

Support Vector Machines

Dr. Fayyaz ul Amir Afsar Minhas

PIEAS Biomedical Informatics Research Lab Department of Computer and Information Sciences Pakistan Institute of Engineering & Applied Sciences PO Nilore, Islamabad, Pakistan <u>http://faculty.pieas.edu.pk/fayyaz/</u>

Preliminaries: Introduction to COP

 The constrained optimization problem (COP) can be expressed in its general form as follows

$$\begin{array}{ll}
\min_{\mathbf{x}} & f(\mathbf{x}) \\
\text{subject to} \\
g_i(\mathbf{x}) \leq 0 \qquad i = 1...m
\end{array}$$

 $h_i(\mathbf{x}) = 0 \qquad i = 1 \dots p$

• Example

Constrained Optimization: Example

- You are given a string of length L. You are to tie it around a rectangular gift box (I x w x h) with a constant or fixed width 'w' using any length of this string (up to L).
- Can you find the dimensions (length and height only since the width is fixed) of the box with the largest volume that you can tie with this string?
- Mathematically
 - Optimize the objective function
 - Max V=wlh
 - Subject to constraints
 - 2(l+h)≤L



3

© 2008, Fayyaz A. Afsar, DCIS, PIEAS.

CIS 621: Machine Learning

PIEAS Biomedical Informatics Research Lab

Lagrangian Formulation

 Lagrange proposed a method for the solution of COP

 $\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{subject to} \\ g_i(\mathbf{x}) \leq 0 & i = 1...m \\ h_i(\mathbf{x}) = 0 & i = 1...p \end{array}$

- f(x) and $g_i(x)$ are convex functions - $h_i(x)$ are affine functions

Lagrangian Formulation...

• The solution proposed by Lagrange is based on the following unconstrained minimization

$$\min_{\mathbf{x}} \quad f(\mathbf{x}) + \sum_{i=1}^{m} L_{i}^{g}(\mathbf{x}) + \sum_{i=1}^{p} L_{i}^{h}(\mathbf{x})$$

- Where $L_i^g(\mathbf{x})$ and $L_i^h(\mathbf{x})$ have the following properties $L_i^g(\mathbf{x}) = \begin{cases} \infty & g_i(\mathbf{x}) > 0\\ 0 & else \end{cases}, \quad i = 1...m$ $L_i^h(\mathbf{x}) = \begin{cases} \infty & h_i(\mathbf{x}) \neq 0\\ 0 & else \end{cases}, \quad i = 1...p$
 - Large penalties added when the constraints are not satisfied
 - Unconstrained optimization now leads to satisfaction of the constrains and then optimization of the original objective function

Lagrangian Formulation...

 One possible way of achieving the above mentioned properties for the two penalty functions is as follows



Lagrangian Formulation...

• Thus we can write the optimization problem as

$$\min_{\mathbf{x}} f(\mathbf{x}) + \sum_{i=1}^{m} \max_{\alpha_i \ge 0} \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \max_{\beta_i} \beta_i h_i(\mathbf{x})$$

$$\equiv \min_{\mathbf{x}} f(\mathbf{x}) + \max_{\alpha_i \ge 0, \beta_i} \left(\sum_{i=1}^{m} \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \beta_i h_i(\mathbf{x}) \right)$$

$$\equiv \min_{\mathbf{x}} \max_{\alpha_i \ge 0, \beta_i} f(\mathbf{x}) + \left(\sum_{i=1}^{m} \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^{p} \beta_i h_i(\mathbf{x}) \right)$$

α_i and β_i are called Lagrange multipliers (or dual variables) and the function (below) is called the Lagrange Function

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{x}) + \left(\sum_{i=1}^{m} \alpha_{i} g_{i}(\mathbf{x}) + \sum_{i=1}^{p} \beta_{i} h_{i}(\mathbf{x})\right)$$

7

Lagrangian function: Gift example

- The problem can be rewritten as
 - Min $f(x) = -x_1x_2$
 - Subject to constraints $2(x_1+x_2) \le L$, OR $g(x) = 2(x_1+x_2) - L \le 0$

This implies

$$L(\mathbf{x}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \alpha_1 g_1(\mathbf{x})$$
$$= -x_1 x_2 + \alpha_1 (2(x_1 + x_2) - L))$$

Lagrangian function: example...

This can be solved as

$$\min_{\mathbf{x}} \max_{\alpha_i \ge 0} L(x_1, x_2, \alpha) = -x_1 x_2 + \alpha_1 \left(2(x_1 + x_2) - L \right)$$

$$\frac{\partial L}{\partial \alpha} = 2(x_1 + x_2) - L = 0 \qquad \Rightarrow \qquad x_1 + x_2 = \frac{L}{2}$$

$$\frac{\partial L}{\partial x_1} = -x_1 + 2\alpha_1 = 0 \qquad \Rightarrow \qquad x_1 = 2\alpha_1$$

$$\frac{\partial L}{\partial x_2} = -x_2 + 2\alpha_1 = 0 \qquad \Rightarrow \qquad x_2 = 2\alpha_1$$

$$\Rightarrow \qquad x_1 = x_2 = \frac{L}{4}$$

Classification

- Before moving on with the discussion let us restrict ourselves to the following problem
 - $-T = Given Training Set = {(\underline{x}^{(i)}, y_i), i = 1...N}$
 - <u>x</u>⁽ⁱ⁾ε R^m {Data Point i }
 - y_i: class of data point i (+1 or -1)



