

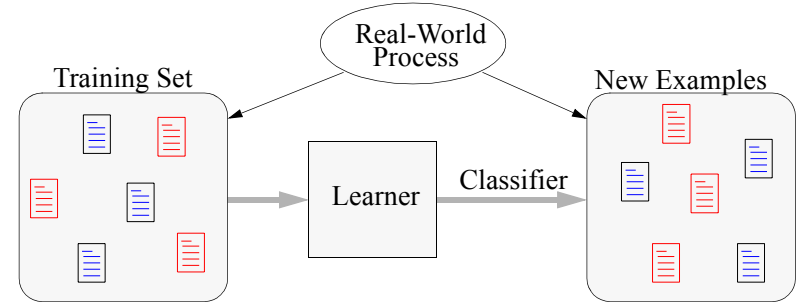
Statistical Learning Theory and PAC-Learning

CS678 Advanced Topics in Machine Learning
Thorsten Joachims
Spring 2003

Outline:

- What is the true (prediction) error of classification rule h ?
- How to bound the true error given the training error?
- Finite hypothesis space and zero training error
- Finite hypothesis space and non-zero training error
- Infinite hypothesis spaces: VC-Dimension and Growth Function

Learning Classifiers



Goal:

- Learner uses training set to find classifier with low prediction error.

Learning Classifiers from Examples (Scenario)

Scenario:

- Generator: Generates descriptions \vec{x} according to distribution $P(\vec{x})$.
- Teacher: Assigns a value y to each description \vec{x} based on distribution $P(y|\vec{x})$.

Given:

- Training examples $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n) \sim P(\vec{x}, y)$ $\vec{x}_i \in \mathfrak{X}^N$ $y_i \in \{1, -1\}$
- Set H of classification rules h (hypotheses) that map descriptions \vec{x} to values y ($h: \vec{x} \rightarrow y$).

Goal of Learner:

- Classification rule h from H that classifies new examples (again from $P(\vec{x}, y)$) with low error rate!

$$P(h(\vec{x}) \neq y) = \int \Delta(h(\vec{x}) \neq y) dP(\vec{x}, y) = Err_P(h)$$

Principle: Empirical Risk Minimization (ERM)

Learning Principle:

Find the decision rule $h^\circ \in H$ for which the training error is minimal:

$$h^\circ = \operatorname{argmin}_{h \in H} \{Err_S(h)\}$$

Training Error:

$$Err_S(h) = \frac{1}{n} \sum_{i=1}^n \Delta(y_i \neq h(\vec{x}_i))$$

==> Number of misclassifications on training examples.

Central Problem of Statistical Learning Theory:
When does a low training error lead to a low generalization error?

Sources of Variation

Learning Task:

- Generator: Generates descriptions \vec{x} according to distribution $P(\vec{x})$.
- Teacher: Assigns a value y to each description \vec{x} based on $P(y|\vec{x})$.

=> Learning Task: $P(\vec{x}, y) = P(y|\vec{x})P(\vec{x})$

Process:

- Select task $P(\vec{x}, y)$
- Training sample S (depends on $P(\vec{x}, y)$)
- Train learning algorithm A (e.g. randomized search)
- Test sample V (depends on $P(\vec{x}, y)$)
- Apply classification rule h (e.g. randomized prediction)

What is the true error of classification rule h?

Includes variation from different test sets.

Problem Setting:

- given rule h
- given (independent) test sample $S = (\vec{x}_1, y_1), \dots, (\vec{x}_k, y_k)$ of size k estimate

$$P(h(\vec{x}) \neq y) = \int \Delta(h(\vec{x}) \neq y) dP(\vec{x}, y) = Err_P(h)$$

Approach: measure error of h on test set

$$Err_V(h) = \frac{1}{k} \sum_{i=1}^k \Delta(y_i \neq h(\vec{x}_i))$$

Binomial Distribution

The probability of observing x heads in a sample of n independent coin tosses, when the probability of heads is p in each toss, is

$$P(X = x|p, n) = \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r}$$

Confidence interval:

Given x observed heads, with at least 95% confidence the true value of p fulfills

$$P(X \geq x|p, n) \geq 0.025 \quad \text{and} \quad P(X \leq x|p, n) \geq 0.025$$

Cross-Validation Estimation

Given:

- training set S of size n

Method:

- partition S into m subsets of equal size
- for i from 1 to m
 - train learner on all subsets except the i' th
 - test learner on i' th subset
 - record error rates on test set

=> Result: average over recorded error rates

Bias of estimate: see leave-one-out

Warning: Test sets are independent, but not the training sets!
=> no strictly valid hypothesis test is known for general learning algorithms (see [Dietterich/97])

Psychic Game

- I guess a 4 bit code
- You all guess a 4 bit code

=> The student who guesses my code clearly has telepathic abilities - right!?

How can You Convince Me of Your Psychic Abilities?

Setting:

- n bits
- $|H|$ players

Question: For which n and $|H|$ is prediction of zero-error player significantly different from random ($p = 0.5$) with probability $1 - \delta$?

=> Hypothesis test for

$$P(h_1 \text{ correct} \vee \dots \vee h_{|H|} \text{ correct}, \text{allnonpsychic}) < \delta$$

or

$$P(\exists h \in H; \text{Err}_s(h) = 0, \forall h \in H; \text{Err}_p(h) = 0.5) < \delta$$

PAC Learning

Definition:

- C = class of concepts $c; X \rightarrow \{1, -1\}$ (functions to be learned)
- H = class of hypotheses $h; X \rightarrow \{1, -1\}$ (functions used by learner A)
- S = training set (of size n)
- ϵ = desired error rate of learned hypothesis
- δ = probability, with which the learner A is allowed to fail

C is PAC-learnable by Algorithm A using H and n examples, if

$$P(\text{Err}(h_{A(S)}) \leq \epsilon) \geq (1 - \delta)$$

for all $c \in C$, ϵ , δ , and $P(X)$ so that A runs in polynomial time dependent on ϵ , δ , the size of the training examples and the size of the concepts.

