CSCE 990 Lecture 7:	Introduction
SVMs for Classification*	Introduction
	• Finally, we get to put everything together!
Stephen D. Scott	 Much of this lecture is material we've covered previously, but now we'll make it specific to SVMs
February 14, 2006	• We'll also formalize the notion of the margin, introduce soft margin, and argue why we want to minimize $\ \mathbf{w}\ ^2$
*Most figures ©2002 MIT Press, Bernhard Schölkopf, and Alex Smola. 1	2
Outline	Canonical Hyperplanes
Outline • Canonical hyperplanes	 Canonical Hyperplanes Any hyperplane in a dot product space H can be written as
Outline Canonical hyperplanes The (geometrical) margin and the margin error bound 	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $H = \{\mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$
Outline • Canonical hyperplanes • The (geometrical) margin and the margin error bound • Optimal margin hyperplanes	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $H = \{\mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$ • $\langle \mathbf{w}, \mathbf{x} \rangle$ is the length of \mathbf{x} in the direction of \mathbf{w} , multiplied by $\ \mathbf{w}\ $, i.e. each $\mathbf{x} \in H$ has the same length in the direction of \mathbf{w}
Outline • Canonical hyperplanes • The (geometrical) margin and the margin error bound • Optimal margin hyperplanes • Adding kernels	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $H = \{\mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$ • $\langle \mathbf{w}, \mathbf{x} \rangle$ is the length of \mathbf{x} in the direction of \mathbf{w} , multiplied by $ \mathbf{w} $, i.e. each $\mathbf{x} \in H$ has the same length in the direction of \mathbf{w}
Outline • Canonical hyperplanes • The (geometrical) margin and the margin error bound • Optimal margin hyperplanes • Adding kernels • Soft margin hyperplanes	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $H = \{\mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$ • $\langle \mathbf{w}, \mathbf{x} \rangle$ is the length of \mathbf{x} in the direction of \mathbf{w} , multiplied by $\ \mathbf{w}\ $, i.e. each $\mathbf{x} \in H$ has the same length in the direction of \mathbf{w}
Outline • Canonical hyperplanes • The (geometrical) margin and the margin error bound • Optimal margin hyperplanes • Adding kernels • Soft margin hyperplanes	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $H = \{\mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$ • $\langle \mathbf{w}, \mathbf{x} \rangle$ is the length of \mathbf{x} in the direction of \mathbf{w} , multiplied by $ \mathbf{w} $, i.e. each $\mathbf{x} \in H$ has the same length in the direction of \mathbf{w} • $\mathbf{w}, \mathbf{x} + b < 0$
Outine • Canonical hyperplanes • The (geometrical) margin and the margin error bound • Optimal margin hyperplanes • Adding kernels • Soft margin hyperplanes • Multi-class classification • Application: handwritten digit recognition	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $H = \{x \in \mathcal{H} \mid \langle w, x \rangle + b = 0\}, w \in \mathcal{H}, b \in \mathbb{R}$ • $\langle w, x \rangle$ is the length of x in the direction of w, multiplied by $ w $, i.e. each $x \in H$ has the same length in the direction of w • $(w, x) + b < 0$
Outline • Canonical hyperplanes • The (geometrical) margin and the margin error bound • Optimal margin hyperplanes • Adding kernels • Soft margin hyperplanes • Multi-class classification • Application: handwritten digit recognition	Canonical Hyperplanes • Any hyperplane in a dot product space \mathcal{H} can be written as $\mathcal{H} = \{ \mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0 \}, \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$ • $\langle \mathbf{w}, \mathbf{x} \rangle$ is the length of \mathbf{x} in the direction of \mathbf{w} , multiplied by $\ \mathbf{w}\ $, i.e. each $\mathbf{x} \in H$ has the same length in the direction of \mathbf{w} • $(\mathbf{w}, \mathbf{x} + b < 0)$

Canonical Hyperplanes

(cont'd)

