

Kernel Methods and Support Vector Machines

Oliver Schulte - CMPT 726

Bishop PRML Ch. 6

Support Vector Machines

Defining Characteristics

- Like logistic regression, good for continuous input features, discrete target variable.
- Like nearest neighbor, a *kernel method*: classification is based on weighted similar instances. The kernel defines similarity measure.
- Sparsity: Tries to find a few important instances, the *support vectors*.
- Intuition: Netflix recommendation system.

SVMs: Pros and Cons

Pros

- Very good classification performance, basically unbeatable.
- Fast and scaleable learning.
- Pretty fast inference.

Cons

- No model is built, therefore black-box.
- Still need to specify kernel function (like specifying basis functions).
- Issues with multiple classes, can use probabilistic version. (Relevance Vector Machine).

Two Views of SVMs

Theoretical View: linear separator

- SVM looks for linear separator but in *new feature space*.
- Uses a new criterion to choose a line separating classes: *max-margin*.

User View: kernel-based classification

- User specifies a kernel function.
- SVM learns weights for instances.
- Classification is performed by taking average of the labels of other instances, weighted by a) similarity b) instance weight.

Nice demo on web

<http://www.youtube.com/watch?v=3liCbRZPrZA>.

Two Views of SVMs

Theoretical View: linear separator

- SVM looks for linear separator but in *new feature space*.
- Uses a new criterion to choose a line separating classes: *max-margin*.

User View: kernel-based classification

- User specifies a kernel function.
- SVM learns weights for instances.
- Classification is performed by taking average of the labels of other instances, weighted by a) similarity b) instance weight.

Nice demo on web

<http://www.youtube.com/watch?v=3liCbRZPrZA>.

Two Views of SVMs

Theoretical View: linear separator

- SVM looks for linear separator but in *new feature space*.
- Uses a new criterion to choose a line separating classes: *max-margin*.

User View: kernel-based classification

- User specifies a kernel function.
- SVM learns weights for instances.
- Classification is performed by taking average of the labels of other instances, weighted by a) similarity b) instance weight.

Nice demo on web

<http://www.youtube.com/watch?v=3liCbRZPrZA>.

Two Views of SVMs

Theoretical View: linear separator

- SVM looks for linear separator but in *new feature space*.
- Uses a new criterion to choose a line separating classes: *max-margin*.

User View: kernel-based classification

- User specifies a kernel function.
- SVM learns weights for instances.
- Classification is performed by taking average of the labels of other instances, weighted by a) similarity b) instance weight.

Nice demo on web

<http://www.youtube.com/watch?v=3liCbRZPrZA>.

Outline

Theoretical View

User View

Building Kernels

Computing the Max-Margin Classifier

Non-Separable Data

Linear Classification Revisited

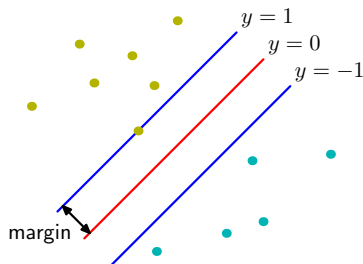
- Consider a two class classification problem
- Use a linear model

$$y(\mathbf{x}) = \mathbf{w} \bullet \mathbf{x} + b$$

followed by a threshold function

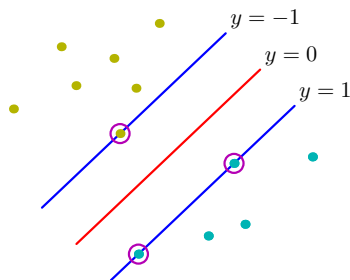
- For now, let's assume training data are linearly separable (possibly after mapping to higher-dimensional space).
 - Recall that the perceptron would converge to a perfect classifier for such data
 - But there are many such perfect classifiers

Max Margin Classifiers



- We can define the **margin** of a classifier as the minimum distance to any example
- In **support vector machines** the decision boundary which maximizes the margin is chosen.
- Intuitively, this is the line “right in the middle” between the two classes.

Support Vectors



- The **support vectors** are the points at minimum distance to the decision boundary.
- The max-margin boundary depends only on the support vectors: other data points do not matter for classification and need not be stored.

