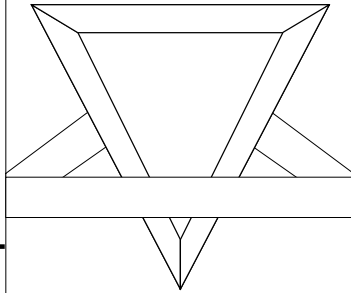
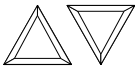


Chapter 4 Statistical Learning Theory



Outline

- ▼ 4.1 Conditions for Consistency and Convergence of ERM
- ▼ 4.2 Growth Function and VC-Dimension
- ▼ 4.3 Bounds on the Generalization
- ▼ 4.4 Structural Risk Minimization
- ▼ 4.5 Case Study: Comparison of Methods for Model Selection
- ▼ 4.6 Summary



Introduction

- ▼ VC (Vapnik-Chrvoenkis) theory: SLT
 - Conditions for consistency of the ERM inductive principle
 - Bounds on the generalization ability of learning machines based on these conditions
 - Principles for inductive inference from small samples based on these bounds
 - Constructive methods for implementing above inductive principles



4.1 Conditions for Consistency and Convergence of ERM

$$R(\omega) = \int Q(z, \omega) dF(z) \text{ or } R(\omega) = \int Q(z, \omega) p(z) dz$$

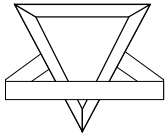
$$Q(z, \omega) = (y - f(x, \omega))^2$$

$$R_{\text{emp}}(\omega) = \sum_{i=1}^n Q(z_i, \omega)$$

solution approaches to the learning problem

- estimate unknown c.d.f from data -> find optimal estimate $f(x, \omega_0)$
- seek an estimate providing minimum of the (known) empirical risk (ERM)





Conditions for Consistency and Convergence of ERM

▼ Consistency

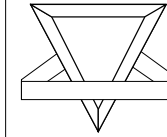
- Requirement that estimates provided by ERM should converge to the true values as the # of samples grows

$$R(\omega^* | n) \rightarrow R(\omega_0) \quad \text{when } n \rightarrow \infty$$

$$R_{emp}(\omega^* | n) \rightarrow R(\omega_0) \quad \text{when } n \rightarrow \infty$$

- We can expect

$$R_{emp}(\omega^* | n) < R(\omega^* | n)$$



Conditions for Consistency and Convergence of ERM

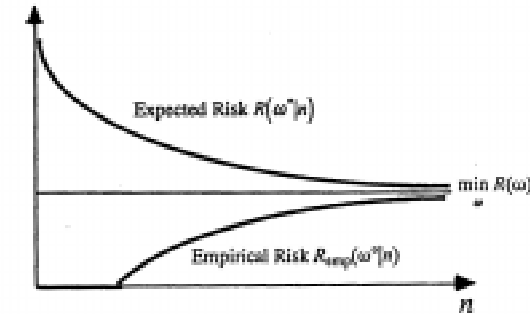
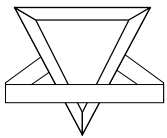


Figure 4.1 Consistency of the ERM.

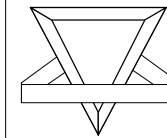


Conditions for Consistency and Convergence of ERM

▼ Nontrivial consistency (Vapnik)

- Consistency should hold for all approximating functions
- Key Theorem of Learning Theory (Vapnik and Chervonenkis, 1989)

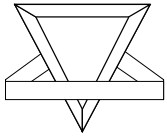
$$\lim_{n \rightarrow \infty} P \left[\sup_{\omega} |R(\omega) - R_{emp}(\omega)| > \varepsilon \right] = 0, \quad \forall \varepsilon > 0$$



