# SVM

Amr

April 1

# Outline

- Margin and VC-dimension
- The separable case
- Non- separable case:  Hing Loss
- Kernels

# SVM: Intuition

- Remember LR
  - Predict Y=1 if $\dfrac{1}{1+\exp(w.x+b)} > 0.5$
  - or if w.x > 0
  - The more prob. >>> 0.5, the more confident we are about our prediction
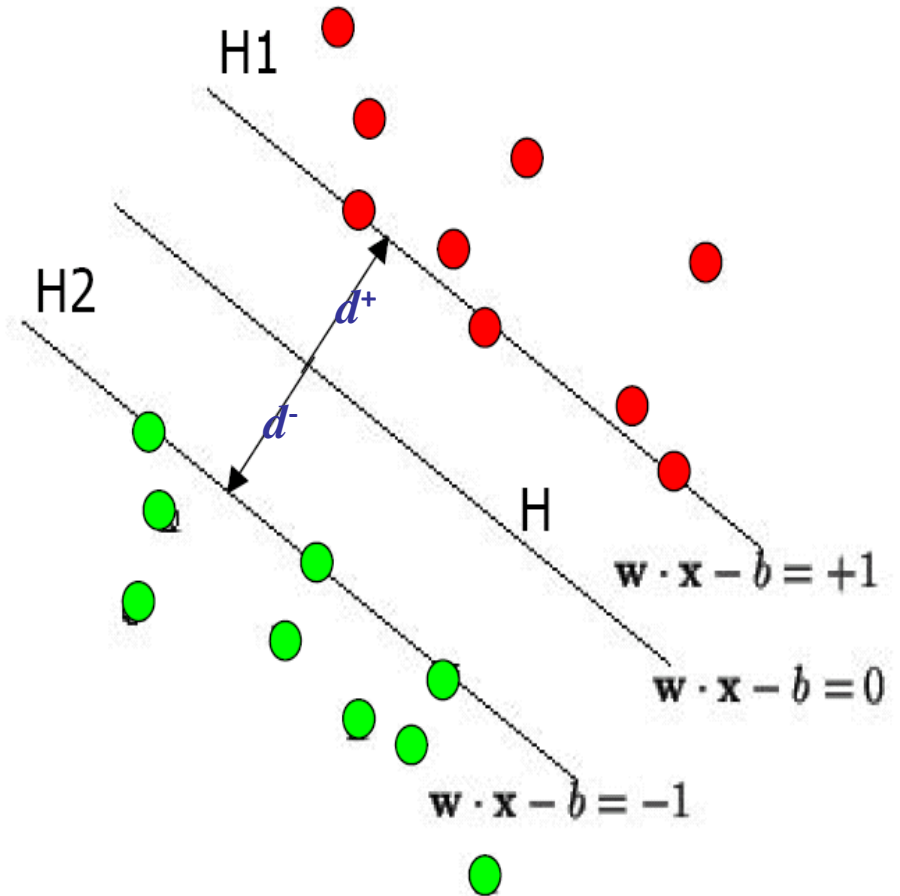  - Or the more w.x+b >>0 (margin), the more we are confident about our prediction
  - Same in SVM

# SVM

- SVM problem

$$\max_{w,b} \quad \frac{1}{\|w\|}$$

$$\text{s.t} \quad y_i(w^T x_i + b) \geq 1, \quad \forall i$$

- Or equivalently

$$\min_{w,b} \quad \frac{1}{2} w^t w$$

$$\text{s.t} \quad y_i(w^T x_i + b) \geq 1, \quad \forall i$$

# SVM using VC-dimension

## VC Theory

(Vapnik, 1982)

Given $x_1, ..., x_n \in \mathbb{R}^d$ iid and $||x_i||_2 \leq D$, if $\mathcal{H}_\gamma$ is the hypothesis space of linear classifiers in $\mathbb{R}^d$ with margin $\gamma$,

$$VC(\mathcal{H}_\gamma) \leq \min\left\{ d, \left\lceil \frac{4D^2}{\gamma^2} \right\rceil \right\}.$$