Outline

Theoretical View

User View

Building Kernels

Computing the Max-Margin Classifier

Non-Separable Data

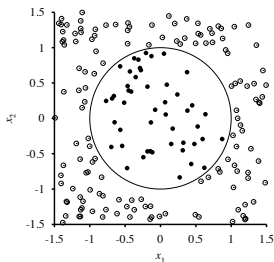
Example: X-OR

- X-OR problem: class of (x_1, x_2) is positive iff $x_1 \cdot x_2 > 0$.
- Not linearly separable in input space, but linearly separable if *we add extra dimensions*.
- Use 6 basis functions
$$\phi(x_1, x_2) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2).$$
- Simple classifier $y(x_1, x_2) = \phi_5(x_1, x_2) = \sqrt{2}x_1x_2$.
 - *Linear* in basis function space.

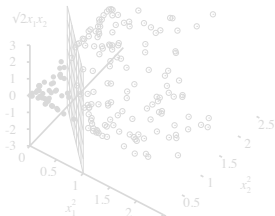
Example: X-OR

- X-OR problem: class of (x_1, x_2) is positive iff $x_1 \cdot x_2 > 0$.
- Not linearly separable in input space, but linearly separable if *we add extra dimensions*.
- Use 6 basis functions
$$\phi(x_1, x_2) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2).$$
- Simple classifier $y(x_1, x_2) = \phi_5(x_1, x_2) = \sqrt{2}x_1x_2$.
 - *Linear* in basis function space.

Linear Separability Example

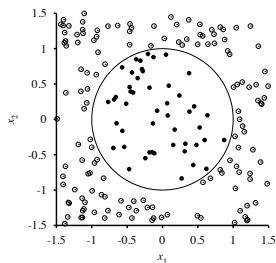


Decision Boundary $x_1^2 + x_2^2 \leq 1$

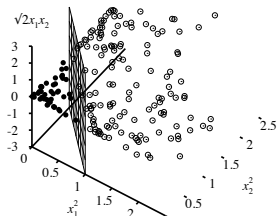


3-D mapping $(x_1^2, x_2^2, \sqrt{2}x_1x_2)$

Linear Separability Example



Decision Boundary $x_1^2 + x_2^2 \leq 1$



3-D mapping $(x_1^2, x_2^2, \sqrt{2}x_1x_2)$

The Kernel Trick

- There can be many extra dimensions, even infinite (see assignment).
 - Don't want to compute basis function mapping $\phi(\mathbf{x})$.
- Key insight 1: Linear classification requires *only the dot product*.
- Key insight 2: The high-dimensional dot product $\phi(\mathbf{x}) \bullet \phi(\mathbf{z})$ can often be computed as a **kernel** function of the input vectors only:

$$\phi(\mathbf{x}) \bullet \phi(\mathbf{z}) = k(\mathbf{x}, \mathbf{z}).$$

Kernel Trick Example

- Consider again the X-OR 6 basis functions
 $\phi(x_1, x_2) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$.
- Exercise: find a closed form expression for $\phi(\mathbf{x}) \bullet \phi(\mathbf{z})$
- Solution: Dot product $\phi(\mathbf{x}) \bullet \phi(\mathbf{z}) = (1 + \mathbf{x} \bullet \mathbf{z})^2 = k(\mathbf{x}, \mathbf{z})$.
- A quadratic kernel.

Kernel Trick Example

- Consider again the X-OR 6 basis functions
 $\phi(x_1, x_2) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$.
- Exercise: find a closed form expression for $\phi(\mathbf{x}) \bullet \phi(\mathbf{z})$
- Solution: Dot product $\phi(\mathbf{x}) \bullet \phi(\mathbf{z}) = (1 + \mathbf{x} \bullet \mathbf{z})^2 = k(\mathbf{x}, \mathbf{z})$.
- A **quadratic kernel**.

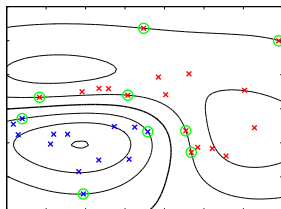
The Kernel Classification Formula

- Suppose we have a kernel function k and N labelled instances with weights $a_n \geq 0, n = 1, \dots, N$.
- As with the perceptron, the target labels $+1$ are for positive class, -1 for negative class.
- Then

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

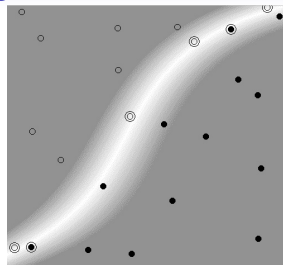
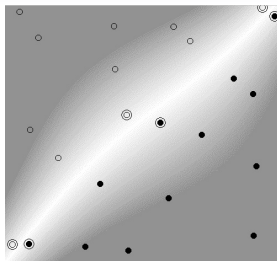
- \mathbf{x} is classified as positive if $y(\mathbf{x}) > 0$, negative otherwise.
- If $a_n > 0$, then \mathbf{x}_n is a **support vector**.
- Don't need to store other vectors.
- \mathbf{a} will be sparse - many zeros.