Conditions for Consistency and Convergence of ERM

- Diversity of a set of functions w.r.t. Z_n (dichotomy case)
 - $N(Z_n)$: # of different dichotomies by $Q(z, \omega)$
- $H(Z_n) = \ln N(Z_n)$: random entropy
- $H(n) = E(\ln N(Z_n))$: VC entropy of the set of indicator functions on a sample of size n from $F(z)$
 - provides a measure of the expected diversity of a set of indicator functions with respect to a sample of a given size
 - depends on the set of indicator funcs and on the (unknown) distribution of samples $F(z)$
- $G(n) = \ln \max_{Z_n} N(Z_n)$
 - growth function: distribution-independent
 - provides an upper bound for the (distribution-dependent) entropy





Conditions for Consistency and Convergence of ERM

$$G(n) \leq n \ln 2$$

- $H_{\text{ann}}(n) = \ln E(N(\mathbf{Z}_n))$: annealed VC entropy
Using Jensen's inequality

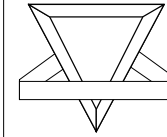
$$\sum_i a_i \ln x_i \leq \ln \left(\sum_i a_i x_i \right)$$

It can be shown that

$$H(n) \leq H_{\text{ann}}(n)$$

- $H(n) \leq H_{\text{ann}}(n) \leq G(n) \leq n \ln 2$
- Necessary and sufficient condition for consistency of the ERM principle (Vapnik and Chervonenkis, 1968)

$$\lim_{n \rightarrow \infty} \frac{H(n)}{n} = 0$$

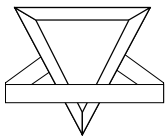


Conditions for Consistency and Convergence of ERM

- asymptotic rate of convergence is called *fast* if for any $n > n_0$ the following holds

$$P(R(\bar{\omega}) - R(\omega^*) < \varepsilon) = e^{-cn\varepsilon^2}$$

($c > 0$ is a positive constant)



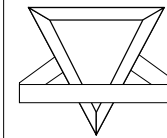
Conditions for Consistency and Convergence of ERM

- Sufficient condition for fast rate of convergence

$$\lim_{n \rightarrow \infty} \frac{H_{\text{ann}}(n)}{n} = 0 \text{ (distribution-dependent condition)}$$

- Distribution-independent condition for consistency of ERM and fast convergence

$$\lim_{n \rightarrow \infty} \frac{G(n)}{n} = 0$$



4.2 Growth Function and VC-Dimension

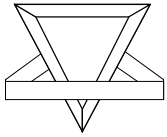
▼ VC-dimension

- Growth function is either linear or bounded by a logarithmic function of # of samples n .
- The point $n = h$ where the growth starts to slow down is called the VC-dimension.
- If it is finite, the growth function is bounded by

$$G(n) \leq h \left(1 + \ln \frac{n}{h} \right)$$

- Finiteness of h provides necessary and sufficient conditions for the fast rate convergence and for distribution-independent consistency of ERM.





VC-Dimension

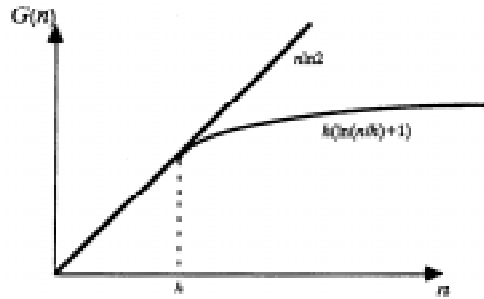
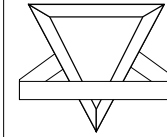
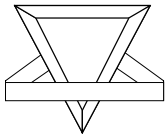


Figure 4.2 Behavior of the Growth Function.



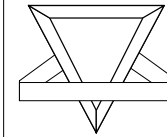
VC-Dimension

- ▼ Criterion for demarcation between true and false (inductive) theories (Popper)
 - The necessary condition for the inductive theory to be true is the feasibility of its falsification, i.e., the existence of certain assertions (facts) that cannot be explained by the theory.
 - e.g.) VC-dimension is infinite -> false model
 - *shattering*: n samples can be separated by a set of indicator func.s in all 2^n .



