

SVM for Overlapping Classes

srihari@buffalo.edu

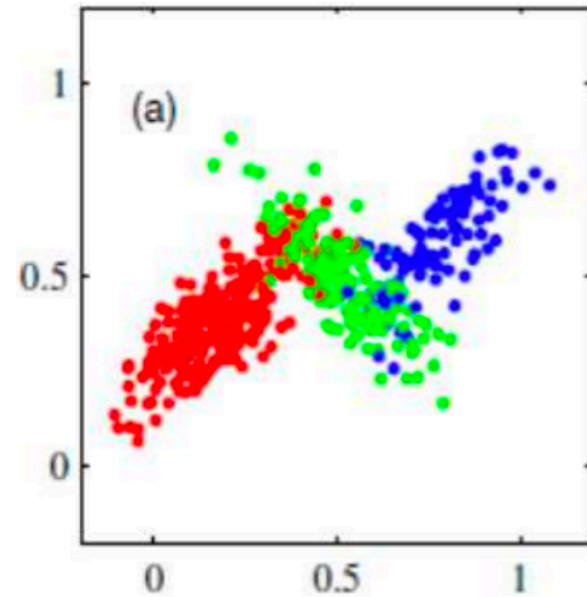
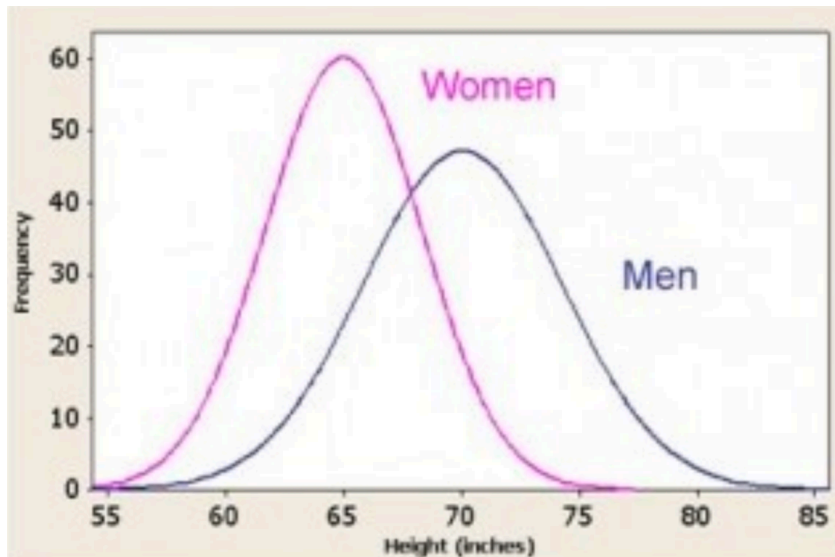
SVM Discussion Overview

1. Importance of SVMs
2. Overview of Mathematical Techniques Employed
3. Margin Geometry
4. SVM Training Methodology
5. Overlapping Distributions
6. Dealing with Multiple Classes
7. SVM and Computational Learning Theory
8. Relevance Vector Machines

Overlapping Class Distributions

- We have assumed training data are linearly separable in mapped space $\phi(\mathbf{x})$
 - Resulting SVM gives exact separation in input space \mathbf{x} although decision boundary will be nonlinear
- In practice class-conditionals overlap
 - Exact separation leads to poor generalization
 - Therefore need to allow SVM to misclassify some training points

Overlapping Distributions



Modifying SVM

- Implicit error function of SVM is

$$\sum_{n=1}^N E_{\infty} \left(y(\mathbf{x}_n) t_n - 1 \right) + \lambda \|\mathbf{w}\|^2$$

- Gives infinite error if point is misclassified
- Zero error if it was classified correctly
- We optimize parameters to maximize margin
- We now modify this approach
 - Points allowed to be on *wrong side* of margin
 - But with penalty that increases with the distance from the boundary
 - Convenient if penalty is a linear function of distance

Definition of Slack Variables ξ_n

- Introduce slack variables, $\xi_n \geq 0$, $n=1, \dots, N$
 - With one slack variable for each training point
- These are defined by
 - $\xi_n = 0$ for data points that are on or inside the correct margin boundary
 - $\xi_n = |t_n - \mathbf{y}(\mathbf{x}_n)|$ for other points