Use of Linear Discriminant in Classification

- Classifiers such as the Single Layer Perceptron (with linear activation function) and SVM use a linear discriminant function to differentiate between patterns of different classes
- The linear discriminant function is given by



11

Use of Linear Discriminant in Classification

 There are a large number of lines (or in general 'hyperplanes') separating the two classes



Use of Linear Discriminant in Classification



Margin of a linear classifier

 The width by which the boundary of a linear classifier can be increased before hitting a data point is called the margin of the linear classifier



Support Vector Machines (SVM)

- Support Vector Machines are linear classifiers that produce the optimal separating boundary (hyper-plane)
 - Find w and b in a way so as to maximize the margin while classifying all the training patterns correctly (for linearly separable problem)

Finding Margin of a Linear Classifier

• Consider a linear classifier with the boundary

 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$ for all x on the boundary

- We know that the vector w is perpendicular to the boundary
 - Consider two points x⁽¹⁾ and x⁽²⁾ on the boundary

$$f\left(\mathbf{x}^{(1)}\right) = \mathbf{w}^{T}\mathbf{x}^{(1)} + b = 0 \tag{1}$$
$$f\left(\mathbf{x}^{(2)}\right) = \mathbf{w}^{T}\mathbf{x}^{(2)} + b = 0 \tag{2}$$

Subtracting (1) from (2)

$$\mathbf{w}^{T}\left(\mathbf{x}^{(2)}-\mathbf{x}^{(1)}\right)=0 \qquad \Rightarrow \qquad \mathbf{w}\perp\left(\mathbf{x}^{(2)}-\mathbf{x}^{(1)}\right)$$

$$\begin{array}{c|c} \mathbf{x}_{2} \\ \mathbf{x}_{2} \\ \mathbf{x}_{2} \\ \mathbf{x}_{2} \\ \mathbf{x}_{2} \\ \mathbf{x}_{2} \\ \mathbf{x}_{1} \end{array}$$

 $x^{(1)}$

>0

Finding Margin of a Linear Classifier

 Let x^(s) be a point in the feature space with its projection x^(p) on the boundary

$$f\left(\mathbf{x}^{(p)}\right) = \mathbf{w}^T \mathbf{x}^{(p)} + b = 0$$

• We know that,



 X_2

 $\mathbf{X}^{(p)}$

CIS 621: Machine Learning

 $\mathbf{X}^{(s)}$

Example

- Consider the line
 - $x_1 + 2x_2 + 3 = 0$
- The distance of (4,2) is
 r = 4.92



Finding Margin of a Linear Classifier

- The points in the training set that lie closest (having minimum perpendicular distance) to the separating hyper-plane are called support vectors
- Margin is twice of perpendicular distances of the hyper-plane from the nearest support vector of the two classes





Geometric vs. Functional Margin

- Functional Margin
 - This gives the position of the point with respect to the plane, which does not depend on the magnitude.
- Geometric Margin
 - This gives the distance between the given training example and the given plane.



Finding Margin of a Linear Classifier

 Assume that all data is at least distance 1 from the hyperplane, then the following two constraints follow for a training set {(x⁽ⁱ⁾, y_i)}

```
\mathbf{w}^{T} \mathbf{x}^{(i)} + b \ge 1 \quad \text{if } y_{i} = 1\mathbf{w}^{T} \mathbf{x}^{(i)} + b \le -1 \quad \text{if } y_{i} = -1\Rightarrow \rho = 2|r| = 2 \left| \frac{f(\mathbf{x}^{*})}{\|\mathbf{w}\|} \right| = 2 \left| \frac{\pm 1}{\|\mathbf{w}\|} \right|\rho = \frac{2}{\|\mathbf{w}\|}
```

Support Vector Machines

- Support Vector Machines, in their basic form, are linear classifiers that are trained in a way so as to maximize the margin
- Principles of Operation
 - Define what an optimal hyper-plane is (in way that can be identified in a computationally efficient way)
 - Maximize margin
 - Allows noise tolerance
 - Extend the above definition for non-linearly separable problems
 - have a penalty term for misclassifications
 - Map data to an alternate space where it is easier to classify with linear decision surfaces
 - reformulate problem so that data is mapped implicitly to this space (using kernels)

Margin Maximization in SVM

• We know that if we require

 $\mathbf{w}^{\mathrm{T}}\mathbf{x}^{(\mathbf{i})} + b \ge 1 \quad \text{if } y_i = 1$ $\mathbf{w}^{\mathrm{T}}\mathbf{x}^{(\mathbf{i})} + b \le -1 \quad \text{if } y_i = -1$

• Then the margin is

$$\rho = \frac{2}{\|\mathbf{w}\|}$$

Margin maximization can be performed by reducing the norm of the w vector

 We can present SVM as the following optimization problem



 Combining the objective function and the constraints (represented as losses) as follows

$$\min_{\mathbf{w},b} \quad M = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N L_i <$$

