

ECE595 / STAT598: Machine Learning I

Lecture 21 Support Vector Machine: Soft & Kernel

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

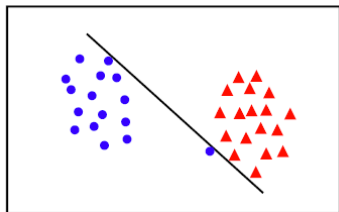
Support Vector Machine

- Lecture 19 SVM 1: The Concept of Max-Margin
- Lecture 20 SVM 2: Dual SVM
- **Lecture 21 SVM 3: Soft SVM and Kernel SVM**

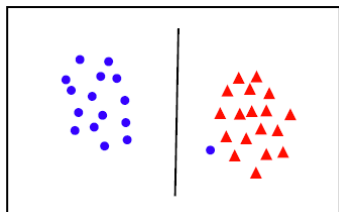
This lecture: Support Vector Machine: Soft and Kernel

- Soft SVM
 - Motivation
 - Formulation
 - Interpretation
- Kernel Trick
 - Nonlinearity
 - Dual Form
 - Kernel SVM

Linearly Not Separable



- the points can be linearly separated but there is a very narrow margin



- but possibly the large margin solution is better, even though one constraint is violated

<http://www.robots.ox.ac.uk/~az/lectures/ml/lect2.pdf>

Soft Margin

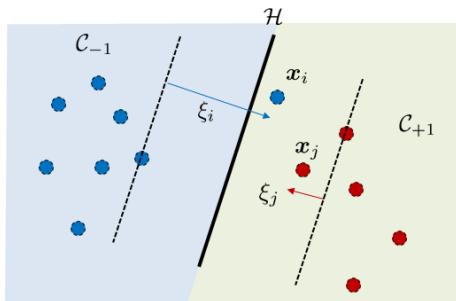
- We want to allow data points to stay inside the margin.
- How about change

$$y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1$$

to this one:

$$y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1 - \xi_j, \quad \text{and} \quad \xi_j \geq 0.$$

- If $\xi_j > 1$, then \mathbf{x}_j will be misclassified.



Soft Margin

- We can consider this problem

$$\begin{aligned} & \underset{\mathbf{w}, w_0, \xi}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{subject to} && y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1 - \xi_j, \\ & && \xi_j \geq 0, \quad \text{for } j = 1, \dots, n, \end{aligned}$$

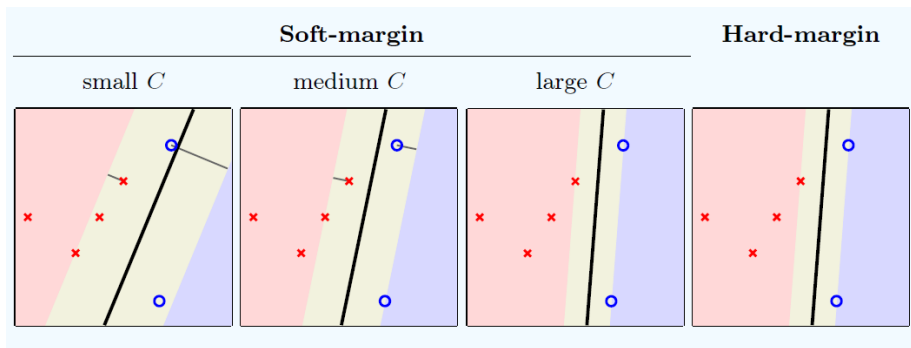
- But we need to control ξ , for otherwise the solution will be $\xi = \infty$.
- How about this:

$$\begin{aligned} & \underset{\mathbf{w}, w_0, \xi}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 + C \|\xi\|^2 \\ & \text{subject to} && y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1 - \xi_j, \\ & && \xi_j \geq 0, \quad \text{for } j = 1, \dots, n, \end{aligned}$$

- Control the energy of ξ .

Role of C

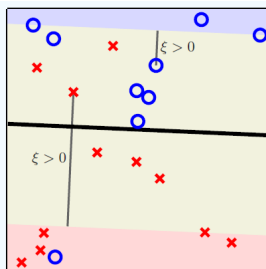
- If C is big, then we enforce ξ to be small.
- If C is small, then ξ can be big.



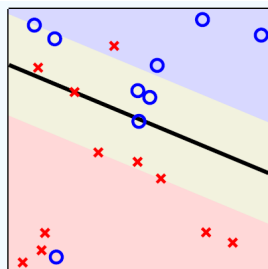
No Misclassification?