VC-Dimension

- VC-dimension of a set of indicator functions
 - VC-dimension $h \leftrightarrow$ if there exist h samples that can be shattered by this set of functions but not $h+1$
- ▼ VC-dimension of the set of real-valued func.
 - Indicator function
 - $A \leq Q(\mathbf{z}, \omega) \leq B$
 - $I(\mathbf{z}, \omega, \beta) = I[Q(\mathbf{z}, \omega) - \beta > 0]$
 - VC-dimension of real function Q is equal to the VC-dimension of indicator function.



VC-Dimension

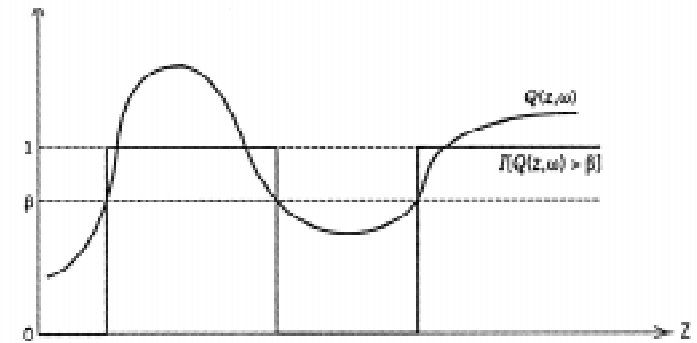
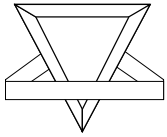


Figure 4.3 VC-dimension of the set of real-valued functions.





VC-Dimension

- ▼ VC-dimension for Classification and regression problems

Classification

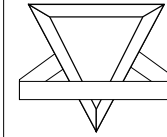
$$Q(\mathbf{z}, \omega) = \begin{cases} f(\mathbf{x}, \omega) & \text{if } y = 0 \\ 1 - f(\mathbf{x}, \omega) & \text{if } y = 1 \end{cases}$$

Regression

$$Q(\mathbf{z}, \omega) = (y - f(\mathbf{x}, \omega))^2$$

$$h_f \leq h \leq ch_f$$

$$h \approx h_f$$



VC-Dimension

- ▼ Examples of calculating VC-dimension

- VC-Dimension of a set of linear indicator func

- $h = d + 1$

$$Q(\mathbf{z}, \mathbf{w}) = I\left(\sum_{i=1}^d w_i z_i + w_0 > 0\right)$$

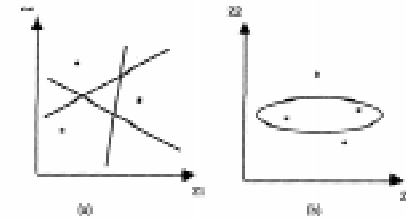
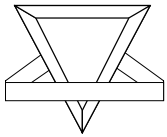


Figure 4.4 VC-dimension of linear indicator functions. (a) Linear functions can separate any three points in a two-dimensional space. (b) Linear functions cannot split five points into two classes as shown.



VC-Dimension

- Set of univariate functions with a single parameter

- $f(x, \omega) = I(\sin wx > 0)$

- $h = \infty$

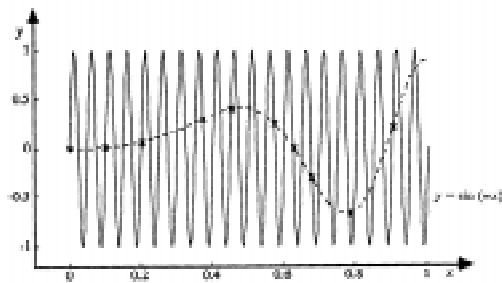
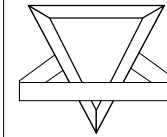


Figure 4.5 Set of indicator functions with infinite VC-dimension.



