



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 15: Duality & Kernels

Slides adapted from Chris Ketelsen, Jordan Boyd-Graber,
and Noah Smith

Administrivia

- HW4 is released
- Final project started!

Outline

- Duality
- Kernels

Outline

- Duality
- Kernels

Binary classification

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

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)

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

Lagrange Multiplier

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

Lagrange Multiplier

p^* : the optimal value in the primal problem

We claim that

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

Handwritten red notes:
min_x f(x)
st. g_i(x) ≥ 0

Lagrange Multiplier

p^* : the optimal value in the primal problem

We claim that

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

This is because

$$\alpha \geq 0$$

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

Lagrange Multiplier

What happens if we reverse min and max:

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

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

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

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

min
w, b *max*
α

Dual problem

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

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

$$\sum_i \alpha_i y_i = 0$$

max
α *min*
w, b

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

$$y_i(\underline{w \cdot x_i + b}) \geq 1, \underline{\alpha_i \geq 0}$$

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

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

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

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

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

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

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

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

Stationarity

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

Complementary slackness

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

What is the dual problem of soft-margin SVM?

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

f(x)

subject to

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

g(x)

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

2m

New Lagrangian

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(w \cdot x_i + \underline{b}) - 1 + \xi_i] - \sum_{i=1}^m \beta_i \xi_i$$

$$\frac{\partial \mathcal{L}}{\partial w} = 0$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi} = 0$$

$$w - \sum_i \alpha_i y_i x_i = 0 \Rightarrow w = \sum_i \alpha_i y_i x_i$$

$$-\sum_i \alpha_i y_i = 0 \Rightarrow \sum_i \alpha_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi} = C - \alpha_i - \beta_i = 0 \Rightarrow \alpha_i + \beta_i = C$$

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

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

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

$$\alpha_i + \beta_i = C$$

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

Dual Problem

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

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

$$\sum_i \alpha_i y_i = 0$$

$$\text{dit } \beta_i = C$$

$$\beta_i \geq 0$$

Dual Problem

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

Karush-Kuhn-Tucker (KKT) conditions

Primal and dual feasibility

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

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

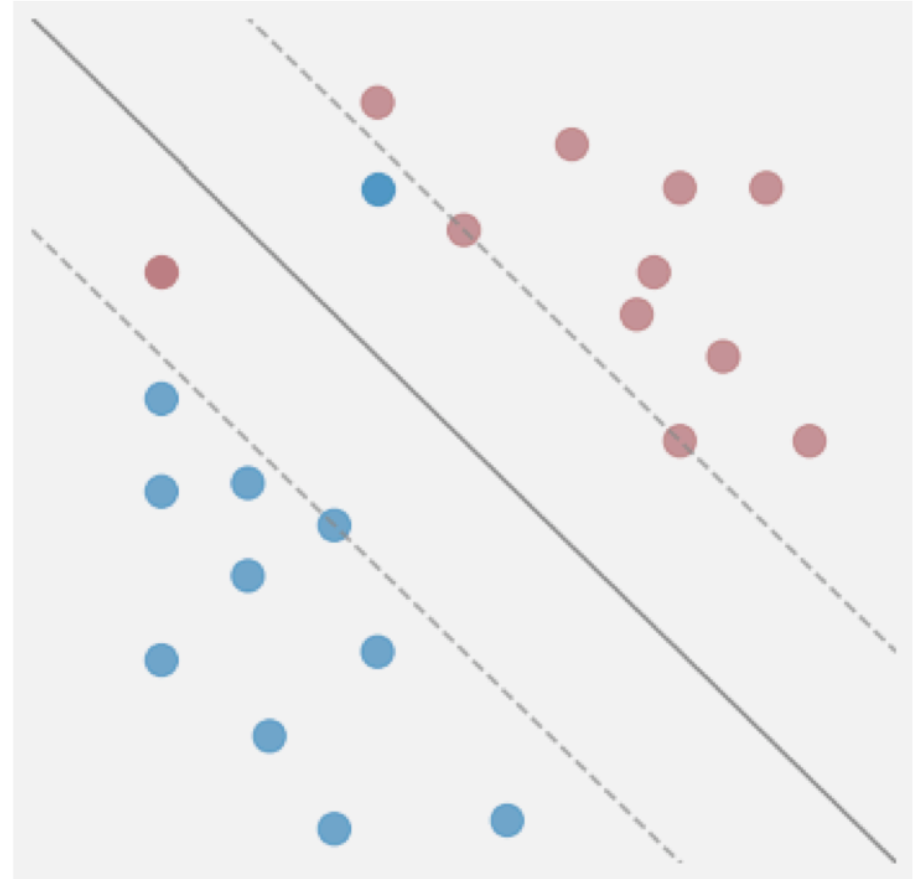
$$\alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

$$\beta_i = 0 \quad \checkmark \quad \xi_i = 0$$

More on Complementary Slackness

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

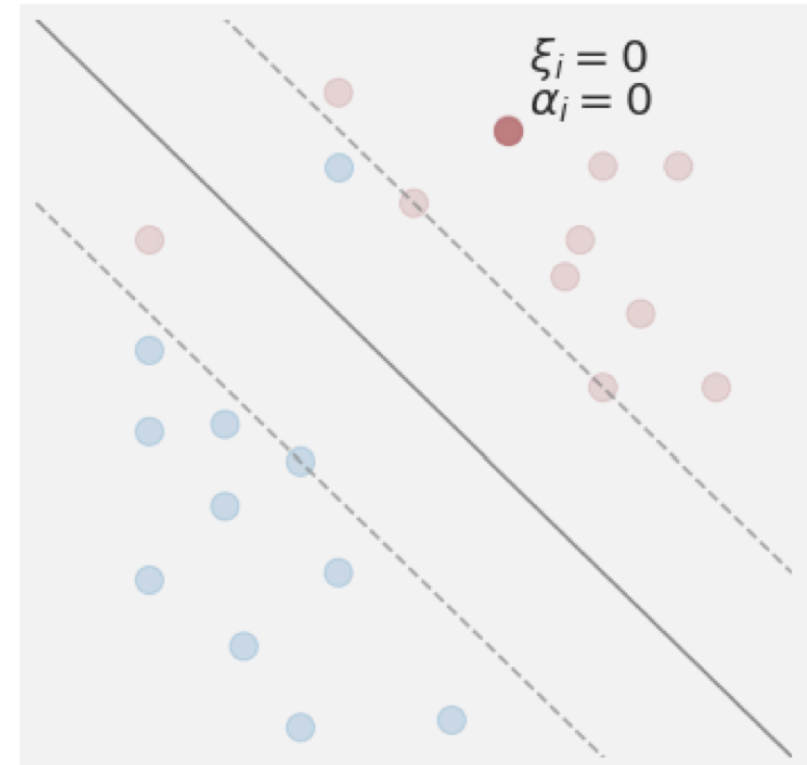
$$\text{Also, } \alpha_i + \beta_i = C$$