Produces a large +ive value when the ith constraint is violated

The constraint cost function can be written as

$$L_{i} = \max_{\alpha_{i}} \alpha_{i} \left\{ 1 - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \right\} = 0 \quad 1 - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \leq 0$$

$$\alpha_{i} \geq 0$$
When the constraint is being strictly satisfied, α_{i} must be constraint is being violated constraint is being violated

• Combining

Let $P = \min_{\mathbf{w},b} \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^{N} \left[\max_{\alpha_i \ge 0} \alpha_i \left\{ 1 - y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \right\} \right] \right\}$ $= \min_{\mathbf{w},b} \max_{\alpha_i \ge 0} \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^{N} \left[\alpha_i \left\{ 1 - y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \right\} \right] \right\}$

This is the Lagrangian (α_i are the Lagrange Multipliers)

- This is called the primal form of the optimization problem
- The corresponding dual can be written as

Let
$$D = \max_{\alpha_i \ge 0} \min_{\mathbf{w}, b} \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^{N} \left[\alpha_i \left\{ 1 - y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \right\} \right] \right\}$$

- Since the optimization is convex therefore if the Karush-Kuhn-Tucker conditions are satisfied then the primal and dual optimal values will be equal
- In simple words the optimization problem in SVM can be interpreted as

$$\max_{\alpha_i \ge 0} \min_{\mathbf{w}, b} \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \left[\alpha_i \left\{ 1 - y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \right\} \right] \right\}$$

s.t. $\alpha_i \ge 0$

$$\max_{\alpha_{i}\geq 0} \min_{\mathbf{w},b} \left\{ \frac{1}{2} \mathbf{w}^{T} \mathbf{w} + \sum_{i=1}^{N} \left[\alpha_{i} \left\{ 1 - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \right\} \right] \right\} = \max_{\alpha_{i}\geq 0} \left\{ \theta_{D} \left(\mathbf{w}, b \right) \right\}$$

s.t. $\alpha_{i} \geq 0$
 $\theta_{D} \left(\mathbf{w}, b \right) = \min_{\mathbf{w},b} \left\{ \frac{1}{2} \mathbf{w}^{T} \mathbf{w} + \sum_{i=1}^{N} \left[\alpha_{i} \left\{ 1 - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \right\} \right] \right\}$ [1]
Finding the saddle point of $\theta_{D} \left(\mathbf{w}, b \right)$
 $\frac{\partial \theta_{D} \left(\mathbf{w}, b \right)}{\partial \mathbf{w}} = \mathbf{w} + \sum_{i=1}^{N} \left[\alpha_{i} \left\{ -y_{i} \mathbf{x}^{(i)} \right\} \right] = 0 \qquad \Rightarrow \mathbf{w}^{*} = \sum_{i=1}^{N} \alpha_{i} y_{i} \mathbf{x}^{(i)}$
 $\frac{\partial \theta_{D} \left(\mathbf{w}, b \right)}{\partial b} = -\sum_{i=1}^{N} \alpha_{i} y_{i} = 0 \qquad \Rightarrow \sum_{i=1}^{N} \alpha_{i} y_{i} = 0$



• Thus the problem can be written as

$$\max_{\alpha_i \ge 0} \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}^{(i)^T} \mathbf{x}^{(j)} \right\}$$

s.t $\alpha_i \ge 0$
 $\sum_{i=1}^N \alpha_i y_i = 0$

 This quadratic optimization problem can be solved for α_i using standard optimization packages

- The training points with their corresponding α_i greater than zero lie on the boundary and are thus called support vectors as they support the boundary
- Once α_i have been found, the weight can be calculated as

$$\mathbf{w}^* = \sum_{i=1}^{N} \alpha_i^* y_i \mathbf{x}^{(i)} \qquad or \qquad b^* = \frac{1}{n_{SV}} \sum_{\alpha_i > 0}^{N} \left(y_i - \mathbf{w}^{*^T} \mathbf{x}^{(i)} \right)$$

Classification

CIS 62

- The label of an unknown point can be determined by

T

*

$$y = \mathbf{sgn} \left(w^* x + b \right)$$

PIEAS Biomedical Informatics Research Lab 31



Example problem



Matrix formulation of SVM Problem

$$\max_{\alpha_i \ge 0} \left\{ \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \mathbf{x}^{(i)^T} \mathbf{x}^{(j)} \right\}$$

s.t $\alpha_i \ge 0$
 $\sum_{i=1}^{N} \alpha_i y_i = 0$

 $max_{\alpha} \mathbf{1}^{T} \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^{T} \mathbf{X}^{T} \mathbf{X} \boldsymbol{\alpha}$ Subject to: $\boldsymbol{\alpha} \ge \mathbf{0}$ $\mathbf{y}^{T} \boldsymbol{\alpha} = \mathbf{0}$

• If we define the following:

$$-\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 & \alpha_2 & \dots & \alpha_N \end{bmatrix}^T$$

$$-\boldsymbol{y} = \begin{bmatrix} y_1 & y_2 & \dots & y_N \end{bmatrix}^T$$

$$-\boldsymbol{1}_N = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}^T$$

$$-\boldsymbol{X}_{(d \times N)} = \begin{bmatrix} \boldsymbol{x}_1 y_1 & \boldsymbol{x}_2 y_2 & \dots & \boldsymbol{x}_N y_N \end{bmatrix}$$

