

# VC Dimension

Note Title

4/5/2010

Additional notes on VC dimension (from Stat 232)

Some changes in notation

VC-dim  $h$

Risk - Error

These notes contain:

- The VC bound theorem.
- Discussion of the capacity term  $\phi(h, n, \delta)$
- VC's for margins and kernels  
(too advanced, kernels haven't been mentioned yet.)
- The "trends" of VC bounds.
- Structural risk minimization

## VC-Dimension

Note Title

11/19/2006

$$\underline{\text{Risk}}: R(\alpha) = \sum_{x,y} L(\alpha(x), y) P(x, y)$$

$$\underline{\text{Empirical Risk}} \quad R_{\text{emp}}(\alpha) = \frac{1}{N} \sum_{i=1}^N L(\alpha(x_i), y_i)$$

Minimize the empirical risk by constraining the set  $\mathcal{M}$  of classifiers.

Vapnik-Chervonenkis (VC) dimension of a set  $\mathcal{M}$  of classifiers.

VC = largest number of datapoints which can be shattered by the classifier set.

Shattered means that all possible dichotomies of the dataset can be expressed by a classifier in the set.

E.g. Hyperplanes can shatter any set provided  $n \leq d+1$ ,  $d = \text{dim space}$

Hence, VC-dimension of hyperplanes is  $d+1$ .

## (2) Probably Approximately Correct (PAC)

Let  $h$  be the VC dimension of a set  $\Gamma$  of classifiers.

If number of data samples  $n \leq h$ , then learning (i.e. generalization) is impossible - you can only memorize.

VC Bound Theorem: Suppose  $n > h$

then, with probability at least  $1 - \delta$ :

$$R[\alpha] \leq R_{\text{emp}}[\alpha] + \phi(h, n, \delta), \quad \forall \alpha \in \Gamma$$

where  $\phi(h, n, \delta) = \sqrt{\left(\frac{8}{n}\right) \left\{ h \log\left(\frac{n}{h}\right) + h + \log\left(\frac{4}{\delta}\right) \right\}}$

(Several different  $\phi$ 's appear in different proofs)

Comment: require  $\phi(h, n, \delta)$  to be small to ensure generalization.

requires:  $n \gg h$        $n \gg |\log \delta|$        $\left\| \begin{array}{l} \text{small } h/n \\ \text{small } |\log \delta|/n \end{array} \right.$

note: to require certainty,  $\delta \rightarrow 0 \Rightarrow n \rightarrow \infty$ .

(3) Comment: the derivation of VC bounds is complicated. The bounds require many approximations — so they are not tight.

Also, the proofs require ruling out all other possible interpretations of the data → this is too conservative.

### VC Margins for Hyperplanes

Consider hyperplanes,  $\underline{a} \cdot \underline{x} = 0$ ,  
normalized wrt. data  $\{x_\mu : \mu = 1, \dots, n\}$   
s.t.  $\min_{\mu \in \{1, \dots, n\}} |\underline{a} \cdot x_\mu| = 1$ .

Then, the set of hyperplane classifiers of form  $\text{sign}(\underline{x} \cdot \underline{a})$ , with  $|\underline{a}| < \Lambda$

has VC dimension  $h < R^2 \Lambda^2$   $\overset{x_i}{\text{inverse margin}}$

where  $R^2$  is radius of smallest sphere containing data points.

Intuition: low VC, if we require large margin  $\nearrow$  region & data lies in small region

(4) The same result also applies to kernels, because the kernel trick can be used.

The margin  $\Lambda$ , depends only on the kernel.  
The radius of the minimum sphere  $R$  depends only on the kernel.

Proof, the minimum sphere can be found by a quadratic optimization problem.

$$L_p(R, \underline{x}^*) = R^2 - \sum_{\mu=1}^n \lambda_{\mu} (R^2 - |x_{\mu} - \underline{x}^*|)$$

$$\lambda_{\mu} \geq 0.$$

$\lambda_{\mu}$  Lagrange term enforces

$$R \geq |x_{\mu} - \underline{x}^*| \quad \forall \mu$$

$$\frac{\partial L_p}{\partial \underline{x}^*} = 0 \Rightarrow \sum_{\mu} \lambda_{\mu} (x_{\mu} - \underline{x}^*) = 0$$

$$\frac{\partial L_p}{\partial R} = 0 \Rightarrow$$

$$1 = \sum_{\mu=1}^n \lambda_{\mu}.$$

$$\Rightarrow \sum_{\mu} \lambda_{\mu} x_{\mu} - \underline{x}^* = \sum_{\mu} \lambda_{\mu} \underline{x}^* - \underline{x}^*.$$

Hence: 
$$L_d = \sum_{\mu=1}^n \lambda_{\mu} (x_{\mu} - \underline{x}^*)^2$$

$$L_d = \sum_{\mu=1}^n \lambda_{\mu} x_{\mu} \cdot x_{\mu} - 2 \sum_{\mu=1}^n \lambda_{\mu} x_{\mu} \cdot \underline{x}^* + \sum_{\mu=1}^n \lambda_{\mu} (\underline{x}^*)^2$$

$$= \sum_{\mu=1}^n \lambda_{\mu} x_{\mu} \cdot x_{\mu} - \sum_{\mu, \nu=1}^n \lambda_{\mu} \lambda_{\nu} x_{\mu} \cdot x_{\nu}$$