VC-Dimension

- Set of rectangular indicator functions

$$Q(\mathbf{z}, \mathbf{c}, \mathbf{w}) = 1 \text{ if and only if } |z_i - c_i| \leq w_i, i = 1, 2, \dots, d$$

- $h = 2d$

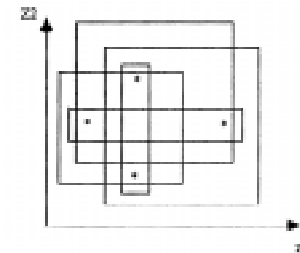
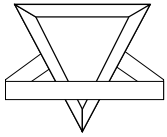


Figure 4.6 VC-dimension of a set of rectangular functions.





VC-Dimension

- ▼ Set of radially symmetric indicator functions

$Q(\mathbf{z}, \mathbf{c}, r) = 1$ if and only if $\|\mathbf{z} - \mathbf{c}\| \leq r$

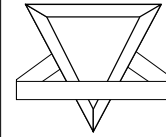
- $h = d+1$

- ▼ Set of real-valued “local” functions

$$I(\mathbf{x}, \mathbf{c}, \alpha) = K\left(\frac{\|\mathbf{x} - \mathbf{c}\|}{\alpha}\right)$$

$$I(\mathbf{x}, \mathbf{c}, \alpha, \beta) = I\left[K\left(\frac{\|\mathbf{x} - \mathbf{c}\|}{\alpha}\right) - \beta\right] \quad (d+2 \text{ free parameter})$$

- $h = d+1$



VC-Dimension

- Linear combination of fixed basis functions

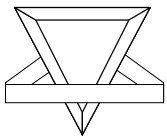
$$Q_m(\mathbf{z}, \mathbf{w}) = \sum_{i=1}^m w_i \mathbf{g}_i(\mathbf{z}) + w_0$$

- $h = m+1$

- Linear combination of adaptive basis functions nonlinear in parameters

$$Q_m(\mathbf{z}, \mathbf{w}, \mathbf{v}) = \sum_{i=1}^m w_i \mathbf{g}_i(\mathbf{z}, \mathbf{v}) + w_0$$

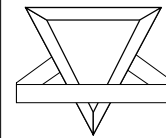
- can be infinite



4.3. Bounds on the Generalization

- ▼ The upper bounds on the rate of uniform convergence of the learning processes

- Evaluate the difference between true risk and the known empirical risk.
- Constructive distribution-independent bounds form the foundation for a new inductive principle (structural risk minimization) and associated constructive procedures.



Bounds on the Generalization

- ▼ Classification

- With probability $1-\eta$ for all Q

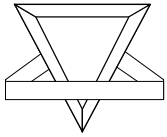
$$R(\omega) \leq R_{emp}(\omega) + \frac{\varepsilon}{2} \left(1 + \sqrt{1 + \frac{4R_{emp}(\omega)}{\varepsilon}}\right)$$

$$\text{where, } \varepsilon = \varepsilon\left(\frac{n}{h}, \frac{-\ln \eta}{n}\right) = a \cdot \frac{h[\ln(a_2 n / h) + 1] - \ln(\eta / 4)}{n}$$

- If the set of loss functions contains a finite number of elements N

$$\varepsilon = 2 \frac{\ln N - \ln \eta}{n}$$



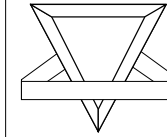


Bounds on the Generalization

- With $1-2\eta$ probability for the func. that minimizes empirical risk

$$R(\omega^* | n) - \min_{\omega} R(\omega) \leq \sqrt{\frac{-\ln \eta}{2n}} + \frac{\epsilon}{2} \left(1 + \sqrt{1 + \frac{4}{\epsilon}}\right)$$

- Confidence level : $1 - \eta$
- There is a trade-off between the accuracy provided by the bounds and the degree of confidence