More on Complementary Slackness

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

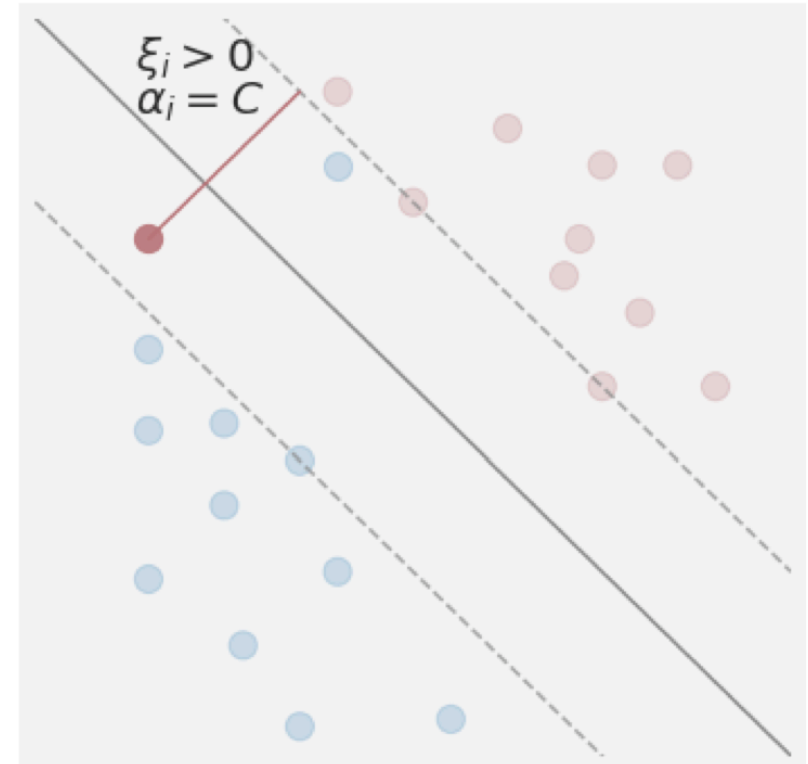
- \mathbf{x}_i satisfies the margin,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) > 1 \Rightarrow \alpha_i = 0 \Rightarrow p_i = C$



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,
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) < 1 \Rightarrow \alpha_i = C$