Example 1



- SVM with Gaussian kernel
- Support vectors circled.
- They are the closest to the other class.
- Note non-linear decision boundary in x space

Example2



- From Burges, *A Tutorial on Support Vector Machines for Pattern Recognition* (1998)
- SVM trained using cubic polynomial kernel

$$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \bullet \mathbf{x}_2 + 1)^3$$
- Left is linearly separable
 - Note decision boundary is almost linear, even using cubic polynomial kernel
- Right is not linearly separable
 - But is separable using polynomial kernel

Learning the Instance Weights

- The max-margin classifier is found by solving the following problem:
- Maximize wrt \mathbf{a}

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

subject to the constraints

- $a_n \geq 0, n = 1, \dots, N$
- $\sum_{n=1}^N a_n t_n = 0$
- It is quadratic, with linear constraints, convex in \mathbf{a}
- Optimal \mathbf{a} can be found
 - With large datasets, local search strategies employed

let's check SVM demo

<http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml>

Learning the Instance Weights

- The max-margin classifier is found by solving the following problem:
- Maximize wrt \mathbf{a}

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

subject to the constraints

- $a_n \geq 0, n = 1, \dots, N$
- $\sum_{n=1}^N a_n t_n = 0$
- It is quadratic, with linear constraints, convex in \mathbf{a}
- Optimal \mathbf{a} can be found
 - With large datasets, local search strategies employed

let's check SVM demo

<http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml>

Outline

Theoretical View

User View

Building Kernels

Computing the Max-Margin Classifier

Non-Separable Data

Valid Kernels

- Valid kernels: if $k(\cdot, \cdot)$ satisfies:
 - Symmetric; $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$
 - Positive definite; for any $\mathbf{x}_1, \dots, \mathbf{x}_N$, the Gram matrix \mathbf{K} must be positive semi-definite:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

- Positive semi-definite means $\mathbf{x} \bullet \mathbf{K} \mathbf{x} \geq 0$ for all \mathbf{x} (like metric)
- then $k(\cdot, \cdot)$ corresponds to a dot product in some space ϕ
- a.k.a. Mercer kernel, admissible kernel, reproducing kernel
 - Theorem of Mercer's 1909!

Examples of Kernels

- Some kernels:
 - Linear kernel $k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \bullet \mathbf{x}_2$
 - Polynomial kernel $k(\mathbf{x}_1, \mathbf{x}_2) = (1 + \mathbf{x}_1 \bullet \mathbf{x}_2)^d$
 - Contains all polynomial terms up to degree d
 - Gaussian kernel $k(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\|\mathbf{x}_1 - \mathbf{x}_2\|^2/2\sigma^2)$
 - Infinite dimension feature space

Constructing Kernels

- Can build new valid kernels from existing valid ones:
 - $k(\mathbf{x}_1, \mathbf{x}_2) = ck_1(\mathbf{x}_1, \mathbf{x}_2)$, $c > 0$
 - $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) + k_2(\mathbf{x}_1, \mathbf{x}_2)$
 - $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2)k_2(\mathbf{x}_1, \mathbf{x}_2)$
 - $k(\mathbf{x}_1, \mathbf{x}_2) = \exp(k_1(\mathbf{x}_1, \mathbf{x}_2))$
- Table on p. 296 gives many such rules

More Kernels

- **Stationary kernels** are only a function of the difference between arguments: $k(\mathbf{x}_1, \mathbf{x}_2) = k(\mathbf{x}_1 - \mathbf{x}_2)$
 - Translation invariant in input space:
 $k(\mathbf{x}_1, \mathbf{x}_2) = k(\mathbf{x}_1 + \mathbf{c}, \mathbf{x}_2 + \mathbf{c})$
- **Homogeneous kernels**, a. k. a. **radial basis functions** only a function of magnitude of difference: $k(\mathbf{x}_1, \mathbf{x}_2) = k(\|\mathbf{x}_1 - \mathbf{x}_2\|)$
- Set subsets $k(A_1, A_2) = 2^{|A_1 \cap A_2|}$, where $|A|$ denotes number of elements in A
- Domain-specific: think hard about your problem, figure out what it means to be similar, define as $k(\cdot, \cdot)$, prove positive definite.