Solving the SVM using QP

- CVXOPT is a Python package that implements a quadratic programming solver
 - http://abel.ee.ucla.edu/cvxopt/userguide/coneprog.html#quadraticprogramming
- cvxopt.solvers.qp(P, q[, R, s[, U, v[, solver[, initvals]]])
 - Solves the following problem for z:

SV/N/ Droblom

Programming the SVM


Example: Solution of the OR problem



Handling Non-Separable Patterns

 All the derivation till now was based on the requirement that the data be linearly separable

– Practical Problems are Non-Separable

Non-separable data can be handled by relaxing the constraint

$$y_i\left(\mathbf{w}^T\mathbf{x}+b\right) \ge 1$$

• This is achieved as $\mathbf{w}^{T}\mathbf{x}^{(i)} + b \ge 1 - \zeta_{i} \qquad \forall y_{i} = +1$ $\mathbf{w}^{T}\mathbf{x}^{(i)} + b \ge -1 + \zeta_{i} \qquad \forall y_{i} = -1$ $\zeta_{i} \ge 0$

CIS 621: Machine Learning

 $y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \ge 1 - \zeta_i$

Slack Variables

Handling Non-Separable Patterns (Soft Margin)



Handling Non-Separable Patterns (Soft Margin)

- Our objective in designing a SVM here is to maximize the margin while minimizing the misclassification error
- The misclassification error can be written as:

$$\Phi(\zeta) = \sum_{i=1}^{N} I(\zeta_i - 1) \qquad I(\zeta) = \begin{cases} 0 & \zeta \le 0\\ 1 & else \end{cases}$$

 Since this function is non-linear and nonconvex, therefore we choose to use an approximate given by

$$\Phi(\zeta) = \sum_{i=1}^{N} \zeta_i$$

The overall optimization problem can be written as



Weight of the penalty due to margin violation

• Writing the constraints as losses, we have

$$\begin{split} \min_{\mathbf{w},b,\zeta} & M = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^{N} \zeta_i + \sum_{i=1}^{N} L_i + \sum_{i=1}^{N} Z_i \\ L_i = & \max_{\alpha_i} & \alpha_i \left\{ 1 - y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \right\} = \begin{cases} 0 & 1 - \zeta_i - y_i \left(\mathbf{w}^T \mathbf{x}^{(i)} + b \right) \le 0 \\ \infty & else \end{cases} \\ \mathcal{Z}_i = & \max_{\beta_i} & \beta_i \left\{ -\zeta_i \right\} = \begin{cases} 0 & -\zeta_i \le 0 \\ \infty & else \end{cases} , \quad \beta_i \ge 0 \\ \text{Thus} \end{split}$$

$$\begin{array}{ll} \min_{\mathbf{w},b,\zeta \quad \alpha_{i},\beta_{i}} & \frac{1}{2}\mathbf{w}^{T}\mathbf{w} + C\sum_{i=1}^{N}\zeta_{i} + \sum_{i=1}^{N}\alpha_{i}\left\{1 - \zeta_{i} - y_{i}\left(\mathbf{w}^{T}\mathbf{x}^{(i)} + b\right)\right\} - \sum_{i=1}^{N}\beta_{i}\zeta_{i} \\
s.t \quad \alpha_{i} \geq 0, \beta_{i} \geq 0
\end{array}$$

 This is the primal form of the optimization problem $\alpha_i \geq 0$

• The Dual Form is

 $\max_{\alpha_{i},\beta_{i}} \min_{\mathbf{w},b,\zeta} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} + C \sum_{i=1}^{N} \zeta_{i} + \sum_{i=1}^{N} \alpha_{i} \left\{ 1 - \zeta_{i} - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \right\} - \sum_{i=1}^{N} \beta_{i} \zeta_{i}$ s.t $\alpha_{i} \ge 0, \beta_{i} \ge 0$

- Since the optimization problem is convex, therefore the optimal values of the dual and primal are equal
- Therefore the optimization problem can be written solved in its dual form

$$\begin{split} \max_{\alpha_{i},\beta_{i}} & \min_{\mathbf{w},b,\zeta} & \frac{1}{2} \mathbf{w}^{T} \mathbf{w} + C \sum_{i=1}^{N} \zeta_{i} + \sum_{i=1}^{N} \alpha_{i} \left\{ 1 - \zeta_{i} - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \right\} - \sum_{i=1}^{N} \beta_{i} \zeta_{i} \\ &= \max_{\alpha_{i},\beta_{i}} \quad \theta_{D} \left(\boldsymbol{\alpha}, \boldsymbol{\beta} \right) \\ \theta_{D} \left(\boldsymbol{\alpha}, \boldsymbol{\beta} \right) = \min_{\mathbf{w},b,\zeta} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} + C \sum_{i=1}^{N} \zeta_{i} + \sum_{i=1}^{N} \alpha_{i} \left\{ 1 - \zeta_{i} - y_{i} \left(\mathbf{w}^{T} \mathbf{x}^{(i)} + b \right) \right\} - \sum_{i=1}^{N} \beta_{i} \zeta_{i} \\ Solving \\ \frac{\partial \theta_{D}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^{N} \alpha_{i} y_{i} \mathbf{x}^{(i)} = 0 \qquad \Rightarrow \mathbf{w} = \sum_{i=1}^{N} \alpha_{i} y_{i} \mathbf{x}^{(i)} \\ \frac{\partial \theta_{D}}{\partial b} = \sum_{i=1}^{N} \alpha_{i} y_{i} = 0 \qquad \Rightarrow \sum_{i=1}^{N} \alpha_{i} y_{i} = 0 \\ \frac{\partial \theta_{D}}{\partial \zeta_{i}} = C - \alpha_{i} - \beta_{i} = 0 \qquad \Rightarrow \alpha_{i} + \beta_{i} = C \end{split}$$