- Note that if both w and b are multiplied by the same non-zero constant, H is unchanged
- **D7.1** The pair $(\mathbf{w}, b) \in \mathcal{H}$ is called a <u>canonical</u> form of the hyperplane H wrt a set of patterns $\mathbf{x}_1, \ldots, \mathbf{x}_m \in \mathcal{H}$ if it is scaled such that

$$\min_{i=1,\dots,m} |\langle \mathbf{w}, \mathbf{x}_i \rangle + b| = 1$$



• Given a canonical hyperplane (\mathbf{w}, b) , the corresponding <u>decision function</u> is $f_{\mathbf{w},b}(\mathbf{x}) := \operatorname{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$

5

7

Justifications for Large Margin

- Why do we want large margin hyperplanes (that separate the training data)?
- Insensitivity to pattern noise
 - E.g. if each (noisy) test point $(\mathbf{x} + \Delta \mathbf{x}, y)$ is near some (noisy) training point (\mathbf{x}, y) with $\|\Delta \mathbf{x}\| < r$, then if $\rho > r$ we correctly classify all test points



The Margin

D7.2 For a hyperplane $\{\mathbf{x} \in \mathcal{H} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}$, define

$$\rho_{\mathbf{w},b}(\mathbf{x},y) := y(\langle \mathbf{w}, \mathbf{x} \rangle + b) / \|\mathbf{w}\|$$

as the geometrical margin (or simply margin) of the point $(x, y) \in \mathcal{H} \times \{-1, +1\}$. Further,

$$\rho_{\mathbf{w},b} := \min_{i=1,\dots,m} \rho_{\mathbf{w},b}(\mathbf{x}_i, y_i)$$

is the (geometrical) margin of $(x_1, y_1), \ldots, (x_m, y_m)$ (typically the training set)

- In D7.2, we are really using the hyperplane $(\hat{\mathbf{w}}, \hat{b}) := (\mathbf{w}/||\mathbf{w}||, b/||\mathbf{w}||)$, which has unit length
- Further, $\langle \hat{\mathbf{w}}, \mathbf{x} \rangle + \hat{b}$ is x's distance to this hyperplane, and multiplying by y implies that the margin is positive if (\mathbf{x}, y) is correctly classified
- Since canonical hyperplanes have minimum distance 1 to data points, the margin of a canonical hyperplane is $\rho_{{\bf w},b}=1/\|{\bf w}\|$
- \bullet I.e. decreasing $\|\mathbf{w}\|$ increases the margin!

6

Justifications for Large Margin (cont'd)

- Insensitivity to parameter noise
 - If all patterns are at least ρ from $H = (\mathbf{w}, b)$ and all patterns are bounded in length by R, then small changes in the parameters of H will not change classification
 - I.e. can encode H with fewer bits than if we precisely encoded it and still be correct on training set
 - \Rightarrow minimum description length/compression of data



Justifications for Large Margin (cont'd)

T7.3 For decision functions $f(\mathbf{x}) = \operatorname{sgn}\langle \mathbf{w}, \mathbf{x} \rangle$, let $\|\mathbf{w}\| \leq \Lambda$, $\|\mathbf{x}\| \leq R$, $\rho > 0$, and ν be the margin error, i.e. the fraction of training examples with margin $< \rho/\|\mathbf{w}\|$. Then if all training and test patterns are drawn iid, with probability at least $1 - \delta$ the test error is upper bounded by

$$\nu + \sqrt{\frac{c}{m} \left(\frac{R^2 \Lambda^2}{\rho^2} \ln^2 m + \ln(1/\delta)\right)}$$

where \boldsymbol{c} is a constant and \boldsymbol{m} is the training set size

• Related to VC dimension of large-margin classifiers, but not exactly what we covered in Chapter 5; e.g. R_{emp} , which was a prediction error rate, is replaced with ν , which is a margin error rate

Optimal Margin Hyperplanes

- Want hyperplane that correctly classifies all training patters with maximum margin
- When using canonical hyperplanes, implies that we want $y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \ge 1$ for all i = 1, ..., m
- We know that we want to minimize the weight vector's length to maximize the margin, so this yields the following constrained quadratic optimization problem:

 $\begin{array}{ll} \underset{\mathbf{w}\in\mathcal{H},b\in\mathbb{R}}{\text{minimize}} & \tau(\mathbf{w}) = \|\mathbf{w}\|^2/2 \\ \text{s.t.} & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \ge 1, i = 1, \dots, m \end{array}$ (1)