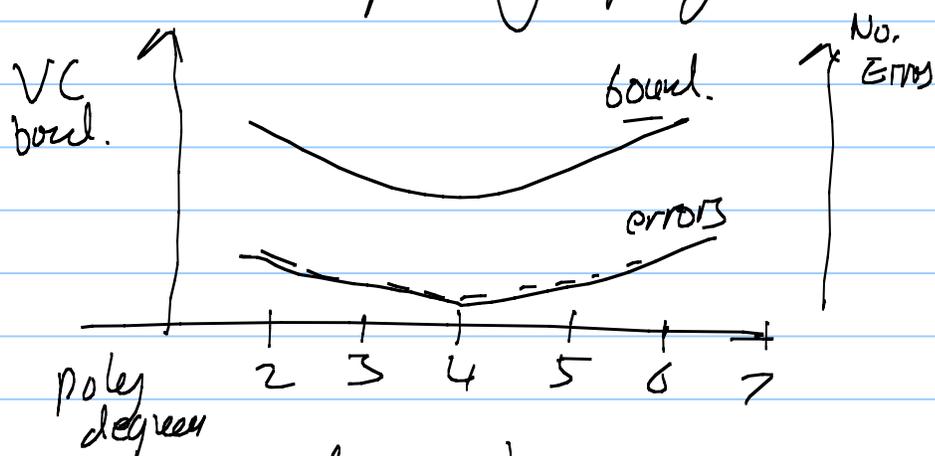
depends on dot product. Hence kernel trick applies.

## (5) Applications of VC-bound.

Use the Margin VC to decide which kernels will generalize best for learning to classify handwritten digits (US Post Office dataset).

→ For each kernel, solve the dual problem to estimate  $R$ .

Schölkopf did this comparing polynomial kernels.



The VC bound predicted the form of the

errors (as a function of poly. degree).

But, the bounds very significantly worse than the errors.

## (6) Structural Risk Minimization

Take the generalization into account when estimating the best classifier.

Statistics says, learn a classifier on a training dataset and use cross-validation to evaluate it.

VC theory says, evaluate the bounded  $\phi(h, n, S)$ . and ensure there are sufficient number of samples to make  $\phi(h, n, S)$  very small.

Alternative - Structural Risk Minimization.

Divide the set of classifiers into a hierarchy of sets  $S_p \subset S_{p+1} \subset \dots$  with VC dimension  $h_p, h_{p+1}, \dots$

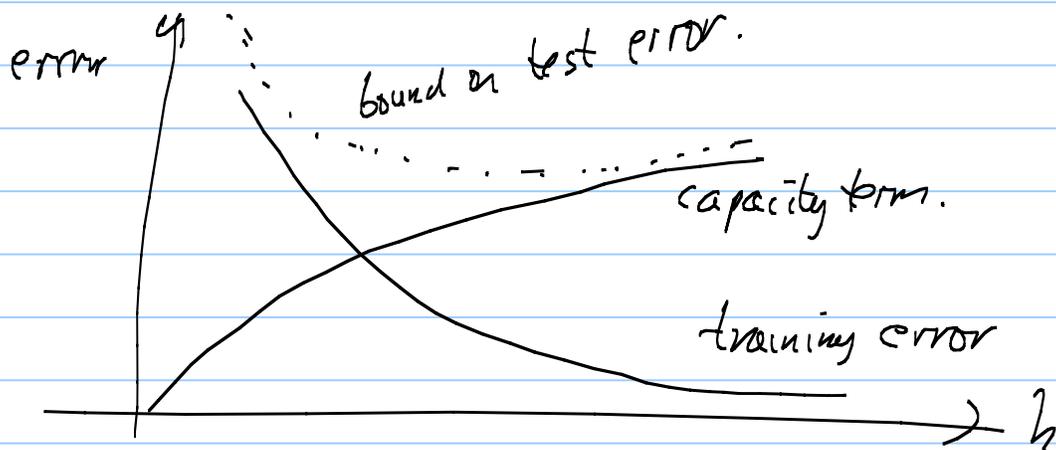
Select classifier to minimize

$$\underbrace{R_{\text{emp}}[\alpha]}_{\text{empirical risk}} + \underbrace{\left(\frac{1}{m}\right) \left( h_p (\log(2m/h_p) + 1) + \log(4/\delta) \right)}_{\text{capacity term}}$$

(7) A small capacity term ensures good generalization (with high probability)

Increasing the amount of data allows you to increase  $\rho$

But, is the bound good enough? Does it give a fair balance between performance on the training set ( $R_{emp}$ ) and generalizability (capacity term)? Does it overweight the importance of generalization?



$$C_{S_{p-1}} \subset C_{S_p} \subset \dots$$