

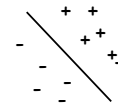
## 15-859(B) Machine Learning Theory

Lecture 11: More on why large margins are good for learning. Kernels and general similarity functions.  $L_1 - L_2$  connection.

Avrim Blum  
02/18/09

### Basic setting

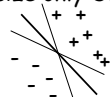
- w Examples are points  $x$  in instance space, like  $\mathbb{R}^n$ .
- w Labeled + or -.
- w Assume drawn from some probability distribution:
  - Distribution  $D$  over  $x$ , labeled by target function  $c$ .
  - Or distribution  $P$  over  $(x, l)$
  - Will call  $P$  (or  $(c, D)$ ) our "learning problem".
- w Given labeled training data, want algorithm to do well on new data.



### Margins

If data is separable by large margin  $\gamma$ , then that's a good thing. Need sample size only  $\tilde{O}(1/\gamma^2)$ .

$$|w \cdot x| / |x| \geq \gamma, |w| = 1$$



Some ways to see it:

1. The perceptron algorithm does well: makes only  $1/\gamma^2$  mistakes.
2. Margin bounds: whp all consistent large-margin separators have low true error.
3. Really-Simple-Learning + boosting...
4. Random projection... Will do 3 then 4 then 2.

### A really simple learning algorithm

Suppose our problem has the property that whp a sufficiently large sample  $S$  would be separable by margin  $\gamma$ . Here is another way to see why this is good for learning.

Consider the following simple algorithm...

1. Pick a random hyperplane.
2. See if it is any good.
3. If it is a weak-learner (error rate  $\leq \frac{1}{2} - \gamma/4$ ), plug into boosting. Else don't. Repeat.

**Claim:** if data has a large margin separator, there's a reasonable chance a random hyperplane will be a weak-learner.

### A really simple learning algorithm

Claim: if data has a separator of margin  $\gamma$ , there's a reasonable chance a random hyperplane will have error  $\leq \frac{1}{2} - \gamma/4$ . [all hyperplanes through origin]

Proof:

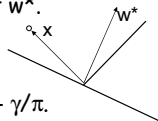
- w Pick a (positive) example  $x$ . Consider the 2-d plane defined by  $x$  and target  $w^*$ .

$$\Pr_h(h \cdot x \leq 0 \mid h \cdot w^* \geq 0) \leq (\pi/2 - \gamma)/\pi = \frac{1}{2} - \gamma/\pi.$$

$$\text{So, } E_h[\text{err}(h) \mid h \cdot w^* \geq 0] \leq \frac{1}{2} - \gamma/\pi.$$

- w Since  $\text{err}(h)$  is bounded between 0 and 1, there must be a reasonable chance of success.

QED



### Another way to see why large margin is good

Johnson-Lindenstrauss Lemma:

Given  $n$  points in  $\mathbb{R}^n$ , if project randomly to  $\mathbb{R}^k$ , for  $k = O(\epsilon^{-2} \log n)$ , then whp all pairwise distances preserved up to  $1 \pm \epsilon$  (after scaling by  $(n/k)^{1/2}$ ).

Cleanest proofs: IM98, DG99

## JL Lemma

Given  $n$  points in  $\mathbb{R}^n$ , if project randomly to  $\mathbb{R}^k$ , for  $k = O(\epsilon^{-2} \log n)$ , then whp all pairwise distances preserved up to  $1 \pm \epsilon$  (after scaling).

Cleanest proofs: IM98, DG99

### Proof intuition:

- w Consider a random unit-length vector  $(x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ . What does  $x_1$  coordinate look like?
- w  $E[x_1^2] = 1/n$ . Usually  $\leq c/n$ .
- w If indep,  $\Pr[|(x_1^2 + \dots + x_k^2) - k/n| \geq \epsilon k/n] \leq e^{-O(k\epsilon^2)}$ .
- w So, at  $k = O(\epsilon^{-2} \log n)$ , with prob  $1 - 1/\text{poly}(n)$ , projection to 1st  $k$  coordinates has length  $(k/n)^{1/2} (1 \pm \epsilon)$ .
- w Now, apply this to vector  $v_{ij} = p_i - p_j$ , projecting onto random  $k$ -diml space.  
Whp all  $v_{ij}$  project to length  $(k/n)^{1/2}(1 \pm \epsilon)|v_{ij}|$