$$\max_{\alpha_{i},\beta_{i}\geq0}\left\{\sum_{i=1}^{N}\alpha_{i}-\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_{i}\alpha_{j}y_{i}y_{j}\mathbf{x}^{(i)^{T}}\mathbf{x}^{(j)}\right\}$$

s.t $\alpha_{i}\geq0,\beta_{i}\geq0$
 $\alpha_{i}+\beta_{i}=C, \qquad \sum_{i=1}^{N}\alpha_{i}y_{i}=0$
Simplifying further

$$\max_{\alpha_{i},\beta_{i}\geq0}\left\{\sum_{i=1}^{N}\alpha_{i}-\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_{i}\alpha_{j}y_{i}y_{j}\mathbf{x}^{(i)^{T}}\mathbf{x}^{(j)}\right\}$$

s.t $0\leq\alpha_{i}\leq C, \qquad \sum_{i=1}^{N}\alpha_{i}y_{i}=0$

• This problem can be solved using standard optimization packages

- Some Observations
 - $-\alpha_i$ will be non-zero (positive) only for the points that are support vectors
 - $-\beta_i = C \alpha_i$ will be zero for the points that violate the margin condition
 - C is the upper bound on α_i
 - C is the weight of the penalty of the term representing margin violation
 - If C is small, then more margin violations will occur
 - If C is large, lesser margin violations will result

C can be found out through cross-validation



Example...



SVMs uptil now

- Vapnik and Chervonenkis:
 - Hard SVM (1962)
 - Theoretical foundations for SVMs
- Corinna Cortes
 Soft SVM (1995)
- Enter: Bernard Scholkopf (1997)
 - Complete Kernel trick!
 - Kernels not only allow nonlinear boundaries but also allow representation of non-vectoral data





R. A. Fisher 1890-1962

Rosenblatt 1928-1971





V. Vapnik 1936 -

Chervonenkis 1938 - 2014





C. Cortes 1961 -

Scholkopf 1968 -

http://www.svms.org/history.html

Reading

- Sections 10.1-10.3
- Sections 13.1-13.3
- Alpaydin, Ethem. Introduction to Machine Learning. Cambridge, Mass. MIT Press, 2010.

Nonlinear Separation through Transformation

 Given a classification problem with a nonlinear boundary, we can, at times, find a mapping or transformation of the feature space which makes the classification problem linear separable in the transformed space



Examples: Transformation

• Let's define the mapping

$$- \phi\left(\begin{bmatrix}x_1\\x_2\end{bmatrix}\right) = \begin{bmatrix}x_1\\x_2\\x_1x_2\end{bmatrix}$$

• With a mathematical proof show that the above mapping makes the XOR classification problem linearly separable.



XOR: Linear Separability



Examples: Transformation

- Does this mapping do it?
 - $\phi\left(\begin{bmatrix}x_1\\x_2\end{bmatrix}\right) = \begin{bmatrix}x_1\\x_2\\1\end{bmatrix}$
- What about this one?

•
$$\phi\left(\begin{bmatrix}x_1\\x_2\end{bmatrix}\right) = (x_1 + x_2 - 1)^2$$

Transformation Examples

• Can you find a transform that makes the following classification problems linear separable? Can you draw the data points in the new transformed feature space?



Transformations

- Transformations can be used to make the data linearly separable
- But it may not always be possible to find a transformation

– Use a soft-SVM

Another look at the SVM

$$\max_{\alpha_i,\beta_i\geq 0} \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}^{(i)^T} \mathbf{x}^{(j)} \right\}$$

s.t $0 \leq \alpha_i \leq C, \qquad \sum_{i=1}^N \alpha_i y_i = 0$

- We can replace the dot product $x^{(i)^T} x^{(j)}$ with:
 - a generalized dot product (inner product)

•
$$\langle x^{(i)}, x^{(j)} \rangle = \boldsymbol{\phi} (x^{(i)})^T \boldsymbol{\phi} (x^{(j)})$$

- Advantage: can implement feature transformations
- or a function (called the kernel function)
 - $K_{ij} = K(i,j)$
 - Advantage: No need for explicit feature representation

Replacement with a feature transformation

• For the XOR problem we defined the transformation

$$- \phi\left(\begin{bmatrix}x_1\\x_2\end{bmatrix}\right) = \begin{bmatrix}x_1\\x_2\\x_1x_2\end{bmatrix}$$