Bounds on the Generalization

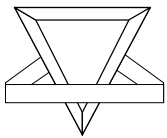
▼ Regression

- For $1 - \eta$ probability, c depending on the “tail of the distribution” of the loss functions

$$R(\omega) \leq \frac{\text{Re } mp(\omega)}{(1 - c\sqrt{\epsilon})^+}$$

- $1 - 2\eta$ for that function that minimizes empirical risk

$$\frac{R(\omega^* | n) - \min_{\omega} R(\omega)}{\min_{\omega} R(\omega)} \leq \frac{c\sqrt{\epsilon}}{1 - c\sqrt{\epsilon}} + O(1/n)$$

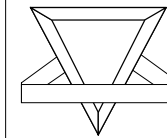


Bounds on the Generalization

- If n and η are held at particular values, it is possible to determine the value of h that leads to the bound approaching

$$\epsilon(h) = a_1 \frac{h[\ln(a_2 n / h) + 1 - \ln(\eta/4)]}{n} \geq 1 \quad \text{with } a_1 = 1, a_2 = 1$$

$$\frac{h}{n} \leq 0.8 \quad \text{for } \eta \geq \min\left(\frac{4}{\sqrt{n}}, 1\right)$$



Bounds on the Generalization

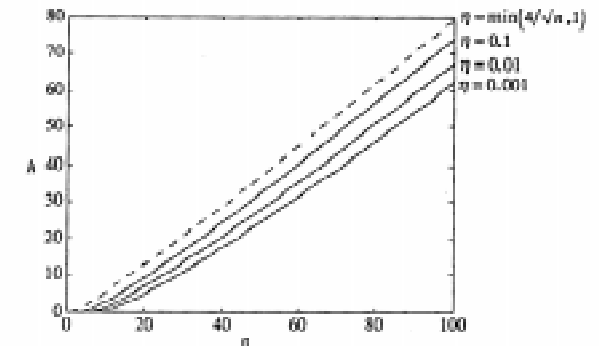
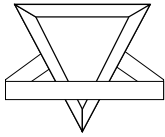


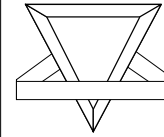
Figure 4.7 Values of n , η , and h that cause the generalization bound to approach infinity under real-life conditions ($a_1 = 1$, $a_2 = 1$).





4.4 Structural Risk Minimization

- If $\frac{h}{n}$ is small, other factors must be minimized.
- First term in (4.22) depends on a particular func. from the set of functions.
- Second term depends mainly on the VC-dimension of the set of functions.
- Structural risk minimization (SRM) provides a formal mechanism for choosing an optimal model complexity for a finite sample.



SRM

- Under SRM the set S of loss functions $Q(\mathbf{z}, \mathbf{w})$, $\mathbf{w} \in \Omega$ has a *structure*, that is, it consists of the nested subsets $S_k = \{Q(\mathbf{z}, \mathbf{w}), \mathbf{w} \in \Omega_k\}$ such that $S_1 \subset S_2 \subset \dots \subset S_k \subset \dots$ where $h_1 < h_2 \dots < h_k \dots$
- Solving a learning problem with finite data
 - requires a priori specification of a structure on a set of approximating functions then
 - 1. selecting an element of a structure (having optimal complexity)
 - 2. estimating the model from the element

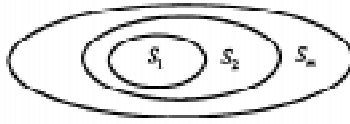
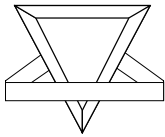


Figure 4.8 Structure on a set of functions.

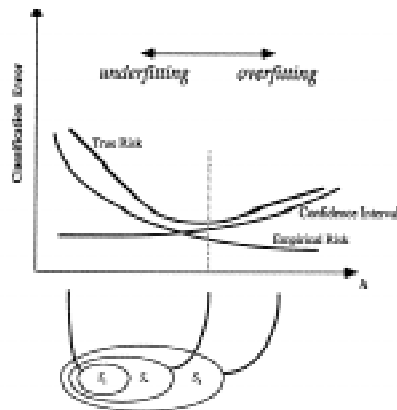
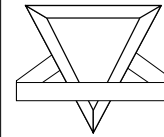


Figure 5.8 An upper bound on the true (expected) risk and the empirical risk as a function of model fixed h .