## JL Lemma, cont

Proof easiest for slightly different projection:

- w Pick  $k$  vectors  $u_1, \dots, u_k$  iid from  $n$ -diml gaussian.
- w Map  $p \rightarrow (p \cdot u_1, \dots, p \cdot u_k)$ .
- w What happens to  $v_{ij} = p_i - p_j$ ?
  - n Becomes  $(v_{ij} \cdot u_1, \dots, v_{ij} \cdot u_k)$
  - n Each component is iid from 1-diml gaussian, scaled by  $|v_{ij}|$ .
  - n For concentration on sum of squares, plug in version of Hoeffding for RVs that are squares of gaussians.
- w So, whp all lengths apx preserved, and in fact not hard to see that whp all angles are apx preserved too.

## Random projection and margins

Natural connection [AV99]:

- w Suppose we have a set  $S$  of points in  $\mathbb{R}^n$ , separable by margin  $\gamma$ .
- w JL lemma says if project to random  $k$ -dimensional space for  $k = O(\gamma^{-2} \log |S|)$ , whp still separable (by margin  $\gamma/2$ ).
  - n Think of projecting points and target vector  $w$ .
  - n Angles between  $p_i$  and  $w$  change by at most  $\pm \gamma/2$ .
- w Could have picked projection before sampling data.
- w So, it's really just a  $k$ -dimensional problem after all. Do all your learning in this  $k$ -diml space.

So, random projections can help us think about why margins are good for learning. [note: this argument does NOT imply uniform convergence in original space]

## Uniform convergence bounds for large margins

Claim: Whp, any linear separator that gets training data correct by margin  $\gamma$  has true error  $\leq \epsilon$  so long as  $|S| \gg (1/\epsilon)[(1/\gamma^2)\log^2(1/(\gamma\epsilon)) + \log(1/\delta)]$

Proof in two steps:

1. What is the maximum number of points that can be shattered by separators of margin at least  $\gamma$ ? (aka "fat-shattering dimension")
  - n Ans:  $O(1/\gamma^2)$ .
  - n Proof: corollary to Perceptron mistake bound (why?) (if dimension is  $d$ , can force Perceptron to make  $\geq d$  mistakes)
2. Now want to use like in VC-dim analysis. Sauer's lemma analog still applies, but there's a complication we'll need to address...

## Uniform convergence bounds for large margins

Claim: Whp, any linear separator that gets training data correct by margin  $\gamma$  has true error  $\leq \epsilon$  so long as  $|S| \gg (1/\epsilon)[(1/\gamma^2)\log^2(1/(\gamma\epsilon)) + \log(1/\delta)]$

Proof in two steps:

2. Now want to use like in VC-dim analysis. Sauer's lemma analog still applies, but there's a complication we'll need to address...
  - n Draw  $2m$  points from  $D$ , split into  $S_1, S_2$  as before.
  - n Argue whp no separator gets  $S_1$  correct by margin  $\gamma$ , but makes  $\geq \epsilon m$  mistakes on  $S_2$ .
  - n To do this, tempting to do union bound over all separators that have no points in  $S$  within margin  $\gamma$  (which we can count using Sauer)
  - n But this is undercounting...

## Uniform convergence bounds for large margins

Claim: Whp, any linear separator that gets training data correct by margin  $\gamma$  has true error  $\leq \epsilon$  so long as  $|S| \gg (1/\epsilon)[(1/\gamma^2)\log^2(1/(\gamma\epsilon)) + \log(1/\delta)]$

Proof in two steps:

2. Now want to use like in VC-dim analysis. Sauer's lemma analog still applies, but there's a complication we'll need to address...
  - n Let  $h(x) = h \cdot x$ , but truncated at  $\pm \gamma$ .
  - n Define  $\text{dist}(h_1, h_2) = \max_{x \in S} |h_1(x) - h_2(x)|$ .
  - n Define  $H$  to be a " $\gamma/2$  cover": for all separators, exists  $h \in H$  within distance  $\gamma/2$ .
  - n For  $h \in H$ , define "correct" as "correct by margin at least  $\gamma/2$ ", else call it a "mistake". Now, run usual union-bound argument on these.
  - n Finally, apply bound of [Alon et al] on cover-sizes...