- You can have misclassification in soft SVM
- ξ_j can be big for a few outliers

$$\begin{aligned} & \underset{\mathbf{w}, w_0, \xi}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 + C \|\xi\|^2 \\ & \text{subject to} && y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1 - \xi_j, \\ & && \xi_j \geq 0, \quad \text{for } j = 1, \dots, N. \end{aligned}$$



(a) $C = 1$



(b) $C = 500$

L1 Regularization

- Instead of ℓ_1 -norm, you can also do

$$\begin{aligned} & \underset{\mathbf{w}, w_0, \boldsymbol{\xi}}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 + C \|\boldsymbol{\xi}\|_1 \\ & \text{subject to} && y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1 - \xi_j, \\ & && \xi_j \geq 0, \quad \text{for } j = 1, \dots, N. \end{aligned}$$

- This enforces $\boldsymbol{\xi}$ to be sparse.
- Only a few entries samples are allowed to live in the margin.
- The problem remains convex.
- So you can still use CVX to solve the problem.

Connection with Perceptron Algorithm

- In soft-margin SVM, $\xi_j \geq 0$ and $y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \geq 1 - \xi_j$ imply that
$$\xi_j \geq 0, \quad \text{and} \quad \xi_j \geq 1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0).$$

- We can combine them to get

$$\begin{aligned}\xi_j &\geq \max \left\{ 0, 1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \right\} \\ &= \left[1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \right]_+\end{aligned}$$

- So if we use SVM with ℓ_1 penalty, then

$$\begin{aligned}J(\mathbf{w}, w_0, \boldsymbol{\xi}) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{j=1}^N \xi_j \\ &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{j=1}^N \left[1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \right]_+\end{aligned}$$

Connection with Perceptron Algorithm

- This means that the training loss is

$$J(\mathbf{w}, w_0) = \sum_{j=1}^N \left[1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \right]_+ + \frac{\lambda}{2} \|\mathbf{w}\|_2^2,$$

if we define $\lambda = 1/C$.

- Now, you can make $\lambda \rightarrow 0$. This means $C \rightarrow \infty$
- Then,

$$\begin{aligned} J(\mathbf{w}, w_0) &= \sum_{j=1}^N \left[1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \right]_+ \\ &= \sum_{j=1}^N \max \left\{ 0, 1 - y_j(\mathbf{w}^T \mathbf{x}_j + w_0) \right\} \\ &= \sum_{j=1}^N \max \left\{ 0, 1 - y_j g(\mathbf{x}_j) \right\} \end{aligned}$$

Connection with Perceptron Algorithm

- **SVM Loss:**

$$J(\mathbf{w}, w_0) = \sum_{j=1}^N \max\{0, 1 - y_j g(\mathbf{x}_j)\}$$

- **Perceptron Loss:**

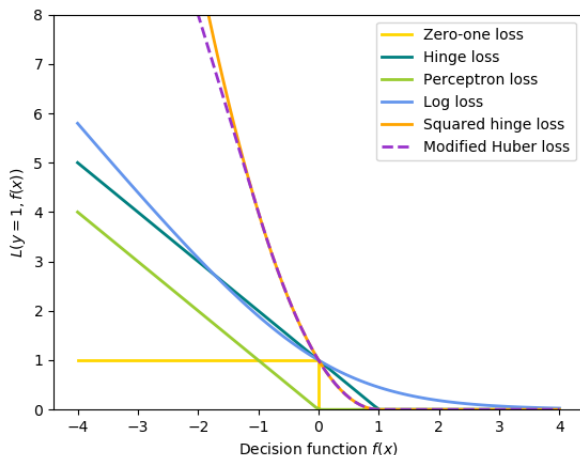
$$J(\mathbf{w}, w_0) = \sum_{j=1}^N \max\{0, -y_j g(\mathbf{x}_j)\}$$

- Therefore: SVM generalizes perceptron by allowing

$$J(\mathbf{w}, w_0) = \sum_{j=1}^N \max\{0, 1 - y_j g(\mathbf{x}_j)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- $\|\mathbf{w}\|_2^2$ regularizes the solution.

Comparing Loss functions



https://scikit-learn.org/dev/auto_examples/linear_model/plot_sgd_loss_functions.html

Outline

Support Vector Machine

- Lecture 19 SVM 1: The Concept of Max-Margin
- Lecture 20 SVM 2: Dual SVM
- **Lecture 21 SVM 3: Soft SVM and Kernel SVM**

This lecture: Support Vector Machine: Soft and Kernel

- Soft SVM
 - Motivation
 - Formulation
 - Interpretation
- Kernel Trick
 - Nonlinearity
 - Dual Form
 - Kernel SVM

The Kernel Trick

- A trick to turn linear classifier to nonlinear classifier.
- Dual SVM

$$\begin{aligned} & \underset{\lambda \geq 0}{\text{maximize}} && -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{j=1}^n \lambda_j \\ & \text{subject to} && \sum_{j=1}^n \lambda_j y_j = 0. \end{aligned}$$

- Kernel Trick

$$\begin{aligned} & \underset{\lambda \geq 0}{\text{maximize}} && -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) + \sum_{j=1}^n \lambda_j \\ & \text{subject to} && \sum_{j=1}^n \lambda_j y_j = 0. \end{aligned}$$

- You have to do this in dual. Primal is hard. See next slide.

