &= C\end{aligned}$$

Simplifying dual objective

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

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

$$\alpha_i + \beta_i = C$$

Simplifying dual objective

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

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

$$\alpha_i + \beta_i = C$$

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$- \sum_{i=1}^m \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i]$$

$$- \sum_{i=1}^m \beta_i \xi_i$$

Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_j \cdot \mathbf{x}_i) \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0, i \in [1, m] \\ & \sum_i \alpha_i y_i = 0 \end{aligned}$$

Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_j \cdot \mathbf{x}_i) \\ \text{s.t.} \quad & \mathbf{C} \geq \alpha_i \geq 0, i \in [1, m] \\ & \sum_i \alpha_i y_i = 0 \end{aligned}$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0, C \geq \alpha_i \geq 0, \beta_i \geq 0$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0, C \geq \alpha_i \geq 0, \beta_i \geq 0$$

Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^m \alpha_i y_i = 0, \alpha_i + \beta_i = C$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0, C \geq \alpha_i \geq 0, \beta_i \geq 0$$

Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^m \alpha_i y_i = 0, \alpha_i + \beta_i = C$$

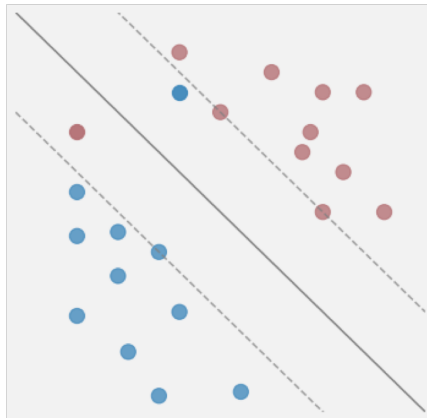
Complementary slackness

$$\alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

More on Complementary Slackness

$$\alpha_i[y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

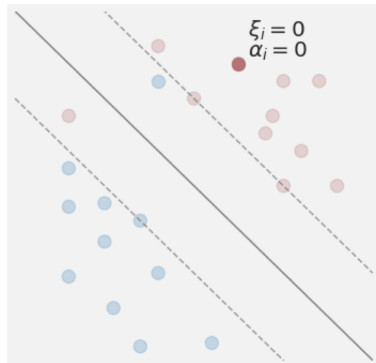
$$\text{Also, } \alpha_i + \beta_i = C$$



More on Complementary Slackness

$$\alpha_i[y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

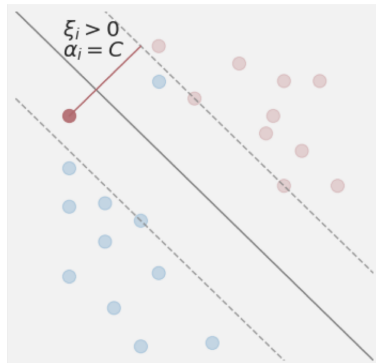
- \mathbf{x}_i satisfies the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) > 1 \Rightarrow \alpha_i = 0$



More on Complementary Slackness

$$\alpha_i[y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

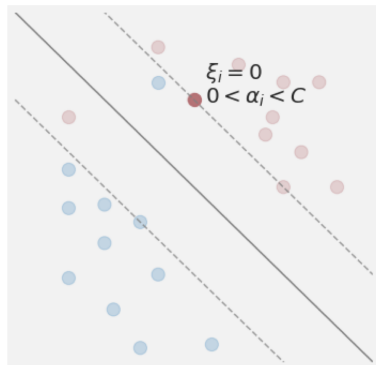
- \mathbf{x}_i satisfies the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) > 1 \Rightarrow \alpha_i = 0$
- \mathbf{x}_i does not satisfy the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) < 1 \Rightarrow \alpha_i = C$



More on Complementary Slackness

$$\alpha_i[y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

- \mathbf{x}_i satisfies the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) > 1 \Rightarrow \alpha_i = 0$
- \mathbf{x}_i does not satisfy the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) < 1 \Rightarrow \alpha_i = C$
- \mathbf{x}_i is on the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 \Rightarrow 0 \leq \alpha_i \leq C$



Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0, C \geq \alpha_i \geq 0, \beta_i \geq 0$$

Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^m \alpha_i y_i = 0, \alpha_i + \beta_i = C$$

Complementary slackness

$$\alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$