OK, now on to kernels...

## Kernel functions

W We have a lot of great algorithms for learning linear separators (perceptron, SVM, ...). But, a lot of time, data is not linearly separable.

in "Old" answer: use a multi-layer neural network.

in "New" answer: use a kernel function!

W Many algorithms only interact with the data via dot-products.

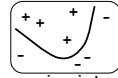
in So, let's just re-define dot-product.

in E.g.,  $K(x,y) = (1 + x \cdot y)^d$ .

-  $K(x,y) = \phi(x) \cdot \phi(y)$ , where  $\phi()$  is implicit mapping into an  $n^d$ -dimensional space.

in Algorithm acts as if data is in " $\phi$ -space". Allows it to produce non-linear curve in original space.

in Don't have to pay for high dimension if data is linearly separable there by a large margin.



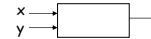
Question: do we need the notion of an implicit space to understand what makes a kernel helpful for learning?

[BB'06][BBS'08]

## Kernel fns have become very popular

...but there's something a little funny:

W On the one hand, operationally a kernel is just a similarity function:  $K(x,y) \in [-1,1]$ , with some extra requirements. [here I'm scaling to  $|\Phi(x)| = 1$ ]



W And in practice, people think of a good kernel as a good measure of similarity between data points.

W But Theory talks about margins in implicit high-dimensional  $\Phi$ -space.  $K(x,y) = \Phi(x) \cdot \Phi(y)$ .

I want to use ML to classify protein structures and I'm trying to decide on a similarity fn to use. Any help?

It should be pos. semidefinite, and should result in your data having a large margin separator in implicit high-diml space you probably can't even calculate.

Umm... thanks, I guess.

It should be pos. semidefinite, and should result in your data having a large margin separator in implicit high-diml space you probably can't even calculate.

Can we develop a more intuitive theory?

w Match intuition that you are looking for a good measure of similarity for the problem at hand?

w Get the power of the standard theory with less of "something for nothing" feel to it?

Can we develop a more intuitive theory?

What would we intuitively want in a good measure of similarity?

### A reasonable idea:

w Say have a learning problem  $P$  (distribution  $D$  over examples labeled by unknown target  $f$ ).

w Sim fn  $K: (\text{img}, \text{img}) \rightarrow [-1,1]$  is good for  $P$  if:  
most  $x$  are on average more similar to random pts of their own label than to random pts of the other label, by some gap  $\gamma$ .

E.g., most images of men are on average  $\gamma$ -more similar to random images of men than random images of women, and vice-versa.

### A reasonable idea:

w Say have a learning problem  $P$  (distribution  $D$  over examples labeled by unknown target  $f$ ).

w Sim fn  $K:(x,y) \rightarrow [-1,1]$  is  $(\epsilon, \gamma)$ -good for  $P$  if at least a  $1-\epsilon$  fraction of examples  $x$  satisfy:

$$E_{y \sim D}[K(x,y) | \ell(y) = \ell(x)] \geq E_{y \sim D}[K(x,y) | \ell(y) \neq \ell(x)] + \gamma$$

E.g., most images of men are on average  $\gamma$ -more similar to random images of men than random images of women, and vice-versa.

### A reasonable idea:

w Say have a learning problem  $P$  (distribution  $D$  over examples labeled by unknown target  $f$ ).

w Sim fn  $K:(x,y) \rightarrow [-1,1]$  is  $(\epsilon, \gamma)$ -good for  $P$  if at least a  $1-\epsilon$  fraction of examples  $x$  satisfy:

$$E_{y \sim D}[K(x,y) | \ell(y) = \ell(x)] \geq E_{y \sim D}[K(x,y) | \ell(y) \neq \ell(x)] + \gamma$$

How can we use it?

### Just do "average nearest-nbr"

At least a  $1-\epsilon$  fraction of  $x$  satisfy:

$$E_{y \sim D}[K(x,y) | \ell(y) = \ell(x)] \geq E_{y \sim D}[K(x,y) | \ell(y) \neq \ell(x)] + \gamma$$

w Draw  $S^+$  of  $O((1/\gamma^2) \ln 1/\delta^2)$  positive examples.

w Draw  $S^-$  of  $O((1/\gamma^2) \ln 1/\delta^2)$  negative examples

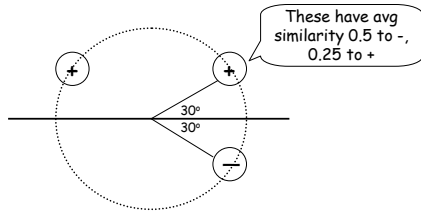
w Classify  $x$  based on which gives better score.