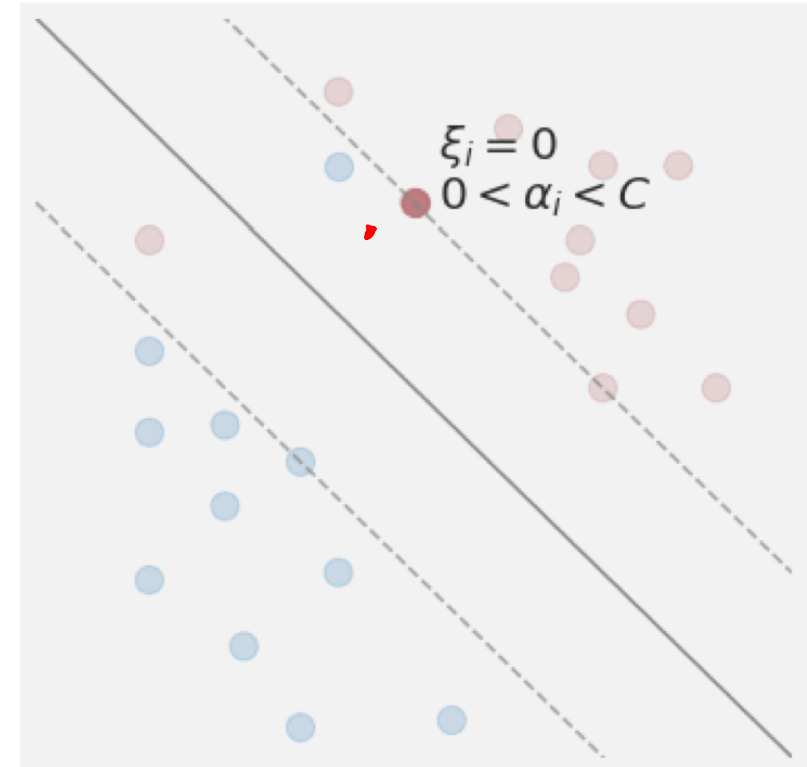


More on Complementary Slackness

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

$$\sum_i \alpha_i y_i = 0$$

- \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

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

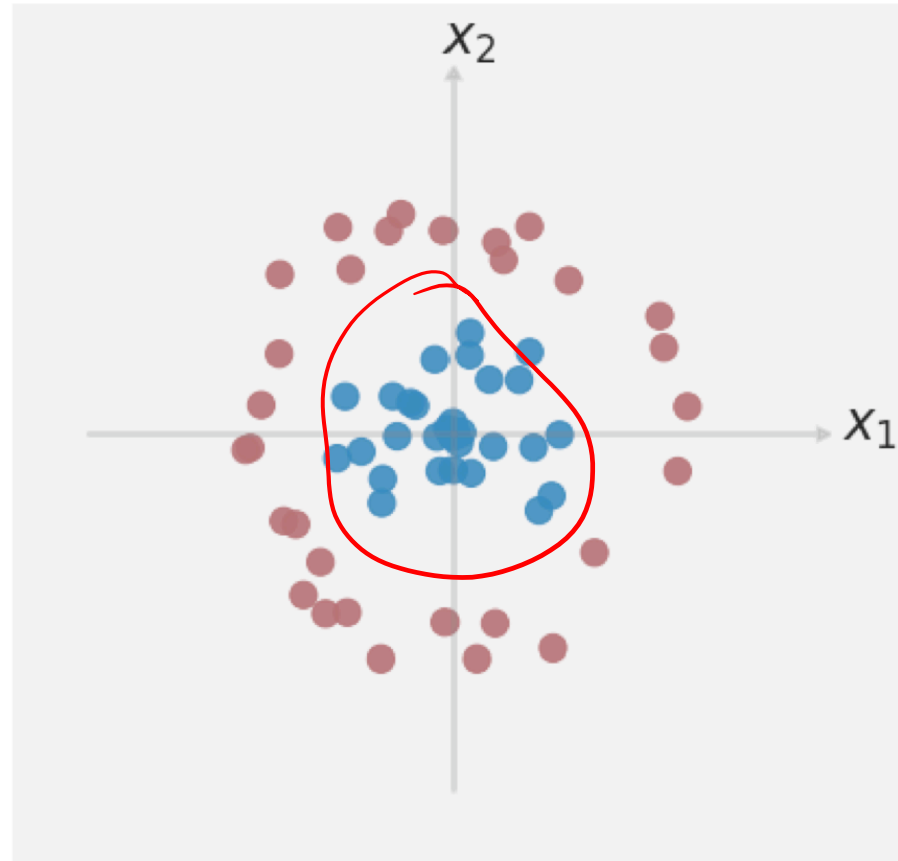
Complementary slackness

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

Outline

- Duality
- **Kernels**

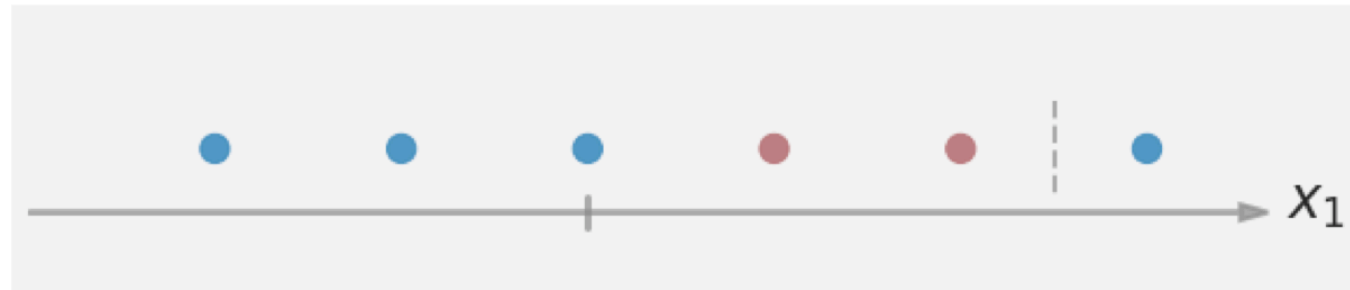
Can you solve this with linear separator?



$$x_1^2 + x_2^2$$

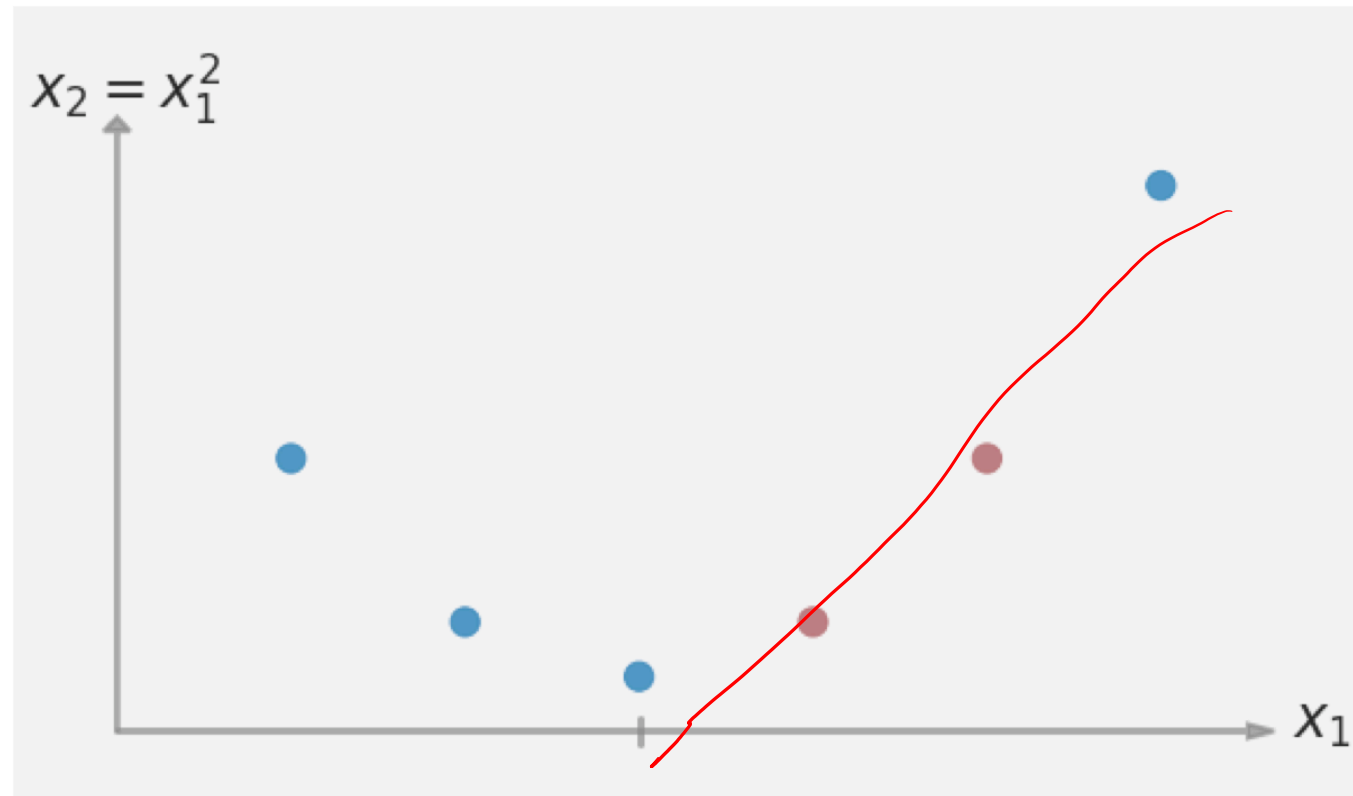
Can you solve this with linear separator?