Outline

Theoretical View

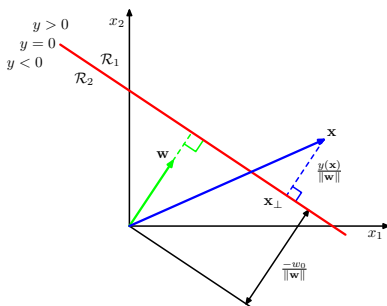
User View

Building Kernels

Computing the Max-Margin Classifier

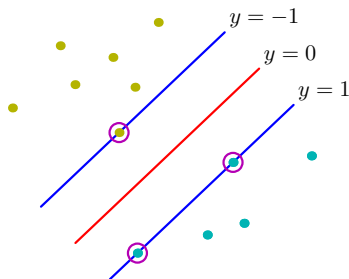
Non-Separable Data

Marginal Geometry



- See assignment.
- Projection of x in w dir. is $\frac{w \bullet x}{\|w\|}$
- $y(x) = 0$ when $w \bullet x = -b$, or $\frac{w \bullet x}{\|w\|} = \frac{-b}{\|w\|}$
- So $\frac{w \bullet x}{\|w\|} - \frac{-b}{\|w\|} = \frac{y(x)}{\|w\|}$ is signed distance to decision boundary

Support Vectors



- Assuming data are separated by the hyperplane, distance to decision boundary is $\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$
- The maximum margin criterion chooses \mathbf{w}, b by:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w} \bullet \mathbf{x}_n + b)] \right\}$$

- Points with this min value are known as **support vectors**

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w} \bullet \mathbf{x}_n) + b] \right\}$$

- Exercise: Prove that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(x_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can set (how?):

$$t_*(\mathbf{w} \bullet \mathbf{x}_* + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w} \bullet \mathbf{x}_n + b) \geq 1$$

- Knowing that *min* distance = 1, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w} \bullet \mathbf{x}_n) + b] \right\}$$

- Exercise: Prove that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(x_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can set (how?):

$$t_*(\mathbf{w} \bullet \mathbf{x}_* + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w} \bullet \mathbf{x}_n + b) \geq 1$$

- Knowing that *min* distance = 1, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w} \bullet \mathbf{x}_n) + b] \right\}$$

- Exercise: Prove that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(x_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can set (how?):

$$t_*(\mathbf{w} \bullet \mathbf{x}_* + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w} \bullet \mathbf{x}_n + b) \geq 1$$

- Knowing that *min* distance = 1, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w} \bullet \mathbf{x}_n) + b] \right\}$$

- Exercise: Prove that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(x_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can set (how?):

$$t_*(\mathbf{w} \bullet \mathbf{x}_* + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w} \bullet \mathbf{x}_n + b) \geq 1$$

- Knowing that *min* distance = 1, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

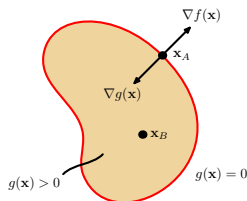
Canonical Representation

- So the optimization problem is now a constrained optimization problem:

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \forall n, t_n(\mathbf{w} \bullet \mathbf{x}_n + b) \geq 1 \end{aligned}$$

- To solve this, we need to use [Lagrange multipliers](#)

Lagrange Multipliers - Inequality Constraints



Consider the problem:

$$\min_{\mathbf{x}} f(\mathbf{x})$$

$$s.t. \quad g_1(\mathbf{x}) \geq 0, \dots, g_N(\mathbf{x}) \geq 0$$

Solutions are stationary points of the **Lagrangian**

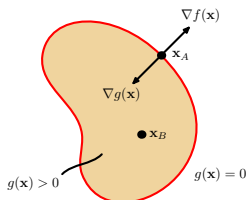
$$L(\mathbf{x}, \mathbf{a}) = f(\mathbf{x}) - \sum_{n=1}^N a_n g_n(\mathbf{x})$$

where for each n we have the KKT conditions,

- $a_n \geq 0$
- $g_n(\mathbf{x}) \geq 0$
- $a_n g_n(\mathbf{x}) = 0$

Therefore $a_n = 0$ (inactive constraint) or $g_n(\mathbf{x}) = 0$ (active constraint).