SRM

- 1. Dictionary representation

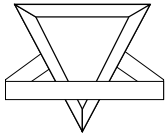
$$f_m(\mathbf{x}, \mathbf{w}, \mathbf{V}) = \sum_{i=0}^m w_i \mathbf{g}(\mathbf{x}, \mathbf{v}_i)$$

$$f_1 \subset f_2 \subset \dots \subset f_k \dots$$

for example,

$$f_m(x, \mathbf{w}) = \sum_{i=0}^m w_i x^i$$





SRM

- 2. Penalization formulation

$$S_k = \{f(x, w), \|w\|^2 \leq c_k\} \text{ where } c_1 < c_2 < c_3 \dots$$

$$R_{pen}(\omega, \lambda_k) = R_{emp}(\omega) + \lambda_k \|w\|^2$$

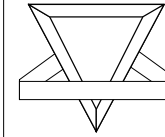
with an appropriately chosen Lagrange multiplier λ_k

such that $\lambda_1 > \lambda_2 > \lambda_3 \dots$

- 3. Input preprocessing

$$z = K(x, \beta)$$

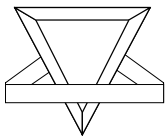
$$S_k = \{f(K(x, \beta), w), \beta \leq c_k\} \text{ where } c_1 > c_2 > c_3 \dots$$



SRM

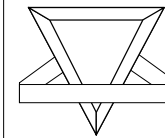
- 4. Initial conditions for training algorithm

$$S_k = \{A : f(x, w), \|w^0\| \leq c_k\}$$



4.5 Case Study : Comparison of Methods for Model Selection

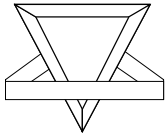
- Objective : to choose the model complexity optimally for a given training sample
- Practical application of either SRM or penalization requires two tasks :
 - Estimation of model parameters (via minimization of the penalized empirical risk)
 - Estimation of the prediction risk
- Two major approaches for estimating prediction risk :
 - Analytic methods
 - Resampling or data-driven methods



Case Study

- Training samples
 - $y = \sin^2(2\pi x) + \varepsilon, x \in [0,1]$
 - sample size : 10, 20, 30, 100
 - noise : defined in terms of SNR as the ratio of the standard deviation of the true output values for given input samples over the standard deviation of the gaussian noise
- Approximating functions
 - class of polynomials of degree m
 - set of functions are linear in parameters : solving linear least squares



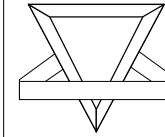


Case Study

• Model selection

- Choosing an optimal polynomial degree m for a given training sample
- Comparison set
 - Final prediction error (fpe)
 - Schwartz criteria (sc)
 - Generalized cross-validation (gcv)
 - Shibata's model selector (sms)
 - Leave-one-out cross-validation (cv)
 - Vapnik's measure (vm)

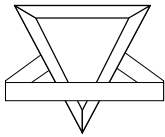
$$g(p, n) = \left(1 - \sqrt{p - p \ln p + \frac{\ln n}{2n}} \right)^{-1}$$



Case Study

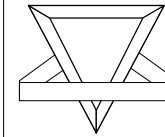
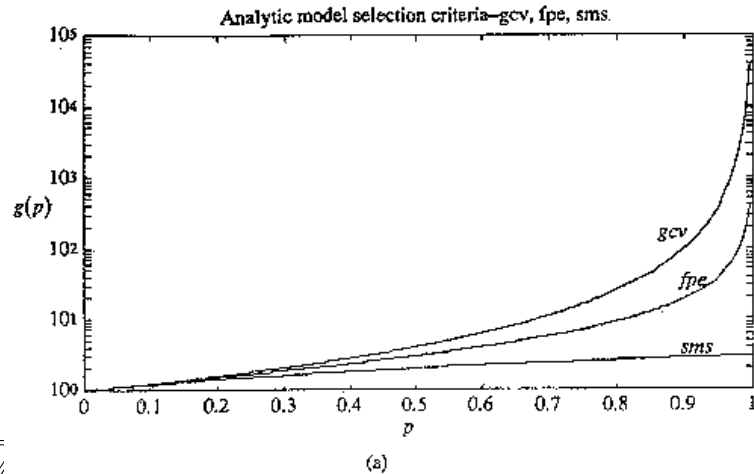
• Comparison strategy