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



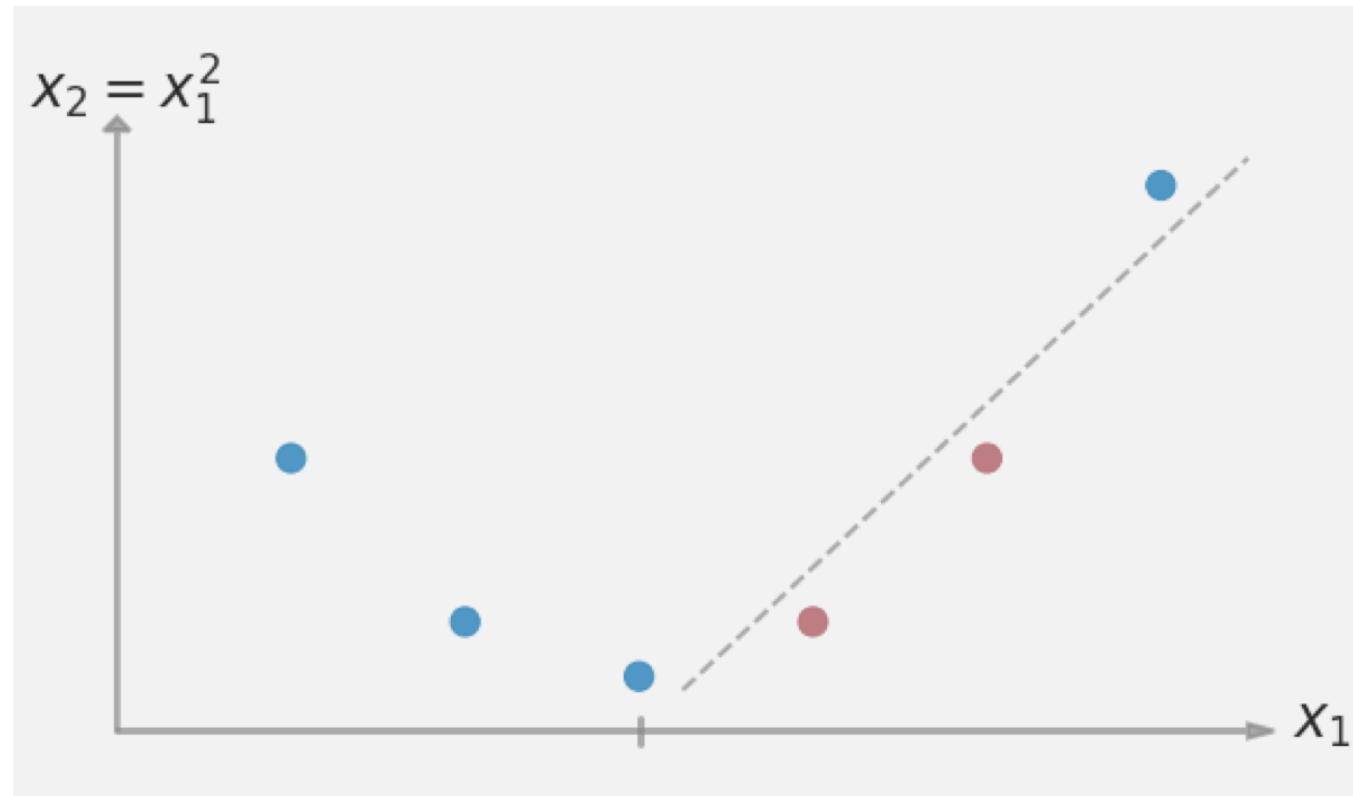
Can you solve this with linear separator?

Add a dimension.



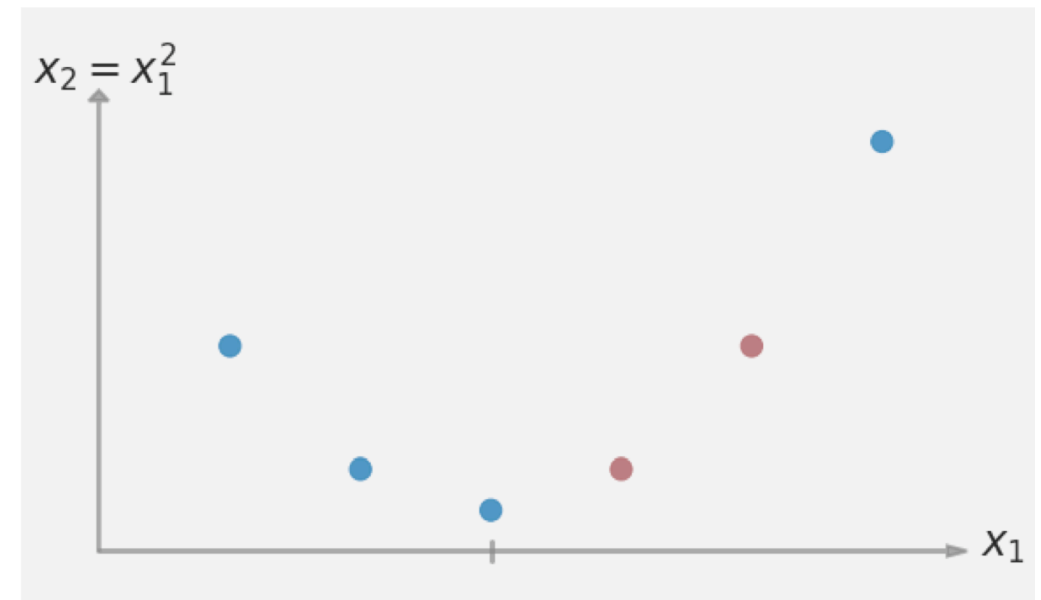
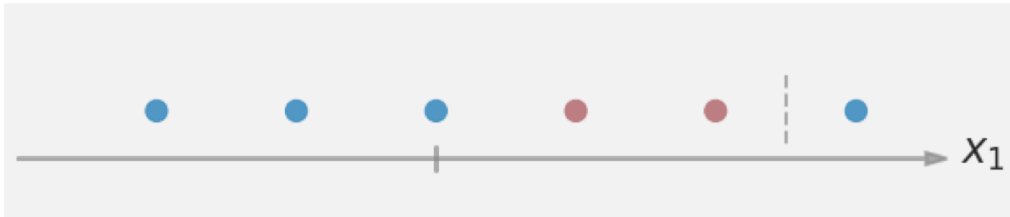
Can you solve this with linear separator?