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln\frac{2m}{VC(H)} + 1) + \ln\frac{4}{\delta}}{m}}$$

# SVM using VC-dimension

- Thus large-margin → small VC-dim → better generalization bound

- Recall that d+1 is the upper bound for a linear classifier in d-space

$$VC(\mathcal{H}_\gamma) \leq \min \left\{ d, \left\lceil \frac{4D^2}{\gamma^2} \right\rceil \right\}.$$

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln \frac{2m}{VC(H)} + 1) + \ln \frac{4}{\delta}}{m}}$$

# Outline

- Margin and VC-dimension
- The separable case
- Kernels
- Non-seperable case: Hing Loss

# Solution Sketch

$$\min_{w,b} \quad \frac{1}{2} w^t w$$

$$\text{s.t} \qquad y_i(w^T x_i + b) \geq 1, \quad \forall i$$

- Form the Langrangian
- Optimize with respect to primal variable
- Subs. Into Lagrangian to get dual problem
- Exploit the KKT condition

# Lagrangian

$$\operatorname*{argmin}_{w,b} \frac{1}{2}||w||^2$$

$$s.t. \quad y_i(w \cdot x_i + b) \geq 1.$$

$$L(w, b, \alpha) = \frac{1}{2} w \cdot w - \sum_i \alpha_i [y_i(w \cdot x_i + b) - 1]$$

- One dual variable per constraints

$$L(w, b, \alpha) = \frac{1}{2} w \cdot w - \sum_i \alpha_i [y_i(w \cdot x_i + b) - 1]$$

$$\frac{\partial}{\partial w} L(w, b, \alpha) = w - \sum_i \alpha_i y_i x_i = 0 \ \rightarrow \ w = \sum_i \alpha_i y_i x_i.$$

$$\frac{\partial}{\partial b} L(w, b, \alpha) = \sum_i \alpha_i y_i = 0.$$

$$\underset{\alpha}{\text{argmax}} \ L(w, b, \alpha) = \underset{\alpha}{\text{argmax}} \ \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$$s.t. \ \ \alpha_i \geq 0,$$
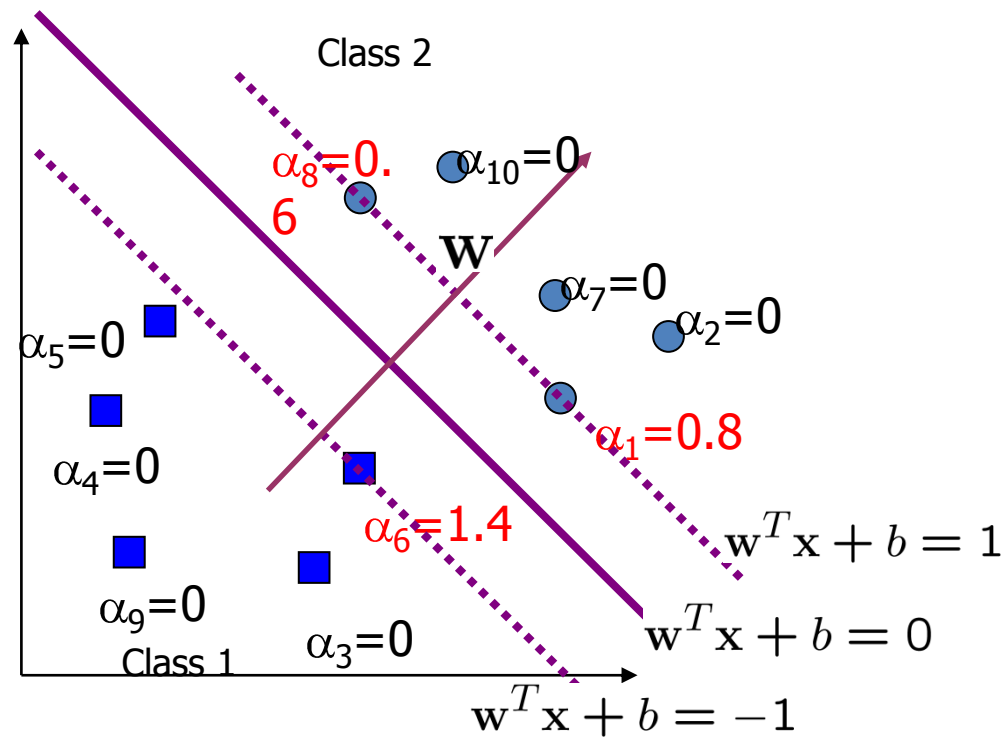
$$\sum_i \alpha_i y_i = 0.$$

