

STOR566: Introduction to Deep Learning

Lecture 5: Kernel Methods

Yao Li
UNC Chapel Hill

Sept 5, 2024

Materials are from *Learning from data* (Caltech) and *Deep Learning* (UCLA)

Outline

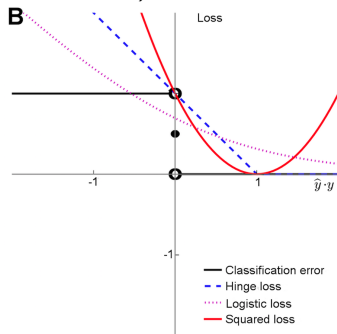
- Linear Support Vector Machines
- Nonlinear SVM, Kernel methods

Support Vector Machines

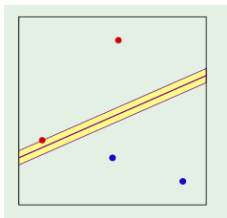
- Given training examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
Consider binary classification: $y_i \in \{+1, -1\}$
- Linear Support Vector Machine (SVM):

$$\arg \min_{\mathbf{w}} \frac{C}{N} \sum_{i=1}^N \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0) + \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

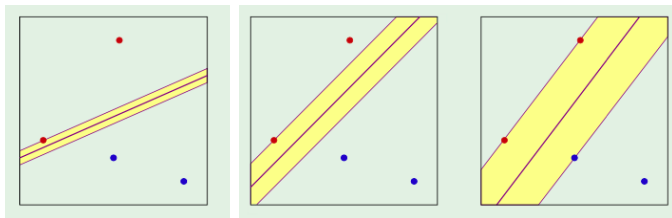
(hinge loss with L2 regularization)



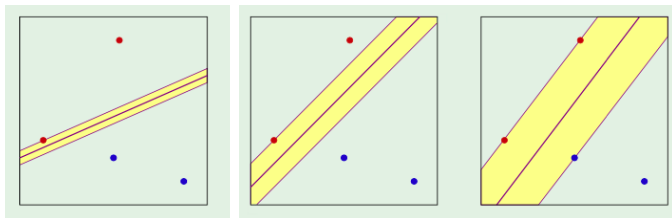
Linear Separation



Linear Separation

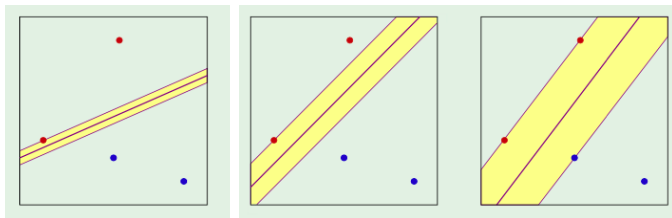


Linear Separation



- Which line is the best?

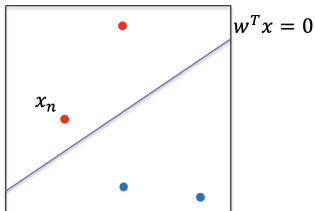
Linear Separation



- Which line is the best?
- Why big margin?
- Which \mathbf{w} maximizes the margin?

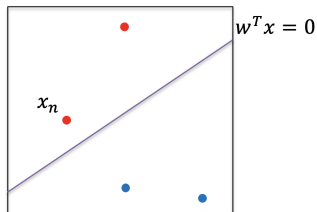
Margin

- $\mathbf{w}^T \mathbf{x} = 0$: the separation line or hyperplane
- \mathbf{x}_n : the nearest data point to the plane



Margin

- $\mathbf{w}^T \mathbf{x} = 0$: the separation line or hyperplane
- \mathbf{x}_n : the nearest data point to the plane



Preliminary:

- Normalize \mathbf{w} :

$$\|\mathbf{w}^T \mathbf{x}_n\| = 1$$

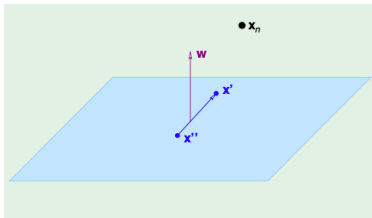
Distance

- The distance between \mathbf{x}_n and the plane $\mathbf{w}^T \mathbf{x} = 0$.
- The vector \mathbf{w} is orthogonal to the plane in the \mathcal{X} space