Add a dimension.



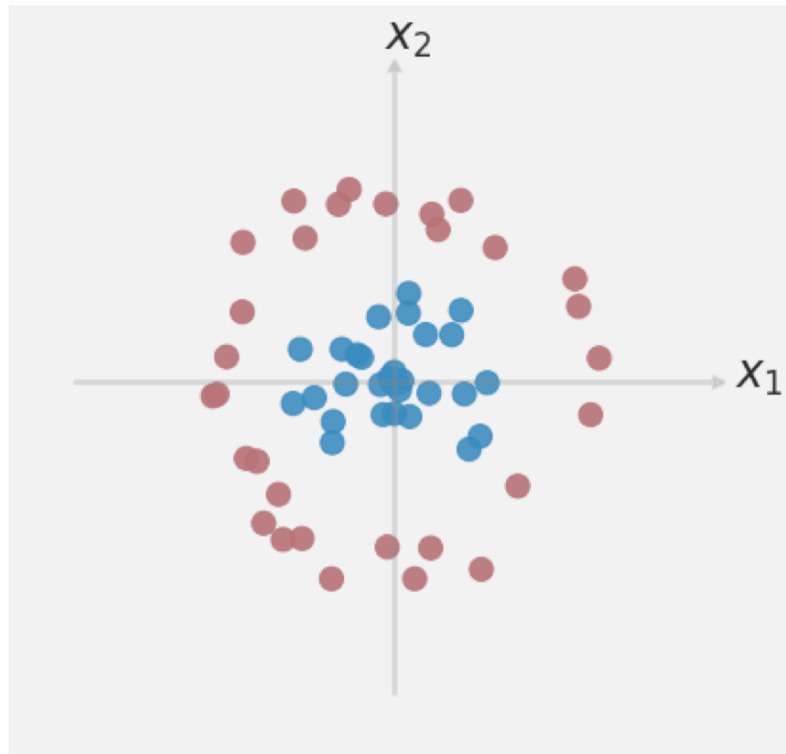
Derived features

We started with the original feature vector, $\mathbf{x} = (x_1)$,
and we created a new derived feature vector, $\phi(\mathbf{x})$ $= (x_1, x_1^2)$.



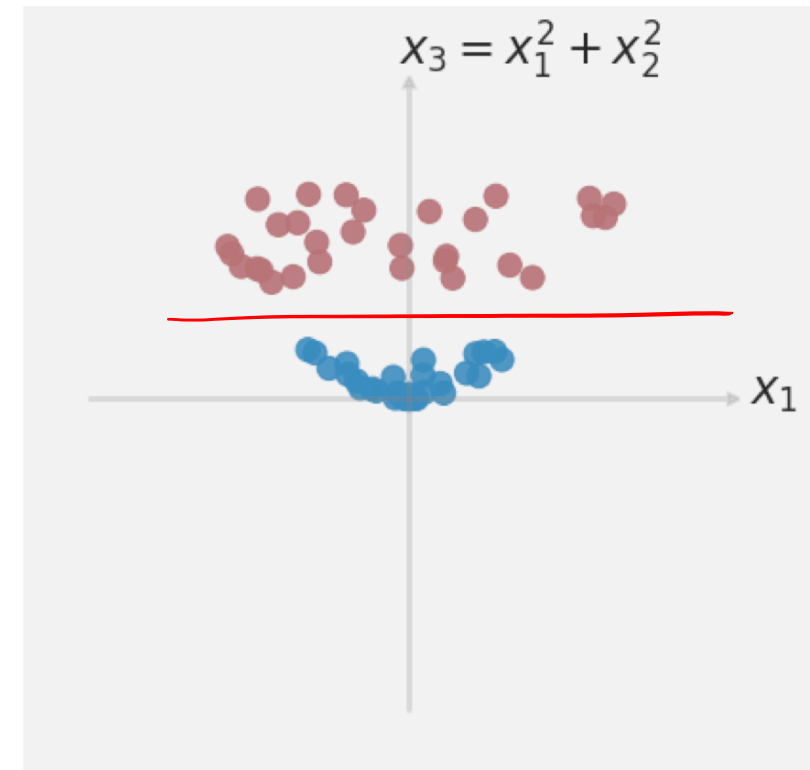
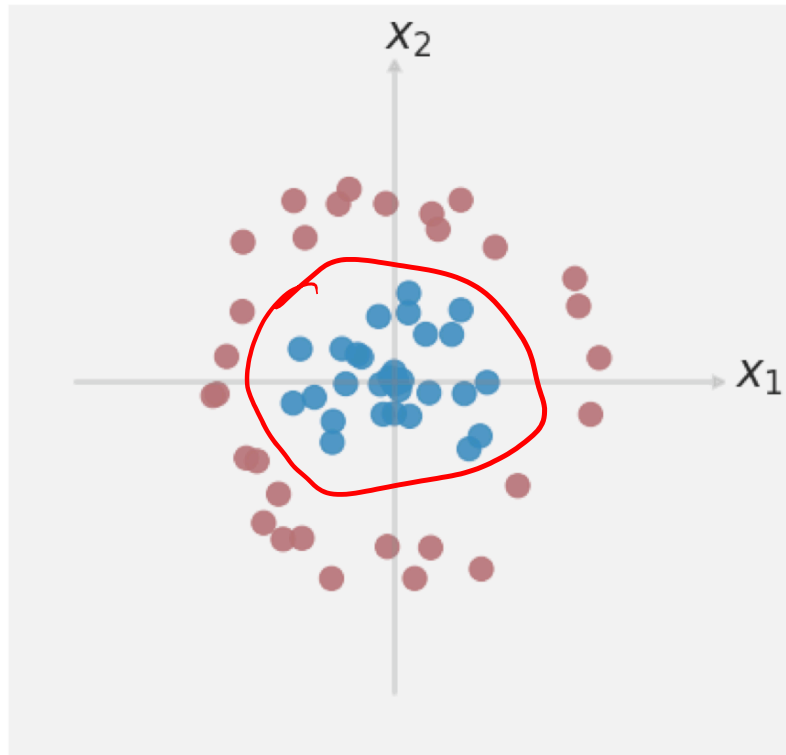
What about the previous problem?