• We can thus define an inner product

$$- \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle = \boldsymbol{\phi} (\mathbf{x}^{(i)})^{T} \boldsymbol{\phi} (\mathbf{x}^{(j)}) = \begin{bmatrix} x_{1}^{(i)} & x_{2}^{(i)} & x_{2}^{(i)} x_{2}^{(i)} \end{bmatrix} \begin{bmatrix} x_{1}^{(j)} \\ x_{2}^{(j)} \\ x_{2}^{(i)} x_{2}^{(j)} \end{bmatrix}$$
$$= x_{1}^{(i)} x_{1}^{(j)} + x_{2}^{(i)} x_{2}^{(j)} + x_{1}^{(i)} x_{2}^{(i)} x_{1}^{(j)} x_{2}^{(j)}$$

• This inner product implements the transformation and can potentially lead to a non-linear boundary

Kernel Functions \leftrightarrow Feature Transformation

- Since \$\langle x^{(i)}, x^{(j)} \rangle\$ is always a scalar, we can actually use a function (called a kernel function \$k(i, j)\$) to map the two examples to a scalar value
 - For the previous transformation
 - $k(i,j) = x_1^{(i)}x_1^{(j)} + x_2^{(i)}x_2^{(j)} + x_1^{(i)}x_2^{(i)}x_1^{(j)}x_2^{(j)}$
 - Thus, the inner product from a feature transformation can be written as a kernel and a <u>valid</u> kernel function can thus be considered as an inner product in some feature space (proven by Moore-Aronszajn Theorem)
 - We'll talk what makes kernel valid later
 - A kernel is thus a generalized dot product
 - A measure of how similar the two examples are
 - not of whether they belong to the same class

Some kernel functions

- We can use arbitrary kernels, for example ...
 - The dot-product kernel
 - $K(a,b) = a^T b$
 - Feature transform: $\mathbf{\phi}(a) = a$
 - The homogenous polynomial kernels

•
$$K(a, b) = (a^T b)^p$$
, **p** is called degree

- For
$$\boldsymbol{p} = \boldsymbol{2}$$
 and 2D data, $\boldsymbol{a} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$: $\boldsymbol{\phi}(\boldsymbol{a}) = \begin{bmatrix} a_1^2 \\ a_2^2 \\ \sqrt{2}a_1a_2 \end{bmatrix}$

- The Radial Basis Function (RBF) Kernel

•
$$K(a, b) = e^{-\frac{\|a-b\|^2}{2\sigma^2}}$$
, σ controls the spread of the Gaussian
- Feature transform?

• The RBF and Homogenous Polynomial kernels implement non-linear boundaries

Kernels as Similarity Functions

- We have been saying that kernels are similarity functions here is how
 - Example
 - For $K(x, y) = \langle x, y \rangle = x^T y$
 - For the dot product, the associated distance measure is
 - $d(x, y)^{2} = ||x y||^{2} = (x y)^{T}(x y) = x^{T}x 2x^{T}y + y^{T}y = ||x||^{2} + ||y||^{2} 2x^{T}y$

- This implies:
$$K(x, y) = \frac{1}{2} (||x||^2 + ||y||^2 - d(x, y))$$

- Thus, the linear kernel measures the similarity between two points as the inverse of the square of the Euclidean distance between them (up to $||x||^2 + ||y||^2$)
 - » Thus the SVM with a linear kernel is no different, in its distance measurements, from a nearest neighbor classifier!!

Another advantage

$$\max_{\alpha_i,\beta_i\geq 0} \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(i,j) \right\}$$

s.t $0 \leq \alpha_i \leq C, \qquad \sum_{i=1}^N \alpha_i y_i = 0$

- Once we replace the dot product with a kernel function (i.e., perform the kernel trick or 'kernelize' the formulation), the above formulation no longer requires any features!
- As long as you have a kernel function, everything works
 - Remember a kernel function is simply a mapping from two examples to a scalar

But how is that an advantage?

- When the number of dimensions is very large, an implicit representation through a kernel is helpful
- Let's say we have a document classification problem
 - We define a M-dimensional feature vector for each document that indicates if it has any of the pre-specified number of words in it (1) or not (0)
 - M can be very large
 - The dot product of two feature vectors is equal to the number of common words between the two documents
 - Why not simply count the number of words?
 - We can now do that with the kernel trick

Kernel Trick

- You can develop a 'kernelized' version of the soft-SVM as well
- Advantage
 - Removes the explicit representation of data
 - Allows non-linear boundaries

• For understanding a 'valid' kernel, we need to introduce the concept of a kernel matrix

Kernel Matrix

- For $\mathbf{x}^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $\mathbf{x}^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $\mathbf{x}^{(3)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $\mathbf{x}^{(4)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ With the kernel: $k(i,j) = x_1^{(i)}x_1^{(j)} + x_2^{(i)}x_2^{(j)} + x_1^{(i)}x_2^{(i)}x_1^{(j)}x_2^{(j)}$ - k(1,1) = 0 - k(1,2) = 0 $- \dots$
- We get this matrix - $K_{ij} = k(i, j)$

0	0	0	0
0	1	0	1
0	0	1	1
0	1	1	3

 If you know the kernel matrix you don't need to know the original features anymore