Slack Variables

- A data point that is on the decision boundary $y(\mathbf{x}_n)=0$ will have $\xi_n = 1$
- Points with $\xi_n \geq 1$ will be misclassified
- Exact classification constraints

$$t_n (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1, \quad n=1, \dots, N$$

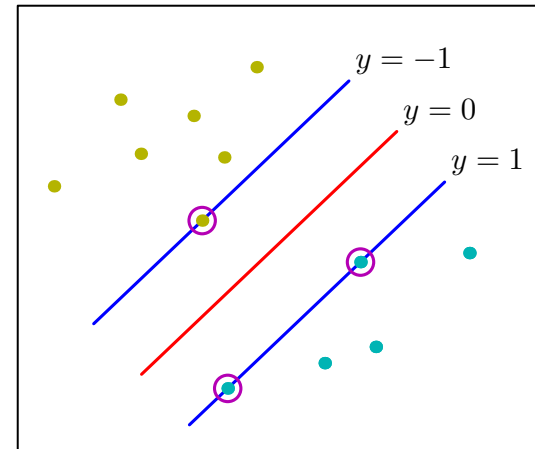
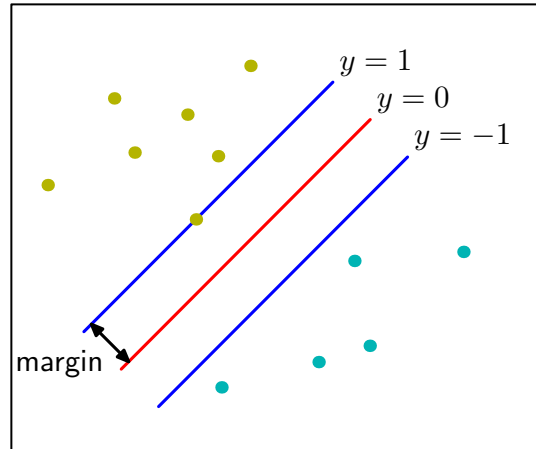
are then replaced with

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n \quad n=1, \dots, N$$

- in which slack variables are constrained to satisfy $\xi_n \geq 0$

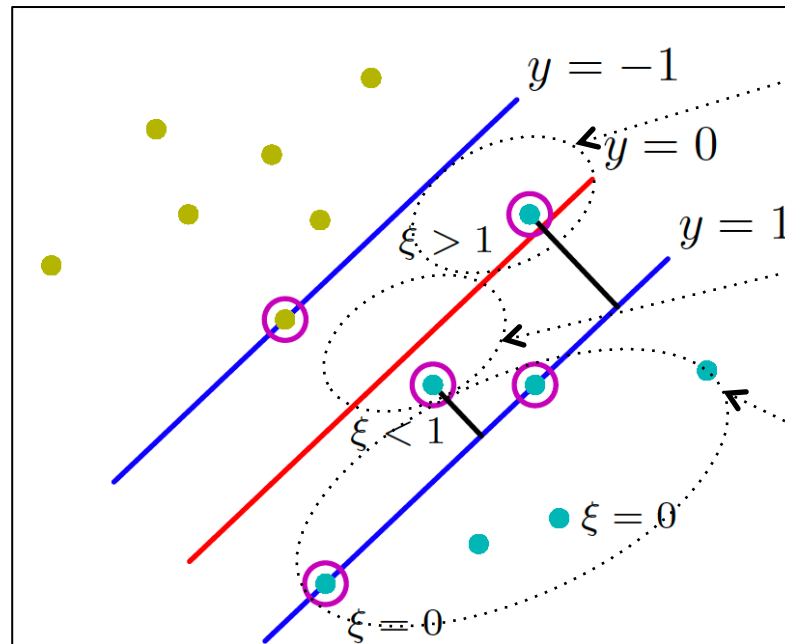
Slack Variables and Margins

Hard Margin



Data points with circles are support vectors

Soft Margin



Data points on wrong side of decision boundary have $\xi > 1$

Correctly classified points inside margin but on correct side of decision boundary have $0 < \xi \leq 1$

Correctly classified points on correct side of margin have $\xi = 0$

Soft Margin Classification

- Allows some of the data points to be misclassified
- While slack variables allow for overlapping class distributions, the framework is still sensitive to outliers
 - because the penalty for misclassification still increases linearly with ξ