- Another optimization problem. Hey! I have a great idea! Let's derive the dual!
- Langrangian:

$$L(\mathbf{w}, b, \alpha) = \|\mathbf{w}\|^2 / 2 - \sum_{i=1}^m \alpha_i (y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1)$$

with $\alpha_i \ge 0$

Justifications for Large Margin Margin Error Bound (cont'd)

- Increasing ρ decreases the square root term, but can increase ν
 - Thus we want to maximize ρ while simultaneously minimizing ν
 - Can instead fix $\rho = 1$ (canonical hyperplanes) and minimize ||w|| while minimizing margin errors
 - In our first quadratic program, we'll set constraints to make $\nu = 0$

10

The Dual Optimization Problem (cont'd)

• Recall that at the saddle point, the partial derivatives of *L* wrt the primal variables must each go to 0:

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = -\sum_{i=1}^{m} \alpha_i y_i = 0$$
$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = \mathbf{w} - \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i = 0$$

which imply $\sum_{i=1}^{m} \alpha_i y_i = 0$ and $\mathbf{w} = \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i$

• Recall from Chapter 6 that for an optimal feasible solution $\bar{\mathbf{w}}$, $\alpha_i c_i(\bar{\mathbf{w}}, \bar{b}) = 0$ for all constraints c_i , so

$$\alpha_i(y_i(\langle \mathbf{x}_i, \bar{\mathbf{w}} \rangle + \bar{b}) - 1) = 0$$

for all $i = 1, \ldots, m$

The Dual Optimization Problem (cont'd)

The Dual Optimization Problem (cont'd)

- The x_i for which $\alpha_i > 0$ are the support vectors, and are the vectors that lie on the margin, i.e. those for which the constraints are tight
 - Other vectors (where $\alpha_i = 0$) are irrelevant to determining the hyperplane w
 - Will be useful later in classification
 - See Prop. 7.8 for relationship between expected number of SVs and test error bound

- Now substitute the saddle point conditions into the Lagrangian
- The $k{\rm th}$ component of the weight vector is $w_k = \sum_{i=1}^m \alpha_i y_i x_{ik},$ so

$$w_k^2 = \left(\sum_{i=1}^m \alpha_i y_i x_{ik}\right) \left(\sum_{i=1}^m \alpha_i y_i x_{ik}\right)$$

• Thus

$$\|\mathbf{w}\|^{2} = \sum_{k} \left(\sum_{i=1}^{m} \alpha_{i} y_{i} x_{ik} \right) \left(\sum_{i=1}^{m} \alpha_{i} y_{i} x_{ik} \right)$$
$$= \sum_{k} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} x_{ik} x_{jk}$$
$$= \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \sum_{k} x_{ik} x_{jk}$$
$$= \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \langle \mathbf{x}_{i}, \mathbf{x}_{j} \rangle$$

14

The Dual Optimization Problem (cont'd)

• Further,

$$\sum_{i=1}^{m} \alpha_i (y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1)$$
$$= \sum_{i=1}^{m} \alpha_i y_i \left(\sum_k x_{ik} w_k \right) - \sum_{i=1}^{m} \alpha_i$$
$$= \sum_{i=1}^{m} \alpha_i y_i \left(\sum_k x_{ik} \sum_{j=1}^{m} \alpha_j y_j x_{jk} \right) - \sum_{i=1}^{m} \alpha_i$$
$$= \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^{m} \alpha_i$$

• Combine them:

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

The Dual Optimization Problem (cont'd)

• Maximizing the Lagrangian wrt α yields the dual optimization problem:

$$\begin{array}{ll} \underset{\boldsymbol{\alpha} \in \mathbb{R}^{m}}{\text{maximize}} & \sum_{i=1}^{m} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \langle \mathbf{x}_{i}, \mathbf{x}_{j} \rangle \\ \text{s.t.} & \alpha_{i} \geq 0, i = 1, \dots, m \\ & \sum_{i=1}^{m} \alpha_{i} y_{i} = 0 \end{array}$$

$$(2)$$

• After optimization, we can label new vectors with the decision function:

$$f(\mathbf{x}) = \operatorname{sgn}\left(\sum_{i=1}^{m} \alpha_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b\right)$$