Lagrange Multipliers - Inequality Constraints



Consider the problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_1(\mathbf{x}) \geq 0, \dots, g_N(\mathbf{x}) \geq 0 \end{aligned}$$

Solutions are stationary points of the **Lagrangian**

$$L(\mathbf{x}, \mathbf{a}) = f(\mathbf{x}) - \sum_{n=1}^N a_n g_n(\mathbf{x})$$

where for each n we have the KKT conditions,

- $a_n \geq 0$
- $g_n(\mathbf{x}) \geq 0$
- $a_n g_n(\mathbf{x}) = 0$

Therefore $a_n = 0$ (inactive constraint) or $g_n(\mathbf{x}) = 0$ (active constraint).

Now Where Were We

- So the optimization problem is now a constrained optimization problem:

$$\arg \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2}$$

$$s.t. \quad \forall n, t_n(\mathbf{w} \bullet \mathbf{x}_n + b) \geq 1$$

- For this problem, the Lagrangian (with N multipliers a_n) is:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{n=1}^N a_n \{t_n(\mathbf{w} \bullet \mathbf{x}_n + b) - 1\}$$

- We can find the derivatives of L wrt \mathbf{w} , b and set to 0:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

$$0 = \sum_{n=1}^N a_n t_n$$

The Dual Formulation

- Recall the condition:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

- Exercise: Show that $\frac{\|\mathbf{w}\|^2}{2} = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m (\mathbf{x}_n \bullet \mathbf{x}_m)$.
- Exercise: Show that

$$\sum_{n=1}^N a_n \{t_n (\mathbf{w} \bullet \mathbf{x}_n + b) - 1\} =$$

$$\sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m (\mathbf{x}_n \bullet \mathbf{x}_m) + b \sum_{n=1}^N a_n t_n - \sum_{n=1}^N a_n$$

Solving The Dual Formulation

- Combining the exercise results, we obtain the **dual Lagrangian** as a function of the Lagrange multipliers only:

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

- The stationary points of \tilde{L} provide lower bounds on the original problem, so we want to maximize \tilde{L} (see http://en.wikipedia.org/wiki/Lagrange_duality).
- Apply the kernel trick to replace \bullet with kernel k , and remember the constraints on the Lagrange multipliers, to arrive at the following problem:
- Maximize $\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$ subject to the constraints that
 - $a_n \geq 0, n = 1, \dots, N$
 - $\sum_{n=1}^N a_n t_n = 0$

From The Dual Solution \mathbf{a} to a Classifier

- Given the solution \mathbf{a} , we have formulas for \mathbf{w} and for b (omitted). Then classify as follows:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

$$y(\mathbf{x}) = \mathbf{w} \bullet \phi(\mathbf{x}) + b = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

- Recall that every constraint is either inactive ($a_n = 0$) or active ($a_n > 0$ and $t_n y(\mathbf{x}_n) = 1$).
- If $a_n > 0$, then \mathbf{x}_n is a **support vector**.
- \mathbf{a} will be sparse - many zeros.
 - Don't need to store \mathbf{x}_n for which $a_n = 0$

Outline

Theoretical View

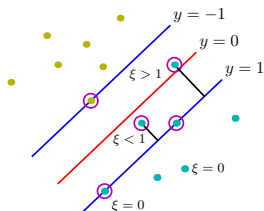
User View

Building Kernels

Computing the Max-Margin Classifier

Non-Separable Data

Non-Separable Data

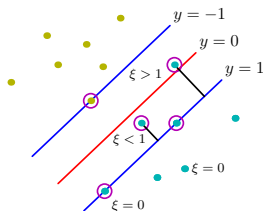


- For most problems, data will not be linearly separable (even in feature space ϕ)
- Can relax the constraints from

$$t_n y(\mathbf{x}_n) \geq 1 \quad \text{to} \quad t_n y(\mathbf{x}_n) \geq 1 - \xi_n$$