Optimization with slack variables

- Our goal is to maximize the margin while softly penalizing points that lie on the wrong side of the margin boundary

- We therefore minimize $C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$

- Parameter $C > 0$ controls trade-off between slack variable penalty and the margin
- Subject to constraints

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n \quad n=1, \dots, N$$

Lagrangian with slack

- The Lagrangian we wish to maximize is

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n [t_n y(\mathbf{x}_n) - 1 + \xi_n] - \sum_{n=1}^N \mu_n \xi_n$$

- where $\{a_n \geq 0\}$, $\{\mu_n \geq 0\}$ are Lagrange multipliers
- The corresponding KKT conditions are

$$\begin{aligned} a_n &\geq 0 \\ t_n y(\mathbf{x}_n) - 1 + \xi_n &\geq 0 \\ a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} &= 0 \\ \mu_n &\geq 0 \\ \xi_n &\geq 0 \\ \mu_n \xi_n &= 0 \end{aligned}$$

- We then optimize out \mathbf{w}, b and $\{\xi_n\}$ by taking derivatives of L wrt \mathbf{w} , b and ξ_n to get

Dual Lagrangian

- The dual Lagrangian has the form

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

- Which is identical to the separable case, except that the constraints are somewhat different

Nu-SVM

- An equivalent formulation of SVM

- It involves maximizing

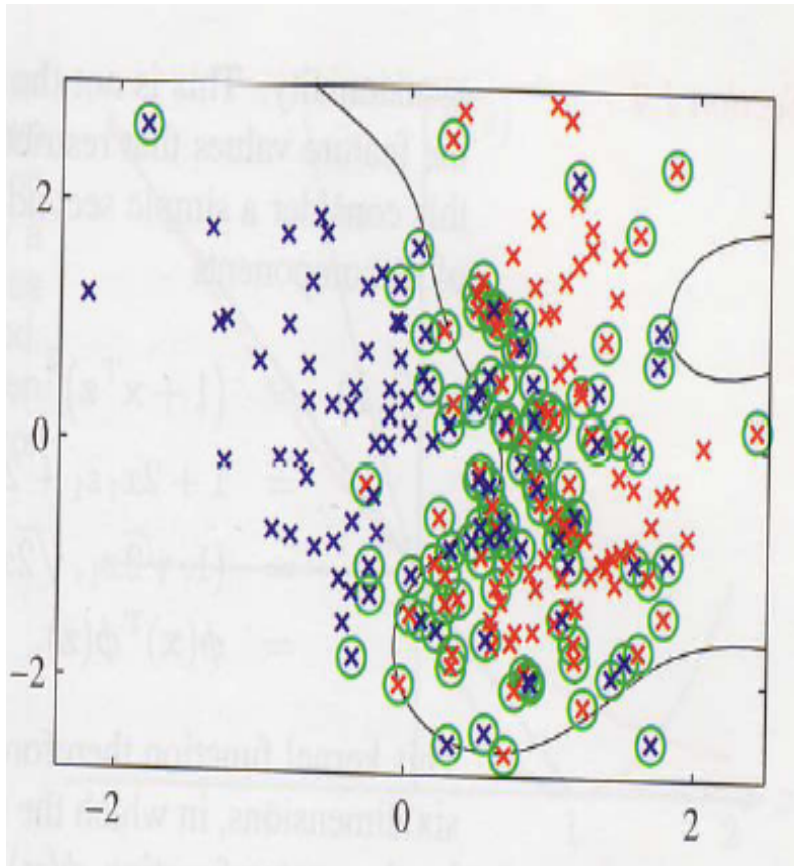
$$\tilde{L}(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

- Subject to the constraints

$$\begin{aligned} 0 &\leq a_n \leq 1/N \\ \sum_{n=1}^N a_n t_n &= 0 \\ \sum_{n=1}^N a_n &\geq \nu \end{aligned}$$

- ν -SVM has the advantage of using a parameter ν for controlling the no. of support vectors

ν -SVM applied to non-separable data



- Support Vectors are indicated by circles
- Done by introducing slack variables
- With one slack variable per training data point
- Maximize the margin while softly penalizing points that lie on the wrong side of the margin boundary
- ν is an upper-bound on the fraction of margin errors (lie on wrong side of margin boundary)