Take \mathbf{x}' and \mathbf{x}'' on the plane

$$\mathbf{w}^T \mathbf{x}' = 0, \mathbf{w}^T \mathbf{x}'' = 0$$

$$\implies \mathbf{w}^T (\mathbf{x}' - \mathbf{x}'') = 0$$



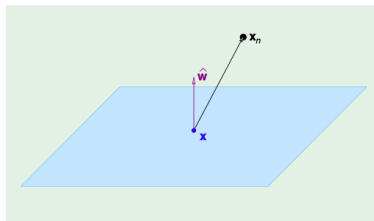
Distance

- The distance between \mathbf{x}_n and the plane $\mathbf{w}^T \mathbf{x} = 0$:

Take any point \mathbf{x} on the plane

Project $\mathbf{x}_n - \mathbf{x}$ on \mathbf{w} :

$$\text{distance} = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T (\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x}| = \frac{1}{\|\mathbf{w}\|}$$



Optimization Problem

- The optimization problem for SVM:

$$\begin{aligned} \max_{\mathbf{w}} \quad & \frac{1}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & \min_{i=1,\dots,N} |\mathbf{w}^T \mathbf{x}_i| = 1, \end{aligned}$$

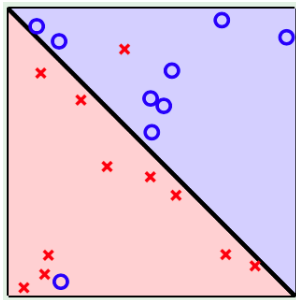
- Equivalent to:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_i \mathbf{w}^T \mathbf{x}_i \geq 1, i = 1, \dots, N, \end{aligned}$$

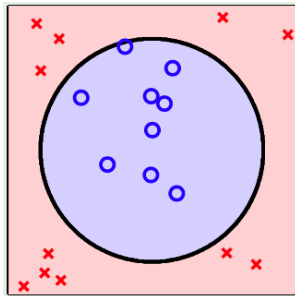
Notice: $|\mathbf{w}^T \mathbf{x}_i| = y_i \mathbf{w}^T \mathbf{x}_i$

Two Types of Non-separable

slightly:



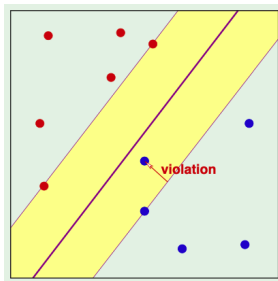
seriously:



Support Vector Machines (Soft)

- Given training data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ with labels $y_i \in \{+1, -1\}$.
- SVM problem with soft constraints:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i, \\ & \xi_i \geq 0, i = 1, \dots, N \end{aligned}$$



Support Vector Machines

- SVM problem:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i) \geq 1 - \xi_i, \\ & \xi_i \geq 0, i = 1, \dots, N. \end{aligned}$$

- Equivalent to

$$\min_{\mathbf{w}} \quad \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{L2 regularization}} + \frac{C}{N} \sum_{i=1}^N \underbrace{\max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)}_{\text{hinge loss}}$$

SVM: Unconstrained

- Unconstrained optimization:

$$\min_{\mathbf{w}} \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{L2 regularization}} + \frac{C}{N} \sum_{i=1}^N \underbrace{\max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)}_{\text{hinge loss}}$$

- Equivalent to:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \xi_i \geq \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i), i = 1, \dots, N. \end{aligned}$$

SVM: Unconstrained

- Unconstrained optimization:

$$\min_{\mathbf{w}} \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{L2 regularization}} + \frac{C}{N} \sum_{i=1}^N \underbrace{\max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)}_{\text{hinge loss}}$$

- Equivalent to:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \xi_i \geq \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i), i = 1, \dots, N. \end{aligned}$$

- Equivalent to:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \xi_i \geq 1 - y_i \mathbf{w}^T \mathbf{x}_i, \xi_i \geq 0, i = 1, \dots, N. \end{aligned}$$

Stochastic Subgradient Method for SVM

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^N \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