Definitely not separable in two dimensions.



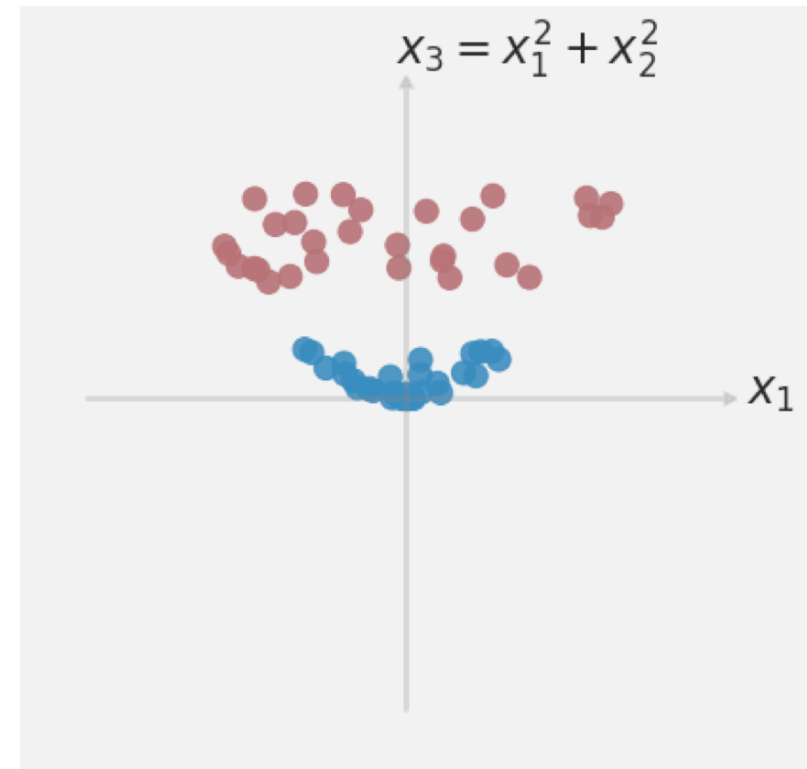
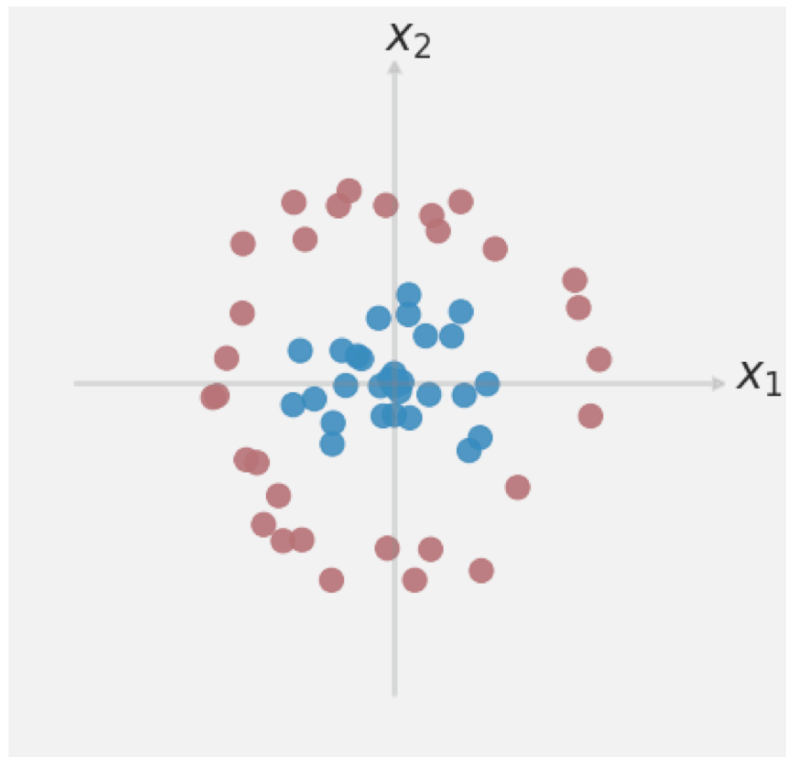
What about the previous problem?

Definitely not separable in two dimensions.
But in three dimensions, it becomes easily separable.