- The $\xi_n \geq 0$ are called **slack variables**
 - $\xi_n = 0$, satisfy original problem, so x_n is on margin or correct side of margin
 - $0 < \xi_n < 1$, inside margin, but still correctly classified
 - $\xi_n > 1$, mis-classified

Loss Function For Non-separable Data



- Non-zero slack variables are bad, penalize while maximizing the margin:

$$\min C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

- Constant $C > 0$ controls importance of large margin versus incorrect (non-zero slack)
 - Set using cross-validation
- Optimization is same quadratic, different constraints, convex

SVM Loss Function

- The SVM for the separable case solved the problem:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

s.t. $\forall n, t_n y_n \geq 1$

- Can write this as:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{\infty}(t_n y_n) + \lambda \|\mathbf{w}\|^2$$

where $E_{\infty}(z) = 0$ if $z \geq 1$, ∞ otherwise

- Non-separable case relaxes this to be:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{SV}(t_n y_n) + \lambda \|\mathbf{w}\|^2$$

where $E_{SV}(z) = [z]_+$ **hinge loss**

- $[z]_+ = z$ if $z \leq 1$, 0 otherwise

SVM Loss Function

- The SVM for the separable case solved the problem:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

s.t. $\forall n, t_n y_n \geq 1$

- Can write this as:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{\infty}(t_n y_n) + \lambda \|\mathbf{w}\|^2$$

where $E_{\infty}(z) = 0$ if $z \geq 1$, ∞ otherwise

- Non-separable case relaxes this to be:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{SV}(t_n y_n) + \lambda \|\mathbf{w}\|^2$$

where $E_{SV}(z) = [z]_+$ hinge loss

- $[z]_+ = z$ if $z \leq 1$, 0 otherwise

SVM Loss Function

- The SVM for the separable case solved the problem:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. \quad \forall n, t_n y_n \geq 1$$

- Can write this as:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{\infty}(t_n y_n) + \lambda \|\mathbf{w}\|^2$$

where $E_{\infty}(z) = 0$ if $z \geq 1$, ∞ otherwise

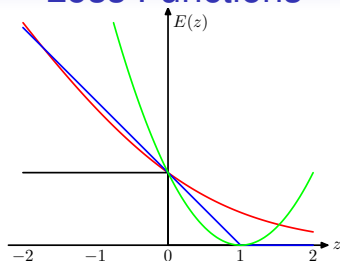
- Non-separable case relaxes this to be:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{SV}(t_n y_n) + \lambda \|\mathbf{w}\|^2$$

where $E_{SV}(z) = [z]_+$ **hinge loss**

- $[z]_+ = z$ if $z \leq 1$, 0 otherwise

Loss Functions



- Linear classifiers, compare **loss function** used for learning
- $z = y_n t_n \leq 0$ iff there is an error (with $t_n \in \{+1, -1\}$).
- $z = y_n t_n \geq 1$ iff the point is on the right side of the margin boundary.
 - Black is misclassification error
 - Transformed simple linear classifier, **squared error**:
 $(y_n - t_n)^2$
 - Transformed logistic regression, **cross-entropy error**: $t_n \ln y_n$
 - SVM, **hinge loss**:
 - positive *only* if the point is on the wrong side of the margin boundary. Sparse solutions.

Two Views of Learning as Optimization

- The original SVM goal was of the form:
 - Find the simplest hypothesis that is consistent with the data, or
 - Maximize simplicity, given a consistency constraint.
- This general idea appears in much scientific model building, in image processing, and other applications.
- Bayesian methods use a criterion of the form
 - Find a trade-off between simplicity and data fit, or
 - Maximize sum of the type (data fit - λ simplicity)
 - e.g., $\ln(P(D|M)) - \lambda \ln(P(M))$ where the model prior M is higher for simpler models.

Pros and Cons of Learning Criteria

- The Bayesian approach has a solid probabilistic foundation in Bayes' theorem.
- Seems to be especially suitable for noisy data.
- The constraint-based approach is often easy for users to understand.
- Often leads to sparser simpler models.
- Suitable for “clean” data.

Conclusion

- Readings: Ch. 6.1-6.2 (pp. 291-297)
- Non-linear features, or domain-specific similarity measurements are useful
- Dot products of non-linear features, or similarity measurements, can be written as kernel functions
 - Validity by positive semi-definiteness of kernel function
- Can have algorithm work in non-linear feature space without actually mapping inputs to feature space
 - Advantageous when feature space is high-dimensional