The Kernel Trick

- Define

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

- The matrix \mathbf{Q} is

$$\mathbf{Q} = \begin{bmatrix} y_1 y_1 \mathbf{x}_1^T \mathbf{x}_1 & \dots & y_1 y_N \mathbf{x}_1^T \mathbf{x}_N \\ y_2 y_1 \mathbf{x}_2^T \mathbf{x}_1 & \dots & y_2 y_N \mathbf{x}_2^T \mathbf{x}_N \\ \vdots & \vdots & \vdots \\ y_N y_1 \mathbf{x}_N^T \mathbf{x}_1 & \dots & y_N y_N \mathbf{x}_N^T \mathbf{x}_N \end{bmatrix}$$

- By Kernel Trick:

$$\mathbf{Q} = \begin{bmatrix} y_1 y_1 K(\mathbf{x}_1, \mathbf{x}_1) & \dots & y_1 y_N K(\mathbf{x}_1, \mathbf{x}_N) \\ y_2 y_1 K(\mathbf{x}_2, \mathbf{x}_1) & \dots & y_2 y_N K(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots \\ y_N y_1 K(\mathbf{x}_N, \mathbf{x}_1) & \dots & y_N y_N K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

Kernel

- The inner product $\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$ is called a **kernel**

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

- Second-Order Polynomial kernel

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^2.$$

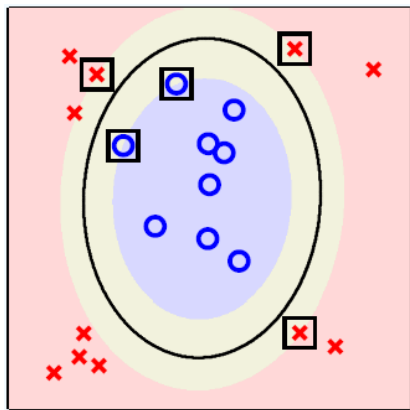
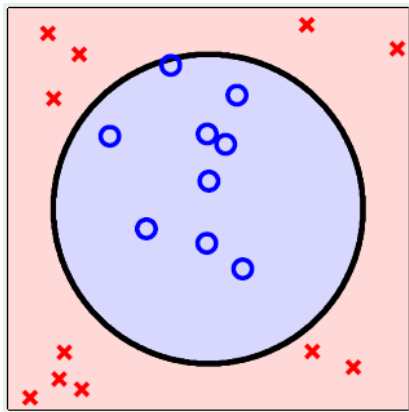
- Degree-Q Polynomial kernel

$$K(\mathbf{u}, \mathbf{v}) = (\gamma \mathbf{u}^T \mathbf{v} + c)^Q.$$

- Gaussian Radial Basis Function (RBF) Kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp \left\{ -\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2} \right\}.$$

SVM with Second Order Kernel

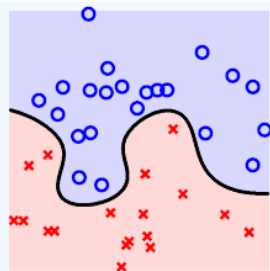


Boxed samples = Support vectors.

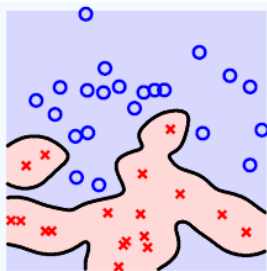
Radial Basis Function

Radial Basis Function takes the form of

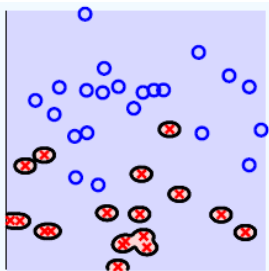
$$K(\mathbf{u}, \mathbf{v}) = \exp \{ -\gamma \|\mathbf{u} - \mathbf{v}\|^2 \}.$$



$$\exp(-1\|\mathbf{x} - \mathbf{x}'\|^2)$$



$$\exp(-10\|\mathbf{x} - \mathbf{x}'\|^2)$$



$$\exp(-100\|\mathbf{x} - \mathbf{x}'\|^2)$$

- Typical $\gamma \in [0, 1]$.
- γ too big: Over-fit.