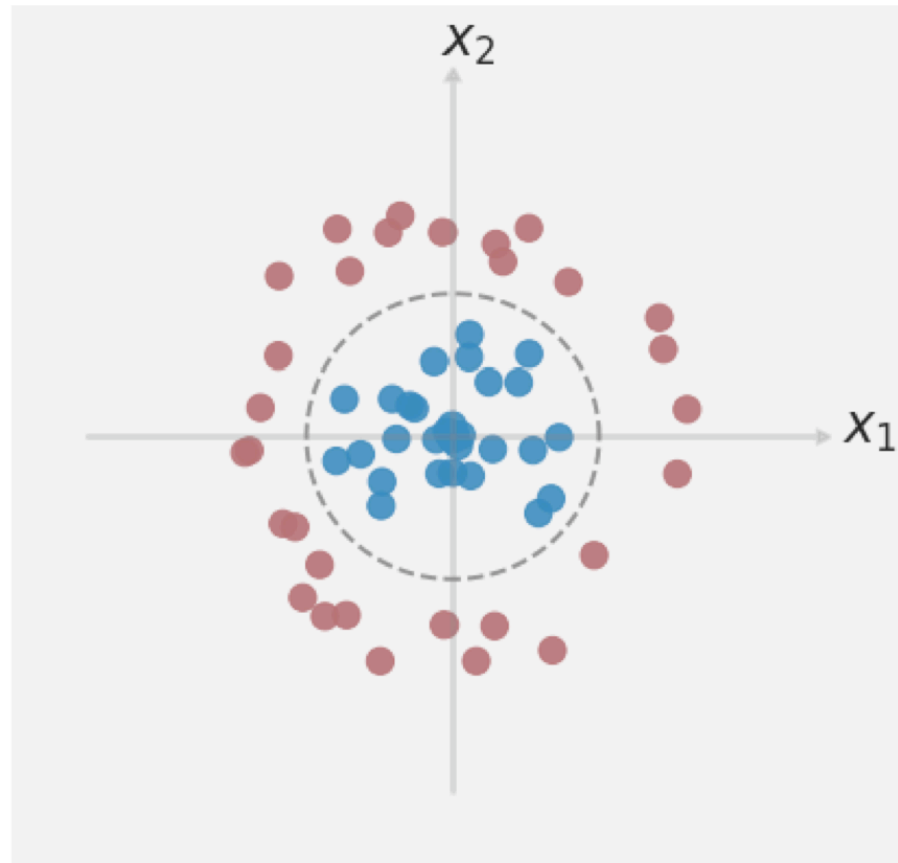
Derived features

We started with the original feature vector, $\mathbf{x} = (x_1, x_2)$, and we created a new derived feature vector, $\phi(\mathbf{x}) = (x_1, x_2, x_1^2 + x_2^2)$.



Derived features

We started with the original feature vector, $\mathbf{x} = (x_1, x_2)$,
and we created a new derived feature vector, $\phi(\mathbf{x}) = (x_1, x_2, x_1^2 + x_2^2)$.



What's special about SVMs?

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

What's special about SVMs?

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

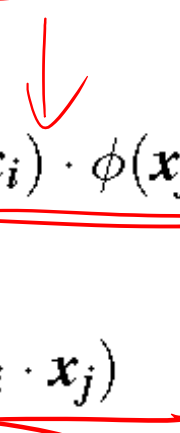
- This dot product is basically just how much x_i looks like x_j . Can we generalize that?

What's special about SVMs?

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

- This dot product is basically just how much x_i looks like x_j . Can we generalize that?
- Kernels!

What's special about SVMs?

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


What does the kernel trick buy us?

Polynomial kernel:

$$K(\mathbf{x}, \mathbf{x}') = (\underbrace{\mathbf{x} \cdot \mathbf{x}'}_{\mathbf{x}=(x_1, x_2)} + c)^d$$

$$\phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$$

$$(\mathbf{x} \cdot \mathbf{x}' + 1)^2$$

$$\phi(\mathbf{x})$$

$$\phi(\mathbf{x}) = (x_1, x_2, 1)$$

$$= (x_1 x_1' + x_2 x_2' + 1)^2$$

$$= x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x_1' x_2 x_2' + 2x_1 x_1' + 2x_2 x_2' + 1$$

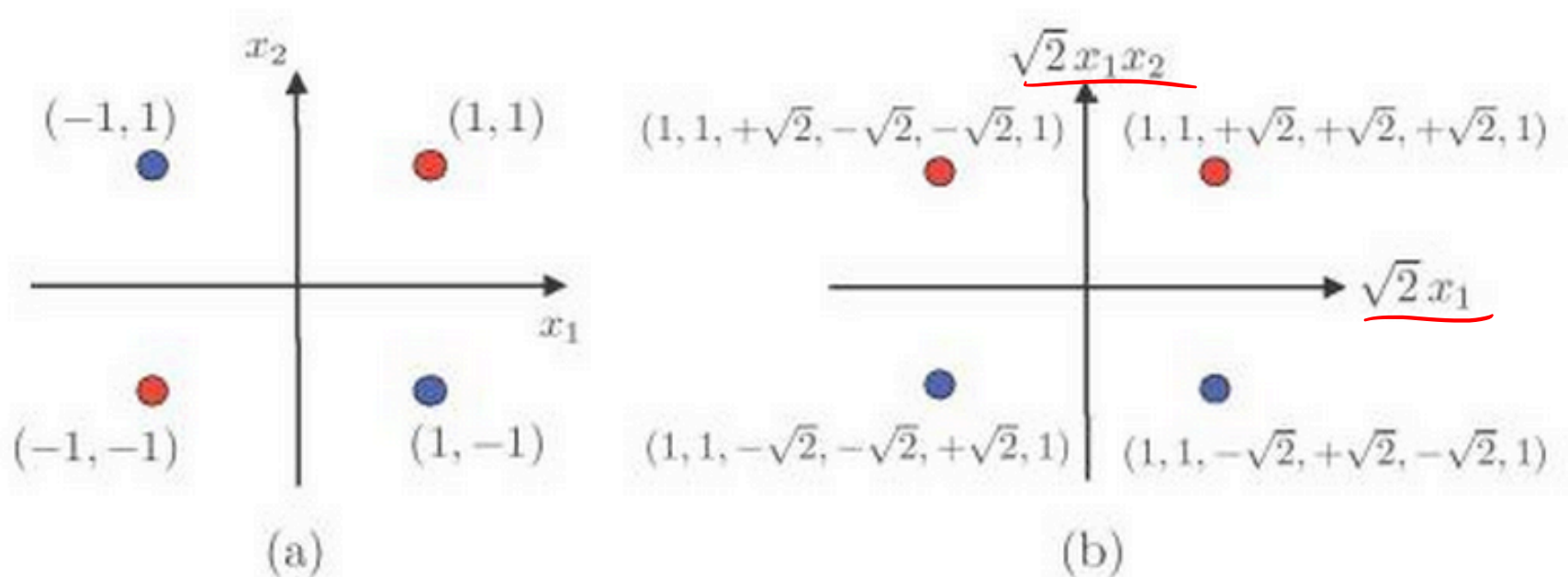
$$\phi(\mathbf{x}) = (\underbrace{x_1^2, x_2^2, \sqrt{2}x_1 x_2}, \underbrace{\sqrt{2}x_1, \sqrt{2}x_2, 1})$$

What does the kernel trick buy us?