(later we'll discuss finding b)

Adding Kernels

- As discussed before, using kernels is an effective way to introduce nonlinearities to the data
 - Nonlinear remapping might make data (almost) linearly separable in the new space
 - Cover's theorem implies that simply increasing the dimension improves the probability of linear separability
- For given remapping Φ , simply replace x with $\Phi(x)$
- Thus in dual optimization problem and in decision function, replace $\langle \mathbf{x}, \mathbf{x}_i \rangle$ with $k(x, x_i)$, where k is the PD kernel corresponding to Φ
- If k is PD, then we still have a convex optimization problem
- Once α is found, can e.g. set b to be the average over all $\alpha_j > 0$ of $y_j \sum_{i=1}^m y_i \alpha_i k(x_j, x_i)$ (derived from KKT conditions)

17

Soft Margin Hyperplanes

(cont'd)

• To relax each constraint from (1), add <u>slack</u> <u>variable</u> $\xi_i \ge 0$:

 $y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \ge 1 - \xi_i, \ i = 1, \dots, m$

- Also need to penalize large ξ_i in the objective function to prevent trivial solutions
 - C-SV classifier
 - v-SV classifier

Soft Margin Hyperplanes

- Under a given mapping Φ, the data might not be linearly separable
- There always exists a Φ that will yield separability, but is it a good idea to find one just for the sake of separating?
- If we choose to keep the mapping that corresponds to our favorite kernel, what are our options?
 - Instead of finding a hyperplane that is perfect on the training set, find one that minimizes training errors
 - Computationally intractable to even approximate
 - Instead, we'll <u>soften</u> the margin, allowing for some vectors to get too close to the hyperplane (i.e. margin errors)

18

Soft Margin Hyperplanes

C-SV Classifier

 Weight with C > 0 (e.g. C = 10m) the importance of minimizing sum of ξ variables:

 $\begin{array}{ll} \underset{\mathbf{w} \in \mathcal{H}, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^m}{\text{minimize}} & \tau(\mathbf{w}, \boldsymbol{\xi}) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{m} \sum_{i=1}^m \xi_i \\ \text{s.t.} & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \geq 1 - \xi_i, \ i = 1, \dots, m \\ & \xi_i \geq 0, \ i = 1, \dots, m \end{array}$

- First term of τ decreases ||w||, second term focuses on margin error rate ν, thus together they focus on T7.3
- The dual is similar to that for hard margin:

$$\begin{array}{ll} \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{maximize}} & W(\boldsymbol{\alpha}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{s.t.} & 0 \leq \alpha_i \leq C/m, \ i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{array}$$

Once α is found, can e.g. set b to be the average over all α_j ∈ (0, C) of y_j − Σ^m_{i=1} y_iα_ik(x_j, x_i)

Soft Margin Hyperplanes u-SV Classifier (cont'd)

Soft Margin Hyperplanes

 ν -SV Classifier

- A more intuitable way to weight the emphasis on reducing margin errors
- Primal:

 $\begin{array}{ll} \underset{\mathbf{w} \in \mathcal{H}, \rho, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^m}{\text{minimize}} & \tau(\mathbf{w}, \boldsymbol{\xi}, \rho) = \frac{1}{2} \|\mathbf{w}\|^2 - \nu\rho + \frac{1}{m} \sum_{i=1}^m \xi_i \\ \text{s.t.} & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \ge \rho - \xi_i, i = 1, \dots, m \\ & \rho \ge 0, \xi_i \ge 0, \ i = 1, \dots, m \end{array}$

• ρ is similar to that in T7.3: for ξ to be 0, all vectors must be at least $\rho/||\mathbf{w}||$ from the hyperplane

21

Soft Margin Hyperplanes *v*-SV Classifier

(cont'd)

• Derivation of dual form (details omitted) yields:

$$\begin{array}{ll} \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{maximize}} & W(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{s.t.} & 0 \leq \alpha_i \leq 1/m, \\ & \sum_{i=1}^m \alpha_i y_i = 0 \\ & \sum_{i=1}^m \alpha_i \geq \nu \end{array}$$