- A subgradient of $\max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$:

$$\begin{cases} -y_i \mathbf{x}_i & \text{if } 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0 \\ \mathbf{0} & \text{if } 1 - y_i \mathbf{w}^T \mathbf{x}_i < 0 \\ \mathbf{0} & \text{if } 1 - y_i \mathbf{w}^T \mathbf{x}_i = 0 \end{cases}$$

- Stochastic Subgradient descent for SVM:

For $t = 1, 2, \dots$

Randomly pick an index i

If $y_i \mathbf{w}^T \mathbf{x}_i < 1$, then

$$\mathbf{w} \leftarrow \mathbf{w} - \eta_t (\mathbf{w} - C y_i \mathbf{x}_i)$$

Else (if $y_i \mathbf{w}^T \mathbf{x}_i \geq 1$):

$$\mathbf{w} \leftarrow \mathbf{w} - \eta_t \mathbf{w}$$

Kernel SVM

Non-linearly separable problems

- What if the data is not linearly separable?

Solution: map data \mathbf{x}_i to higher dimensional(maybe infinite) feature space $\varphi(\mathbf{x}_i)$, where they are linearly separable.

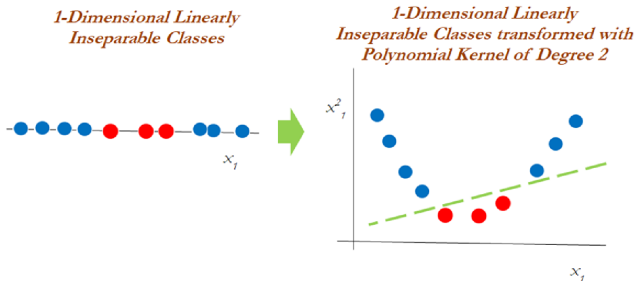
- Example: $\varphi(x) = (x, x^2)^T$

Non-linearly separable problems

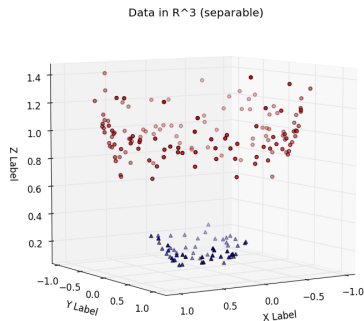
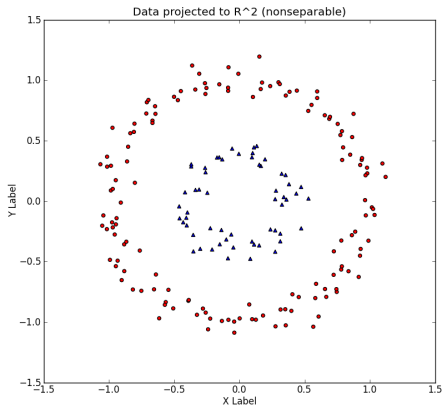
- What if the data is not linearly separable?

Solution: map data \mathbf{x}_i to higher dimensional (maybe infinite) feature space $\varphi(\mathbf{x}_i)$, where they are linearly separable.

- Example: $\varphi(\mathbf{x}) = (x, x^2)^T$



Non-linearly separable problems



- $\varphi(\mathbf{x}) = \varphi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix}$

SVM with nonlinear mapping

- SVM with nonlinear mapping $\varphi(\cdot)$:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{N} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \varphi(\mathbf{x}_i)) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, n, \end{aligned}$$

- Hard to solve if $\varphi(\cdot)$ maps to **very high or infinite dimensional space**

Support Vector Machines

- Primal problem:

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

$$\text{s.t. } y_i \mathbf{w}^T \varphi(\mathbf{x}_i) - 1 + \xi_i \geq 0, \text{ and } \xi_i \geq 0 \quad \forall i = 1, \dots, N$$

- Equivalent to (Dual problem):

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_i \xi_i - \sum_i \alpha_i (y_i \mathbf{w}^T \varphi(\mathbf{x}_i) - 1 + \xi_i) - \sum_i \beta_i \xi_i$$

Support Vector Machines (dual)

- Reorganize the equation:

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i y_i \mathbf{w}^T \varphi(\mathbf{x}_i) + \sum_i \xi_i \left(\frac{C}{N} - \alpha_i - \beta_i \right) + \sum_i \alpha_i$$

- For any given α, β , the minimizer will satisfy

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \varphi(\mathbf{x}_i) = 0 \Rightarrow \mathbf{w}^* = \sum_i y_i \alpha_i \varphi(\mathbf{x}_i)$$

$$\frac{\partial L}{\partial \xi} = \left(\frac{C}{N} - \alpha_1 - \beta_1, \dots, \frac{C}{N} - \alpha_N - \beta_N \right)^T = 0 \Rightarrow \frac{C}{N} = \alpha_i + \beta_i, \forall i$$

- Substitute these two equations back we get

$$\max_{\alpha \geq 0, \beta \geq 0, \frac{C}{N} = \alpha + \beta} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j) + \sum_i \alpha_i$$

Support Vector Machines (dual)

- Therefore, we get the following dual problem

$$\max_{0 \leq \alpha \leq \frac{c}{N}} \left\{ -\frac{1}{2} \alpha^T Q \alpha + \mathbf{1}^T \alpha \right\} := D(\alpha),$$

where Q is an N by N matrix with $Q_{ij} = y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$

- Based on the derivations,
 - ① We can solve the dual problem instead of the primal problem.
 - ② Let α^* be the dual solution and \mathbf{w}^* be the primal solution, we have

$$\mathbf{w}^* = \sum_i y_i \alpha_i^* \varphi(\mathbf{x}_i)$$

- To solve the dual, we only need to know $\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$.

Kernel Trick

- Do **not** directly define $\varphi(\cdot)$
- Instead, define “kernel”

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$$

This is all we need to know for Kernel SVM!

- Examples:
 - Gaussian kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}$
 - Polynomial kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i^T \mathbf{x}_j + c)^d$

The Trick

- Can we compute $K(\mathbf{x}, \mathbf{x}')$ **without** transforming \mathbf{x} and \mathbf{x}' ?
- Example:

Consider a transformation $\varphi(\mathbf{x}) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$

The inner product between $\varphi(\mathbf{x})$ and $\varphi(\mathbf{x}')$

$$K(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$$

- Computation cost: $(1 + \mathbf{x}^T \mathbf{x}')^2$ vs. $\varphi(\mathbf{x})^T \varphi(\mathbf{x}')$

The Trick

- Can we compute $K(\mathbf{x}, \mathbf{x}')$ **without** transforming \mathbf{x} and \mathbf{x}' ?
- Example:

Consider a transformation $\varphi(\mathbf{x}) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$

The inner product between $\varphi(\mathbf{x})$ and $\varphi(\mathbf{x}')$

$$K(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$$

- Computation cost: $(1 + \mathbf{x}^T \mathbf{x}')^2$ vs. $\varphi(\mathbf{x})^T \varphi(\mathbf{x}')$
- Example: Simple one-dimensional Gaussian kernel maps to infinite-dimensional space

$$\begin{aligned} K(x_i, x_j) &= \exp\left(-\frac{1}{2}(x_i - x_j)^2\right) \\ &= \exp\left(-\frac{1}{2}x_i^2\right) \exp\left(-\frac{1}{2}x_j^2\right) \sum_{k=0}^{\infty} \frac{x_i^k x_j^k}{k!} \end{aligned}$$

Solution

- Training: compute $\alpha = [\alpha_1, \dots, \alpha_N]$ by solving the quadratic optimization problem:

$$\min_{0 \leq \alpha \leq C} \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha$$

where $Q_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$

Solution

- Training: compute $\alpha = [\alpha_1, \dots, \alpha_N]$ by solving the **quadratic optimization problem**:

$$\min_{0 \leq \alpha \leq C} \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha$$

where $Q_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$

- Prediction: for a test data \mathbf{x} ,

$$\begin{aligned} \mathbf{w}^T \varphi(\mathbf{x}) &= \sum_{i=1}^N y_i \alpha_i \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}) \\ &= \sum_{i=1}^N y_i \alpha_i K(\mathbf{x}_i, \mathbf{x}) \end{aligned}$$

Conclusions

- SVM, Kernel SVM

Questions?