# KKT conditions

- For convex objective and affine constraints we have

$$\alpha_i \left[ 1 - y_i \left( w.x_i + b \right) \right] = 0, \quad i = 1, \ldots, m$$

- Only a few $\alpha_i$ can be non-zero.

$$w = \sum_{i \in SV} \alpha_i y_i \mathbf{x}_i$$

- How about $b$?

$$y*(z) = \text{sign}\left( \sum_{i \in SV} \alpha_i y_i \left( x_i^T z \right) + b \right)$$



Class 2

$\alpha_8 = 0.6$  $\alpha_{10} = 0$

$\mathbf{W}$

$\alpha_7 = 0$  $\alpha_2 = 0$

$\alpha_5 = 0$

$\alpha_1 = 0.8$

$\alpha_4 = 0$

$\alpha_6 = 1.4$

$\mathbf{w}^T \mathbf{x} + b = 1$

$\alpha_9 = 0$

$\alpha_3 = 0$

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

Class 1

$\mathbf{w}^T \mathbf{x} + b = -1$
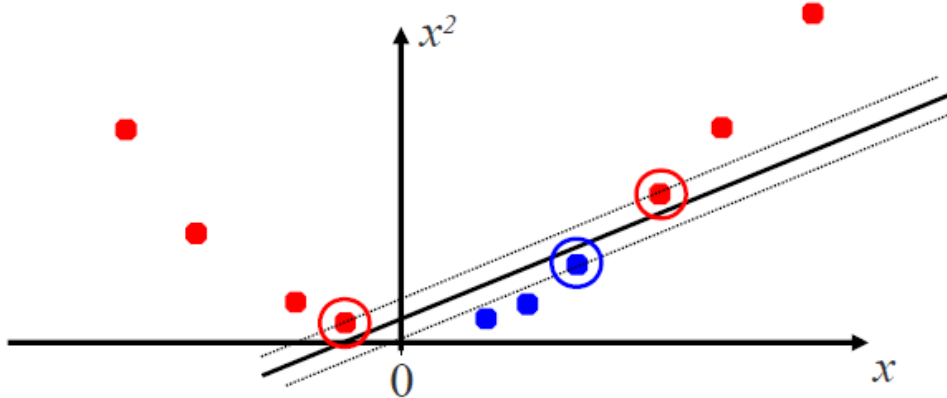
# Kernel Trick

- Is this data linearly-separable?



- How about a quadratic mapping $\phi(x_i)$?

# Kernel Trick

feature space

Input space



- Simply replace $x_i$ with $\phi(x_i)$ !

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

$$\text{s.t.} \quad \alpha_i \geq 0, \quad i = 1, \ldots, k$$

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

$$y^*(z) = \text{sign}\left( \sum_{i \in SV} \alpha_i y_i \phi(\mathbf{x}_i)^t \phi(z) + b \right)$$

- So what is the deal?