Polynomial kernel:

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + c)^d$$

When $d = 2, c = 1$:



What does the kernel trick buy us?

Polynomial kernel:

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + 1)^2$$


What is the corresponding $\phi(\mathbf{x})$, where $\mathbf{x} \in \mathbb{R}^k$?

What is the complexity of storing $\phi(\mathbf{x})$ and computing $\phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$?

What about using the kernel function?

What's a kernel?

- A function $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is a kernel over \mathcal{X} .
- This is equivalent to taking the dot product $\langle \phi(x_1), \phi(x_2) \rangle$ for some mapping
- **Mercer's Theorem:** So long as the function is continuous and symmetric, then K admits an expansion of the form

$$K(x, x') = \sum_{n=0}^{\infty} a_n \phi_n(x) \phi_n(x')$$


What's a kernel?

- A function $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is a kernel over \mathcal{X} .
- This is equivalent to taking the dot product $\langle \phi(x_1), \phi(x_2) \rangle$ for some mapping
- **Mercer's Theorem:** So long as the function is continuous and symmetric, then K admits an expansion of the form

$$K(x, x') = \sum_{n=0}^{\infty} a_n \phi_n(x) \phi_n(x')$$

- The computational cost is just in computing the kernel

Kernel Matrix

The important property of the kernel matrix $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$ is symmetric positive semidefinite.

Kernel Matrix

The important property of the kernel matrix $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$ is symmetric positive semidefinite.

$$\underline{\mathbf{K}^T = \mathbf{K}}$$

Kernel Matrix

The important property of the kernel matrix $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$ is symmetric positive semidefinite.

$$\mathbf{K}^T = \mathbf{K}$$

$$\underbrace{\forall \mathbf{x}, \mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0}_{\sum K_i}$$

Kernel Matrix

The important property of the kernel matrix $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$ is symmetric positive semidefinite.

$$\mathbf{K}^T = \mathbf{K}$$

$$\forall \mathbf{x}, \mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$$

Also known as Gram matrix.

Gaussian Kernel

$$K(x, x') = \exp \left(-\frac{\|x' - x\|^2}{2\sigma^2} \right)$$

Gaussian Kernel

$$K(x, x') = \exp \left(-\frac{\|x' - x\|^2}{2\sigma^2} \right)$$

which can be rewritten as

$$K(x, x') = \sum_n \frac{(x \cdot x')^n}{\sigma^n n!}$$

(All polynomials!)

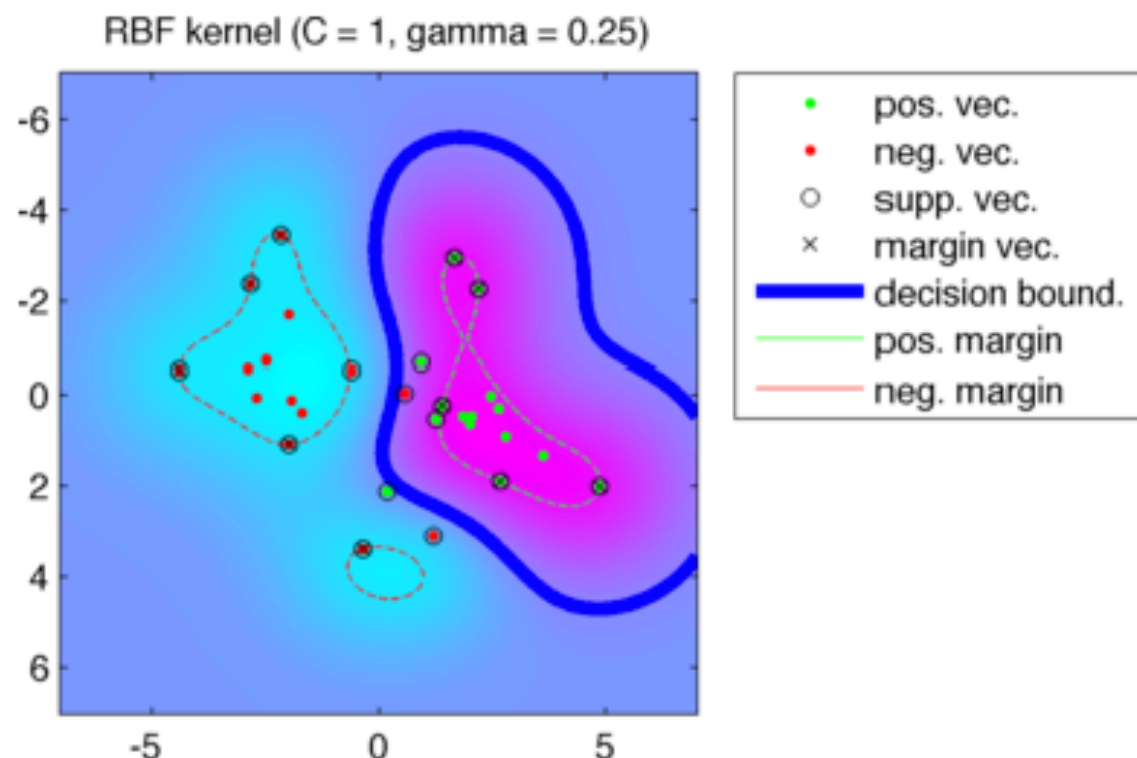
Gaussian Kernel

$$K(x, x') = \exp \left(-\frac{\|x' - x\|^2}{2\sigma^2} \right)$$

which can be rewritten as

$$K(x, x') = \sum_n \frac{(x \cdot x')^n}{\sigma^n n!}$$

(All polynomials!)



How does it affect optimization

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

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

- Replace all dot product with kernel evaluations $K(x_1, x_2)$
- Makes computation more expensive, overall structure is the same

Examples

- Switch to notebooks

Recap

- This completes our discussion of SVMs
- Workhorse method of machine learning
- Flexible, fast, effective