▫ Hoeffding: for any given "good  $x$ ", prob of error over draw of  $S^+, S^-$  at most  $\delta^2$ .

▫ So, at most  $\delta$  chance our draw is bad on more than  $\delta$  fraction of "good  $x$ ".

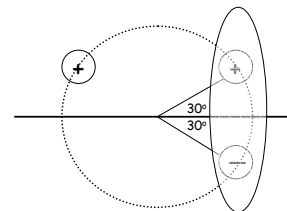
w With prob  $\geq 1-\delta$ , error rate  $\leq \epsilon + \delta$ .

### But not broad enough



- w  $K(x,y)=x \cdot y$  has good separator but doesn't satisfy defn. (half of positives are more similar to negs than to typical pos)

### But not broad enough



- w Idea: would work if we didn't pick y's from top-left.
- w Broaden to say: OK if  $\exists$  large region R s.t. most x are on average more similar to  $y \in R$  of same label than to  $y \in R$  of other label. (even if don't know R in advance)

### Broader defn...

- w Ask that exists a set R of "reasonable" y (allow probabilistic) s.t. almost all x satisfy

$$E_y[\ell(x)\ell(y)K(x,y)|y \in R] \geq \gamma$$

- w Formally, say K is  $(\epsilon', \gamma, \tau)$ -good if have hinge-loss  $\epsilon'$ , and  $\Pr(R) \geq \tau$ .
- w **Thm 1:** this is a legitimate way to think about good kernels:
  - n If kernel has margin  $\gamma$  in implicit space, then for any  $\tau$  is  $(\tau, \gamma^2, \tau)$ -good in this sense. [BBS'08]

### Broader defn...

- w Ask that exists a set R of "reasonable" y (allow probabilistic) s.t. almost all x satisfy

$$E_y[\ell(x)\ell(y)K(x,y)|y \in R] \geq \gamma$$

- w Formally, say K is  $(\epsilon', \gamma, \tau)$ -good if have hinge-loss  $\epsilon'$ , and  $\Pr(R) \geq \tau$ .
- w **Thm 2:** even if not a legal kernel, this is nonetheless sufficient for learning.
  - n If K is  $(\epsilon', \gamma, \tau)$ -good,  $\epsilon' \ll \epsilon$ , can learn to error  $\epsilon$  with  $O((1/\epsilon\gamma^2)\log(1/\epsilon\gamma\tau))$  labeled examples. [and  $\tilde{O}(1/(\gamma^2\tau))$  unlabeled examples]

### How to use such a sim fn?

- w Assume exists R s.t.  $\Pr_y[R] \geq \tau$  and almost all x satisfy

$$E_y[\ell(x)\ell(y)K(x,y)|y \in R] \geq \gamma$$

- n Draw  $S = \{y_1, \dots, y_n\}$ ,  $n \approx 1/(\gamma^2\tau)$ . could be unlabeled
- n View as "landmarks", use to map new data:
 
$$F(x) = [K(x, y_1), \dots, K(x, y_n)]$$
- n Whp, exists separator of good  $L_1$  margin in this space:  $w = [0, 0, 1/n_R, 1/n_R, 0, 0, 0, -1/n_R, 0]$  ( $n_R = \# y_i \in R$ )
- n So, take new set of examples, project to this space, and run good  $L_1$  alg (Winnow).

### Other notes

- w So, large margin in implicit space  $\Rightarrow$  satisfy this defn (with potentially quadratic penalty in margin).
- w This def is really an  $L_1$  style margin, so can also potentially get exponential improvement.
  - n Much like Winnow versus Perceptron.
  - n Can construct class C s.t. for any kernel K, some  $f \in C$  has  $L_2$  margin only  $O(1/|C|^{1/2})$  but there exists a similarity fn satisfying above def with  $\gamma=1$  and  $\tau=1/|C|$ .
- w Interesting to consider other natural properties of similarity functions that motivate other algs.