• Let S_+ and S_- be sets of SVs x_i with labels $y_i = +1$ and -1, $0 < \alpha_i < 1$, and $|S_+| = |S_-| = s > 0$, then set

$$b = -\frac{1}{2s} \sum_{x \in S_+ \cup S_-} \sum_{i=1}^m y_i \alpha_i k(x, x_i)$$

P7.5 If ν -SVC yields a solution with $\rho > 0$, then

- 1. ν is an upper bound on the fraction of margin errors
- 2. ν is a lower bound on the fraction of support vectors
- See Table 7.1, p. 207



Multi-Class Classification

- What if we want to go beyond binary labels ± 1 to M classes?
- Most methods decompose a multi-class problem into a set of binary ones
 - One vs. rest
 - Error-correcting output codes
 - Pairwise classification
 - Kessler's construction/multi-class objective function (doesn't need to decompose into binary cases)

- To handle M classes, train a set of M binary classifiers f^1, \ldots, f^M , where f^i is trained to distinguish patterns from class i from those not in class i
- If (α^i, b^i) is the classifier learned for class i, then a new pattern x is classified as

$$\underset{j=1,\dots,M}{\operatorname{argmax}} \left\{ \sum_{i=1}^{m} y_i \alpha_i^j k(x, x_i) + b^j \right\}$$

i.e. the class with the most confident prediction among the binary classifiers

- Applicable even if the number of classifiers predicting +1 is not exactly 1
- Note that the set of SVs can be different for each class
- Can also let the classifier "punt" if the difference between the top two predictions is small

25

Multi-Class Classification

Error-Correcting Output Codes (ECOC)

- One vs. rest requires *M* classifiers to represent *M* classes
- Is this the minimum amount required?
- E.g. M = 4, so use two linear classifiers:

Class	Binary Encoding	
	Classifier 1	Classifier 2
Class 1	$^{-1}$	-1
Class 2	-1	+1
Class 3	+1	-1
Class 4	+1	+1

and train simultaneously

• <u>Problem</u>: Sensitive to individual classifier errors, so use a <u>set of encodings</u> per class to improve robustness

26

Multi-Class Classification

Error-Correcting Output Codes (ECOC) (cont'd)

- Similar to principle of <u>error-correcting output</u> <u>codes</u> used in communication networks
 - After all classifiers make their predictions, find the code that is nearest to the bit string returned and use that for the predicted class
- Can provably tolerate some mispredictions by individual classifiers, but doesn't use the margin

Multi-Class Classification

Pairwise Classification

- Instead of training one classifier per class as in one vs. rest, train a classifier for each <u>pair</u> of classes
- Now have $\binom{M}{2}$ classifiers to train rather than $\lceil \log_2 M \rceil$ up to M, but each training set is smaller
 - Number of SVs smaller for each classifier due to smaller training set and easier learning problem
- To classify new pattern, evaluate it on all classifiers and choose the class that gets the most votes
 - Can avoid running on all classifiers if votes so far imply that some classes are guaranteed to not win



• All have comparable min error rates, but sensitive to parameter setting

Parameter Setting

• Gaussian kernel: low, med, high values of c



- How to choose parameter settings?
 - Cross validation
 - Settings that work well for similar problems (rescaled)
 - For ν -SVCs, set ν to e.g. test error from other classifiers
 - * $\nu \geq$ margin error \geq train error, which is also \leq test error
 - For $C\text{-}\mathsf{SVCs},\ C\propto 1/R^2,$ where R measures range of data in $\mathcal H$
 - * E.g. R = radius of smallest sphere, max or mean length $k(x_i, x_i)$, or std dev of distance of points to their mean

Overlap of SV Sets

- In the handwritten digit classification experiments, the three kernels typically had 80–93% of their SV sets in common (Table 7.6, p. 220)
- In fact, each kernel got similar error rates when training on SVs of a different kernel rather than the entire training set (Table 7.7)
- Basically, these kernels (dot products) mostly found the same regularities in the data
- Results might vary depending on learning problem/data set

34

Topic summary due in 1 week!