- 1000 times repetition for given small-size training set
- Standard box plot notation describing empirical distribution
- Experimental results
 - Vapnik's measure for model selection shows superior overall performance
- Summary of comparison results
 - Small size training data may cause no guaranteed performance, so that measures like Vapnik's are required.
 - Vapnik's measure guarantees the best worst-case estimates.



Case Study

CASE STUDY: COMPARISON OF METHODS FOR MODEL SELECTION 123



Case Study

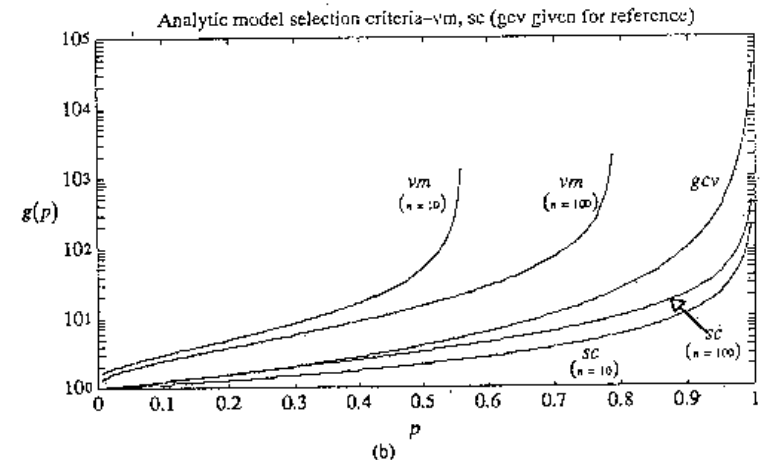
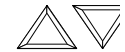
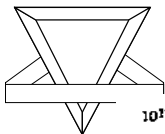
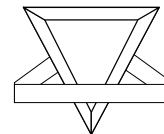
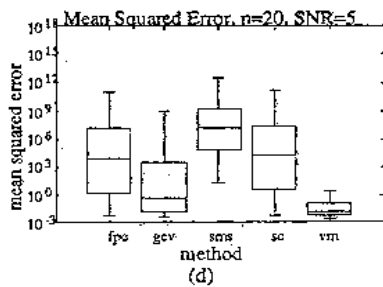
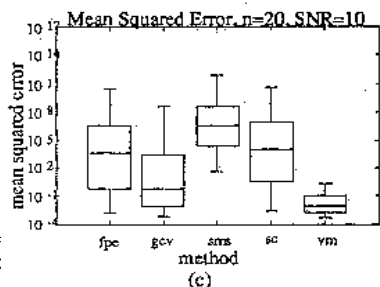
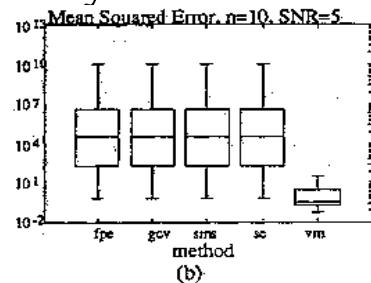
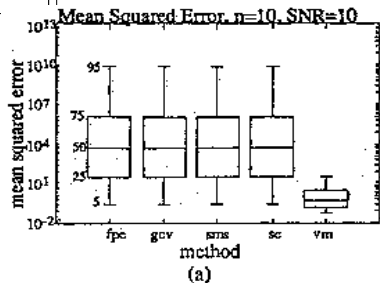


Figure 4.10 Various analytical model selection penalization functions: (a) Generalized cross-validation (gcv), final prediction error (fpe), and Shibata's model selector (sms). (b) Vapnik's measure (vm) and Schwartz criteria (sc) for sample sizes indicated. The parameter p is equal to h/n .