# Kernel Trick

- Computation depends on feature space
  - Bad if its dimension is much larger than input space

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

$$\text{s.t.} \quad \alpha_i \geq 0, \quad i = 1, \ldots, k$$

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

Where $K(x_i, x_j) = \phi(x_i)^t \phi(x_j)$

$$y^*(z) = \text{sign}\left( \sum_{i \in SV} \alpha_i y_i K(\mathbf{x}_i, z) + b \right)$$

# Example Kernel
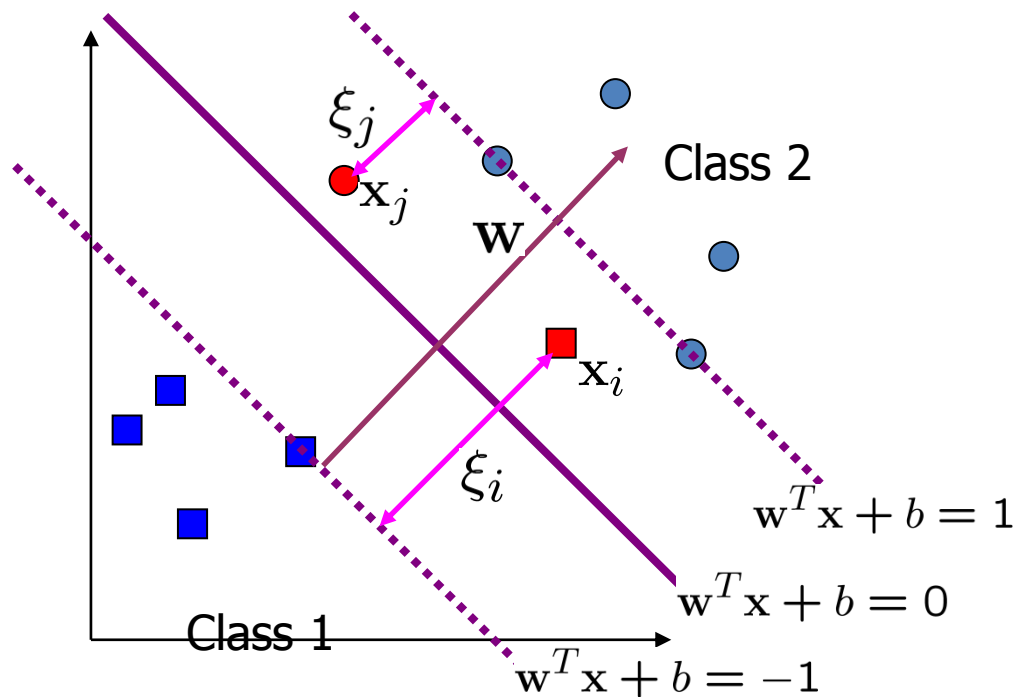
- $x_i$ is a bag of words
- Define $\phi(x_i)$ as a count of every n-gram up to n=k in $x_i$.
  - This is huge space $26^k$
  - What are we measuring by $\phi(x_i)^t \phi(x_j)$?
- Can we compute the same quantity on input space?
  - Efficient linear dynamic program!
- Kernel is a measure of similarity
- Must be positive semi-definite

# Outline

- Margin and VC-dimension
- The separable case
- Kernels
- <span style="color:red">Non-separable case: Hing Loss</span>

# Non-separable case



$$\min_{w,b} \quad \frac{1}{2} w^T w + C \sum_{i=1}^{m} \xi_i$$

$$\text{s.t} \quad y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \forall i$$

$$\xi_i \geq 0, \quad \forall i$$

# Remember Ridge regression

- Min [squared loss + $\lambda$ w$^t$w]
- How about SVM?

$$\text{argmin}_{\{w,b\}} \, w^t w + \lambda \sum_{1}^{m} \max(1 - y_i(w^t x_i + b), 0)$$

regularization                     Loss: hinge loss

$$\min_{w,b} \quad \frac{1}{2}w^T w + C\sum_{i=1}^{m}\xi_i$$

$$\text{s.t} \quad y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \forall i$$

$$\xi_i \geq 0, \quad \forall i$$

$$\xi_i \geq \max\left(0, 1 - y_i(w^T x_i + b)\right)$$

Why?

$$\xi_i = \max\left(0, 1 - y_i(w^T x_i + b)\right)$$

$$\arg\min_{\{w,b\}} w^t w + \lambda \sum_{1}^{m} \max(1 - y_i(w^t x_i + b), 0)$$

regularization

Loss: hinge loss