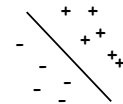# 15-859(B) Machine Learning Theory

Lecture 9: Margins, kernels, and similarity functions

Avrim Blum
02/13/08

---

## Basic Supervised learning setting
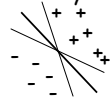
- w Examples are points $x$ in instance space, like $R^n$.
- w Labeled + or -.
- w Assume drawn from some probability distribution:
  - n Distribution D over $x$, labeled by target function c.
  - n Or distribution P over $(x, l)$
  - n Will call P (or (c,D)) our "learning problem".
- w Given labeled 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 $\tilde{O}(1/\gamma^2)$.

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

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…

---

## 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]
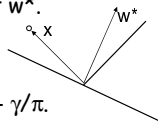
Proof:
- w Pick a (positive) example $x$. Consider the 2-d plane defined by $x$ and target $w^*$.
- w $\Pr_h(h \cdot x \leq 0 \mid h \cdot w^* \geq 0)$
  $\leq (\pi/2 - \gamma)/\pi = \frac{1}{2} - \gamma/\pi$.
- w So, $E_h[err(h) \mid h \cdot w^* \geq 0] \leq \frac{1}{2} - \gamma/\pi$.
- w Since $err(h)$ 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 $R^n$, if project randomly to $R^k$, for $k = O(\varepsilon^{-2} \log n)$, then whp all pairwise distances preserved up to $1 \pm \varepsilon$ (after scaling by $(n/k)^{1/2}$).

Cleanest proofs: IM98, DG99

---

## JL Lemma

Given n points in $R^n$, if project randomly to $R^k$, for $k = O(\varepsilon^{-2} \log n)$, then whp all pairwise distances preserved up to $1\pm\varepsilon$ (after scaling).
Cleanest proofs: IM98, DG99

### Proof intuition:

- w Consider a random unit-length vector $(x_1, x_2, ..., x_n) \in R^n$. What does $x_1$ coordinate look like?
- w $E[x_1^2] = 1/n$. Usually $\leq c/n$.
- w If indep, $Pr[|(x_1^2 + ... + x_k^2) - k/n| \geq \varepsilon k/n] \leq e^{-O(k\varepsilon^2)}$.
- w So, at $k = O(\varepsilon^{-2} \log n)$, with prob $1 - 1/poly(n)$, projection to 1st $k$ coordinates has length $(k/n)^{1/2} (1 \pm \varepsilon)$.
- 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 $(k/n)^{1/2}(1\pm\varepsilon)|v_{ij}|$

## JL Lemma, cont

Proof easiest for slightly different projection:
- w Pick $k$ vectors $u_1, ..., u_k$ iid from n-diml gaussian.
- w Map $p \rightarrow (p \cdot u_1, ..., p \cdot u_k)$.
- w What happens to $v_{ij} = p_i - p_j$?
  - n Becomes $(v_{ij} \cdot u_1, ... , v_{ij} \cdot u_k)$
  - 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 <u>angles</u> are apx preserved too.

## Random projection and margins

Natural connection [AV99]:
- w Suppose we have a set S of points in $R^n$, 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]

OK, now on to kernels…

## Generic problem

- w Given a set of images:  , want to train a classifier to distinguish men from women.
- w Problem: pixel representation not good for LTFs

Classic advice:
- w Use a complicated neural net.
- w But these are hard to train.

Modern advice:
- w Use a Kernel! $K($ ,  $) = \Phi($  $) \cdot \Phi($  $)$. $\Phi$ is implicit, high-dimensional mapping.
- w Many algorithms only interact with data through dot-products, so can be "kernelized". If data is separable in $\Phi$-space by large margin, don't have to pay for dim.