Modification Due to Kernel Function

- Change all inner products to kernel functions
- For training,

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

subject to $C \ge \alpha_i \ge 0, \sum_{i=1}^{n} \alpha_i y_i = 0$

Modification Due to Kernel Function

For testing, the new data z is classified as class
 1 if f >0, and as class 2 if f <0

$$\mathbf{w} = \sum_{j=1}^{s} \alpha_{t_j} y_{t_j} \phi(\mathbf{x}_{t_j})$$
$$f = \langle \mathbf{w}, \phi(\mathbf{z}) \rangle + b = \sum_{j=1}^{s} \alpha_{t_j} y_{t_j} K(\mathbf{x}_{t_j}, \mathbf{z}) + b$$

CIS 621: Machine Learning

The discriminant function

 For a test object z, the discriminant function essentially is a weighted sum of the similarity between z and a pre-selected set of objects (the support vectors)

$$f(\mathbf{z}) = \sum_{\mathbf{x}_i \in S} \alpha_i y_i K(\mathbf{z}, \mathbf{x}_i) + b$$

 $\mathcal S$: the set of support vectors

Vectorization

If we define the following:

$$- \alpha = [\alpha_{1} \quad \alpha_{2} \quad \dots \quad \alpha_{N}]^{T}$$

$$- y = [y_{1} \quad y_{2} \quad \dots \quad y_{N}]^{T}$$

$$- \alpha \circ y = [\alpha_{1}y_{1} \quad \alpha_{2}y_{2} \quad \dots \quad \alpha_{N}y_{N}]^{T}$$

$$\cdot \circ \text{ to mean element wise product}$$

$$- \text{ Hadamard product}$$

$$- \text{ Hadamard product}$$

$$- \text{ I}_{N} = [1 \quad 1 \quad \dots \quad 1]^{T}$$

$$- X_{(d \times N)} = [x_{1} \quad x_{2} \quad \dots \quad x_{N}]$$

$$\cdot \text{ This implies: } K = X^{T}X$$

$$max_{\alpha} \mathbf{1}^{T} \alpha - \frac{1}{2} (\alpha \circ y)^{T} K(\alpha \circ y)$$
Subject to:

$$C \ge \alpha \ge \mathbf{0}$$

$$y^{T} \alpha = \mathbf{0}$$
Notice that any data re here the the transported of the term of term of the term of term of

$$max_{\alpha} \mathbf{1}^{T} \boldsymbol{\alpha} - \frac{1}{2} (\boldsymbol{\alpha} \circ \boldsymbol{y})^{T} \boldsymbol{X}^{T} \boldsymbol{X} (\boldsymbol{\alpha} \circ \boldsymbol{y})$$

Subject to:

$$\boldsymbol{\alpha} \geq \boldsymbol{0} \\ \boldsymbol{y}^{T} \boldsymbol{\alpha} = \boldsymbol{0}$$

This replacement of the dot product with the kernel is called the kernel trick

0

Assignment 2B

Implement a soft kernelized SVM



What is the behavior of an inner product?

- Definition of an inner product
 - Symmetry: $\langle x, y \rangle = \langle y, x \rangle$

- Linearity: $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$

- Principle of superposition
- Positive Definiteness
 - $\langle x, x \rangle \geq 0$
 - $\langle x, x \rangle = 0$ iff x = 0
- These conditions need to be satisfied by the kernel function too

Kernel Matrix

 What's the kernel matrix for this problem with the linear kernel?

0	0	0	0
0	1	0	1
0	0	1	1
0	1	1	2

 With the polynomial kernel (p=2)?


Conditions for a function to be a kernel

- A kernel function must satisfy Mercer's condition
 - The kernel matrix must be symmetric positive semi-definite
 - What does that mean?
 - $K(x^{(1)}, x^{(2)}) = K(x^{(2)}, x^{(1)})$

•
$$\boldsymbol{\gamma}^T \boldsymbol{K} \boldsymbol{\gamma} \geq 0$$
 for any $\boldsymbol{\gamma}$

- The Eigen-values of *K* must be non-negative



Example: Kernel

$$K(\boldsymbol{u},\boldsymbol{v})=\boldsymbol{e}^{-\frac{\|\boldsymbol{u}-\boldsymbol{v}\|^2}{2\sigma^2}}$$

• The kernel matrix is (for $\sigma^2 = 0.5$):



• Eigenvalues are 1.93, 0.73, 0.47, 0.86

Kernel Construction

Proposition 3.22 (Closure properties) Let κ_1 and κ_2 be kernels over $X \times X, X \subseteq \mathbb{R}^n, a \in \mathbb{R}^+, f(\cdot)$ a real-valued function on $X, \phi: X \longrightarrow \mathbb{R}^N$ with κ_3 a kernel over $\mathbb{R}^N \times \mathbb{R}^N$, and **B** a symmetric positive semi-definite $n \times n$ matrix. Then the following functions are kernels:

 $\begin{array}{ll} (\mathrm{i}) & \kappa(\mathbf{x},\mathbf{z}) = \kappa_1(\mathbf{x},\mathbf{z}) + \kappa_2(\mathbf{x},\mathbf{z}), \\ (\mathrm{ii}) & \kappa(\mathbf{x},\mathbf{z}) = a\kappa_1(\mathbf{x},\mathbf{z}), \\ (\mathrm{iii}) & \kappa(\mathbf{x},\mathbf{z}) = \kappa_1(\mathbf{x},\mathbf{z})\kappa_2(\mathbf{x},\mathbf{z}), \\ (\mathrm{iv}) & \kappa(\mathbf{x},\mathbf{z}) = f(\mathbf{x})f(\mathbf{z}), \\ (\mathrm{v}) & \kappa(\mathbf{x},\mathbf{z}) = \kappa_3(\phi(\mathbf{x}),\phi(\mathbf{z})), \\ (\mathrm{vi}) & \kappa(\mathbf{x},\mathbf{z}) = \mathbf{x}'\mathbf{B}\mathbf{z}. \end{array}$

Proposition 3.24 Let $\kappa_1(\mathbf{x}, \mathbf{z})$ be a kernel over $X \times X$, where $\mathbf{x}, \mathbf{z} \in X$, and p(x) is a polynomial with positive coefficients. Then the following functions are also kernels:

(i)
$$\kappa(\mathbf{x}, \mathbf{z}) = p(\kappa_1(\mathbf{x}, \mathbf{z})),$$

(ii) $\kappa(\mathbf{x}, \mathbf{z}) = \exp(\kappa_1(\mathbf{x}, \mathbf{z})),$
(iii) $\kappa(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|^2 / (2\sigma^2)).$

Example

- Suppose we have 5 1D data points
 - $x_1=1$, $x_2=2$, $x_3=4$, $x_4=5$, $x_5=6$, with 1, 2, 6 as class 1 and 4, 5 as class 2 \Rightarrow $y_1=1$, $y_2=1$, $y_3=-1$, $y_4=-1$, $y_5=1$
- We use the polynomial kernel of degree 2

 K(x,y) = (xy+1)²
 C is set to 100
- We first find α_i (*i*=1, ..., 5) by

max.
$$\sum_{i=1}^{5} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{5} \sum_{j=1}^{5} \alpha_{i} \alpha_{j} y_{i} y_{j} (x_{i} x_{j} + 1)^{2}$$

subject to $100 \ge \alpha_{i} \ge 0, \sum_{i=1}^{5} \alpha_{i} y_{i} = 0$

Example

• By using a QP solver, we get

$$-\alpha_1=0, \alpha_2=2.5, \alpha_3=0, \alpha_4=7.333, \alpha_5=4.833$$

- Note that the constraints are indeed satisfied
- The support vectors are $\{x_2=2, x_4=5, x_5=6\}$
- The bias turns out to be b = 9
- The discriminant function is

$$f(z) = 0.6667z^2 - 5.333z + 9$$

Example



CIS 621: Machine Learning

Solving the XOR

- C=1
- With polynomial kernel of degree 2



CIS 621: Machine Learning

Solving the XOR

- C=1
- With RBF Kernel (sigma = 0.5)



CIS 621: Machine Learning

Solving the XOR

- C=1
- With RBF Kernel (sigma = 0.3)



CIS 621: Machine Learning

3D Plot



Using the SVM

- Read:
- Ben-Hur, Asa, and Jason Weston. 2010. "A User's Guide to Support Vector Machines." In *Data Mining Techniques for the Life Sciences*, edited by Oliviero Carugo and Frank Eisenhaber, 223–39. Methods in Molecular Biology 609. Humana Press. <u>http://dx.doi.org/10.1007/978-1-60327-241-4_13</u>
- <u>http://pyml.sourceforge.net/doc/howto.pdf</u>

Steps for Feature based Classification

- Prepare the pattern matrix
- Select the kernel function to use
- Select the parameter of the kernel function and the value of *C*
 - You can use the values suggested by the SVM software, or you can set apart a validation set to determine the values of the parameter
- Execute the training algorithm and obtain the $\alpha_{\rm i}$
- Unseen data can be classified using the α_{i} and the support vectors

Choosing the Kernel Function

- Probably the most tricky part of using SVM.
- The kernel function is important because it creates the kernel matrix, which summarizes all the data
- In practice, a low degree polynomial kernel or RBF kernel with a reasonable width is a good initial try



Choosing C

- Cross-validation
 - To assess how good a classifier is we can use cross-validation
 - Divide the data randomly into k parts
 - If the data is imbalanced, use stratified sampling
 - Use k-1 parts for training
 - And the held-out part for testing to evaluate accuracy or ROC curve or other performance metrics
- To choose C, you can do nested crossvalidation

Handling data imbalance

- If the data is imbalanced (too much of one class and only a small number of examples from the other)
 - You can set an individual C for each example
 - Can also be used to reflect a priori knowledge



Strengths and Weaknesses of SVM

- Strengths
 - Margin maximization and kernelized
 - Training is relatively easy
 - No local optimal, unlike in neural networks
 - It scales relatively well to high dimensional data
 - Tradeoff between classifier complexity and error can be controlled explicitly (through C)
 - Non-traditional data like strings and trees can be used as input to SVM, instead of feature vectors
- Weaknesses
 - Need to choose a "good" kernel function.

Multi-class Classification

- SVM is basically a two-class classifier
- One can change the QP formulation to allow multi-class classification and such SVMs do exist
- But you can also try to do multi-class classification

End of Lecture

We want to make a machine that will be proud of us.

- Danny Hillis

CIS 621: Machine Learning