Non-Linear Transform for RBF?

- Let us consider scalar $u \in \mathbb{R}$.

$$\begin{aligned}K(u, v) &= \exp\{-(u - v)^2\} \\&= \exp\{-u^2\} \exp\{2uv\} \exp\{-v^2\} \\&= \exp\{-u^2\} \left(\sum_{k=0}^{\infty} \frac{2^k u^k v^k}{k!} \right) \exp\{-v^2\} \\&= \exp\{-u^2\} \left(1, \sqrt{\frac{2^1}{1!}} u, \sqrt{\frac{2^2}{2!}} u^2, \sqrt{\frac{2^3}{3!}} u^3, \dots, \right)^T \\&\quad \times \left(1, \sqrt{\frac{2^1}{1!}} v, \sqrt{\frac{2^2}{2!}} v^2, \sqrt{\frac{2^3}{3!}} v^3, \dots, \right) \exp\{-v^2\}\end{aligned}$$

- So Φ is

$$\Phi(x) = \exp\{-x^2\} \left(1, \sqrt{\frac{2^1}{1!}} x, \sqrt{\frac{2^2}{2!}} x^2, \sqrt{\frac{2^3}{3!}} x^3, \dots, \right)$$

So You Need

Example. Radial Basis Function

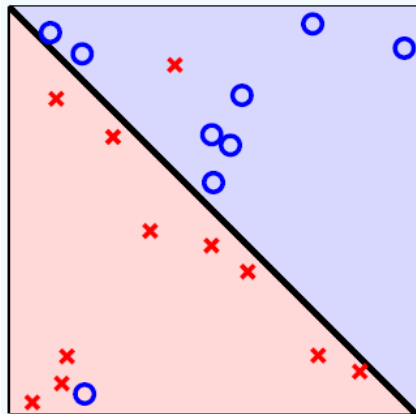
$$K(\mathbf{u}, \mathbf{v}) = \exp \{ -\gamma \|\mathbf{u} - \mathbf{v}\|^2 \}.$$

The non-linear transform is:

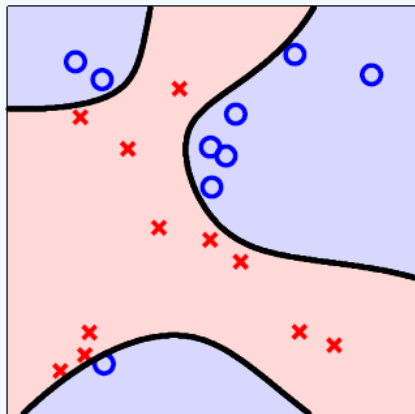
$$\Phi(x) = \exp\{-x^2\} \left(1, \sqrt{\frac{2^1}{1!}}x, \sqrt{\frac{2^2}{2!}}x^2, \sqrt{\frac{2^3}{3!}}x^3, \dots, \right)$$

- You need infinite dimensional non-linear transform!
- But to compute the kernel $K(\mathbf{u}, \mathbf{v})$ you do not need Φ .
- Another Good thing about Dual SVM: You can do infinite dimensional non-linear transform.
- Cost of computing $K(\mathbf{u}, \mathbf{v})$ is bottleneck by $\|\mathbf{u} - \mathbf{v}\|^2$.

Is RBF Always Better than Linear?



(a) linear classifier



(b) Gaussian-RBF kernel

- Noisy dataset: Linear works well.
- RBF: Over fit.

Testing with Kernels

- Recall:

$$\mathbf{w}^* = \sum_{n=1}^N \lambda_n^* y_n \mathbf{x}_n.$$

- The hypothesis function is

$$\begin{aligned} h(\mathbf{x}) &= \text{sign} \left(\mathbf{w}^{*T} \mathbf{x} + w_0^* \right) \\ &= \text{sign} \left(\left(\sum_{n=1}^N \lambda_n^* y_n \mathbf{x}_n \right)^T \mathbf{x} + w_0^* \right) \\ &= \text{sign} \left(\sum_{n=1}^N \lambda_n^* y_n \mathbf{x}_n^T \mathbf{x} + w_0^* \right). \end{aligned}$$

- Now you can replace $\mathbf{x}_n^T \mathbf{x}$ by $K(\mathbf{x}_n, \mathbf{x})$.

Support Vector Machine

- Mustafa, *Learning from Data*, e-Chapter
- Duda-Hart-Stork, *Pattern Classification*, Chapter 5.5
- Chris Bishop, *Pattern Recognition*, Chapter 7.1
- UCSD Statistical Learning
<http://www.svcl.ucsd.edu/courses/ece271B-F09/>