## Generic problem

Modern advice:
- w Use a Kernel! $K($ ,  $) = \Phi($  $) \cdot \Phi($  $)$. $\Phi$ is implicit, high-dimensional mapping.
- w Many algorithms only interact with data through dot-products, so can be "kernelized". If data is separable in $\Phi$-space by large margin, don't have to pay for dim.
- w E.g., $K(x,y) = (1+x_1y_1)(1+x_2y_2)...(1+x_ny_n)$.
  - n $\Phi$:(n-diml space) $\rightarrow$ ($2^n$-diml space).
- w Conceptual warning: You're not really "getting all the power of the high dimensional space without paying for it". (Not enough to just be separable. Need large margin too.) As we saw from JL lemma, assumption of large margin means it's really an $\tilde{O}(1/\gamma^2)$-dimensional problem after all.
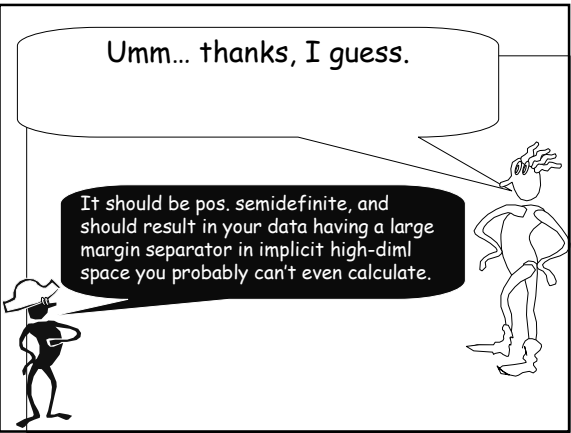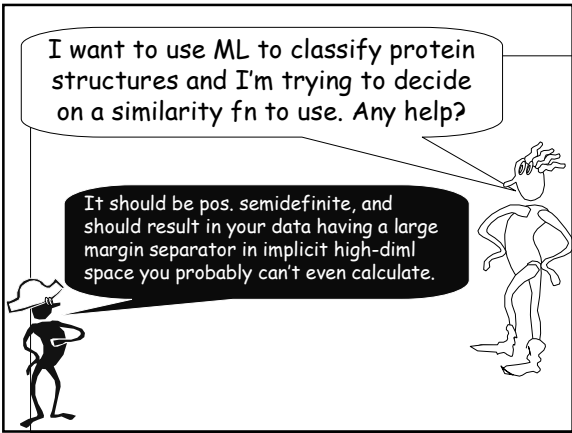
## Slide 1

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

## Slide 2

### 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 $|\Phi(x)| = 1$]

x →
y → [ ] →

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

w But <u>Theory</u> talks about margins in implicit high-dimensional $\Phi$-space. $K(x,y) = \Phi(x)\cdot\Phi(y)$.

## Slide 3

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.

## Slide 4

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.

## Slide 5

### 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 $|\Phi(x)| = 1$]

x →
y → [ ] →

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

w But <u>Theory</u> talks about margins in implicit high-dimensional $\Phi$-space. $K(x,y) = \Phi(x)\cdot\Phi(y)$.

Can we bring these views together?

## Slide 6

### Goal: notion of "good similarity function" that…

1. Talks in terms of more intuitive properties (no implicit high-diml spaces…)

2. If K satisfies these properties for our given problem, then has implications to learning

3. Is broad: includes usual notion of "good kernel" (one that induces a large margin separator in $\Phi$-space).

[Recent work with Nina Balcan, with extensions by Nati Srebro]

## Defn satisfying (1) and (2):

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

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

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

w Note: you can have this property without being a legal kernel.
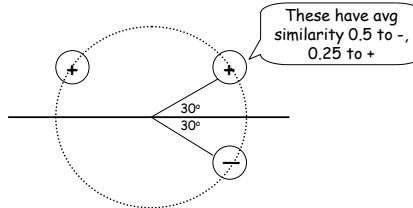
w Q: how could you use this to learn?

---

## How to use it

At least a $1-\epsilon$ prob mass of x satisfy:
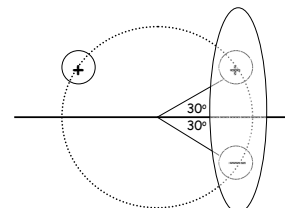$$E_{y\sim D}[K(x,y)|\ell(y)=\ell(x)] \geq E_{y\sim D}[K(x,y)|\ell(y)\neq\ell(x)]+\gamma$$

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.

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

   n 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·y has good separator but doesn't satisfy defn. (half of positives are more similar to negs that 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 ∃ large region R s.t. most x are on average more similar to y∈R of same label than to y∈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

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

w And at least $\epsilon$ probability mass of reasonable positives/negatives.

w But now, how can we use for learning??

---

## Broader defn...

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

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

   n Draw $S = \{y_1,...,y_n\}$, $n\approx 1/(\gamma^2\epsilon)$. _(could be unlabeled)_

   n View as "landmarks", use to map new data: $F(x) = [K(x,y_1), ...,K(x,y_n)]$.

   n Whp, exists separator of good $L_1$ margin in this space: $w=[0,0,1/n_+,1/n_+,0,0,0,-1/n_-,0,0]$

   n $n_+$ [$n_-$] = # reasonable pos [neg] in S.

   n So, take new set of examples, project to this space, and run good $L_1$ alg (Winnow).

## And furthermore

Now, defn is broad enough to include all large margin kernels (some loss in parameters):

- $\gamma$-good margin $\Rightarrow$ apx $(\varepsilon, \gamma^2, \varepsilon)$-good here.

But now, we don't need to think about implicit spaces or require kernel to even have the implicit space interpretation.

If PSD, can also show reverse too:

- $\gamma$-good here & PSD $\Rightarrow$ $\gamma$-good margin.