=> only polynomially many training examples allowed.

Case: Finite H, Zero Error

- The hypothesis space H is finite
- There is always some h with zero training error (A returns one such h)
- Probability that a (single) h with $\text{Err}_p(h) \geq \epsilon$ has training error of zero

$$(1 - \epsilon)^n$$

- Probability that there exists h in H with $\text{Err}_p(h) \geq \epsilon$ that has training error of zero

$$P(\exists h \in H; \text{Err}_s(h) = 0, \text{Err}_p(h) > \epsilon) \leq |H|(1 - \epsilon)^n \leq |H|e^{-\epsilon n}$$

Case: Finite H, Non-Zero Error

Goal:

$$P(|Err_S(h_{A(S)}) - Err_D(h_{A(S)})| \leq \epsilon) \geq (1 - \delta)$$

<=

$$P(\sup_H |Err_S(h_i) - Err_D(h_i)| \leq \epsilon) \geq (1 - \delta)$$

- Probability that for a fixed h, training error and test error differ by more than ϵ (Hoeffding / Chernoff Bound)

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n x_i - p\right| > \epsilon\right) \leq 2e^{-2n\epsilon^2}$$

- Probability over all h in H: union bound => multiply by |H|

Case: Infinite H

- union bound does no longer work.
- maybe not all hypotheses are really different?!

How Many Dichotomies for Fixed Sample?

- Sample S of size n
- Hypothesis class H

$$\Pi_H(S) = \{(h(\vec{x}_1), h(\vec{x}_2), \dots, h(\vec{x}_n)); h \in H\}$$

Definition: H shatters S, if $|\Pi_H(S)| = 2^n$ (i.e. hypotheses from H can split S in all possible ways).

Vapnik/Chervonenkis Dimension

Definition: The VC-dimension of H is equal to the maximal number d of examples that can be split into two sets in all 2^d ways using functions from H (shattering).

	x_1	x_2	x_3	...	x_d
h_1	+	+	+	...	+
h_2	-	+	+	...	+
h_3	+	-	+	...	+
h_4	-	-	+	...	+
...
h_N	-	-	-	...	-

Growth function $\Phi_d(S)$: For all S

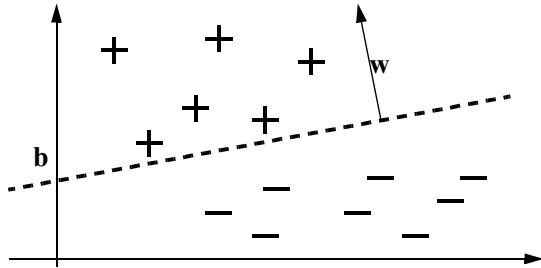
$$|\Pi_H(S)| \leq \Phi_{VCdim(H)}(n) \leq \frac{en}{VCdim(H)} \Bigg\}^{VCdim(H)}$$

Linear Classifiers

Rules of the Form: weight vector \vec{w} , threshold b

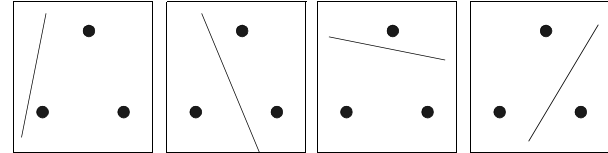
$$h(\vec{x}) = \text{sign} \left[\sum_{i=1}^N w_i x_i + b \right] = \begin{cases} 1 & \text{if } \sum_{i=1}^N w_i x_i + b > 0 \\ -1 & \text{else} \end{cases}$$

Geometric Interpretation (Hyperplane):

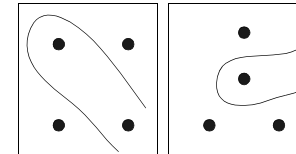


VC-Dimension of Hyperplanes in \mathcal{R}^2

- Three points in \mathcal{R}^2 can be shattered with hyperplanes.



- Four points cannot be shattered.



\Rightarrow Hyperplanes in $\mathcal{R}^2 \rightarrow VCdim=3$

General: Hyperplanes in $\mathcal{R}^N \rightarrow VCdim=N+1$

Error Bound

Question: After n training examples, how close is the training error to the true error?

With probability η it holds for all $h \in H$:

$$Err_P(h) - Err_S(h) \leq \Phi(d, n, \eta)$$

$$\Phi(d, n) = \sqrt{\frac{d \left(\ln \frac{2n}{d} + 1 \right) - \ln \frac{\eta}{4}}{n}}$$

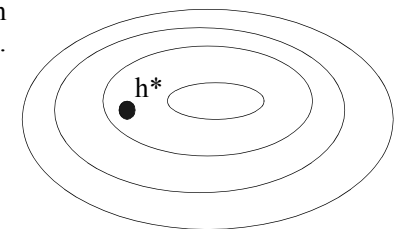
- n number of training examples
- d VC-dimension of hypothesis space H

$$\Rightarrow \boxed{Err_P(h) \leq Err_S(h) + \Phi(d, n, \eta)}$$

SVM Motivation: Structural Risk Minimization

$$\boxed{Err_P(h_i) \leq Err_S(h_i) + \Phi(VCdim(H), n, \eta)}$$

Idea: Structure on hypothesis space.



Goal: Minimize upper bound on true error rate.