Case Study



Case Study

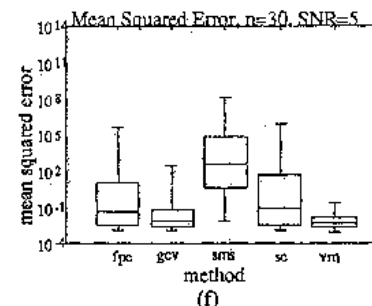
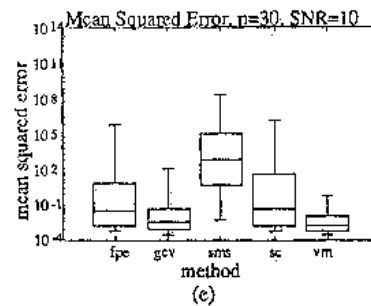
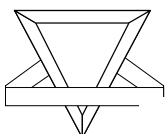
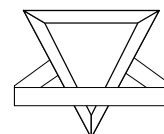
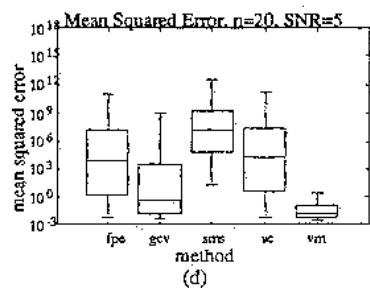
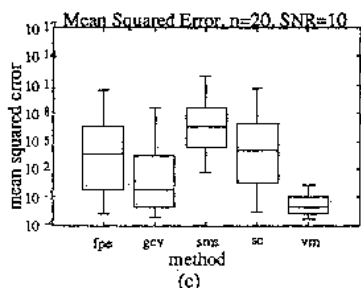
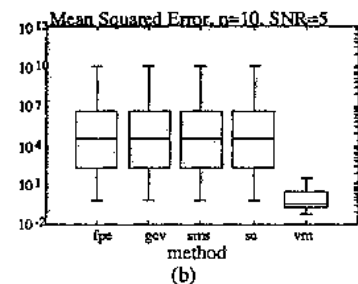
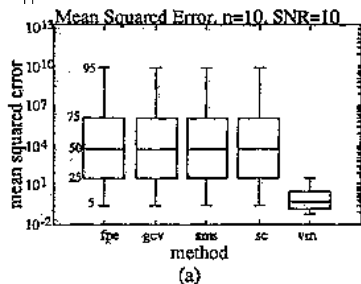


Figure 4.11 Mean squared error for each model selection approach. Training set size and noise level are indicated.



Case Study



Case Study

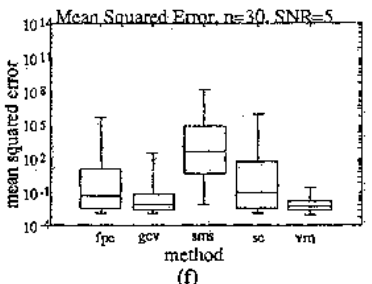
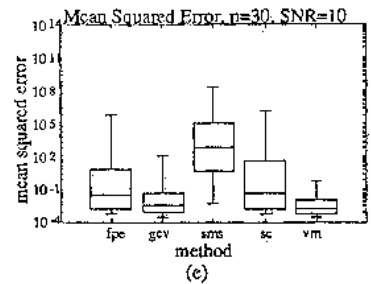
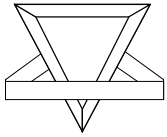
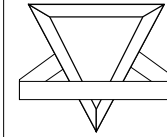
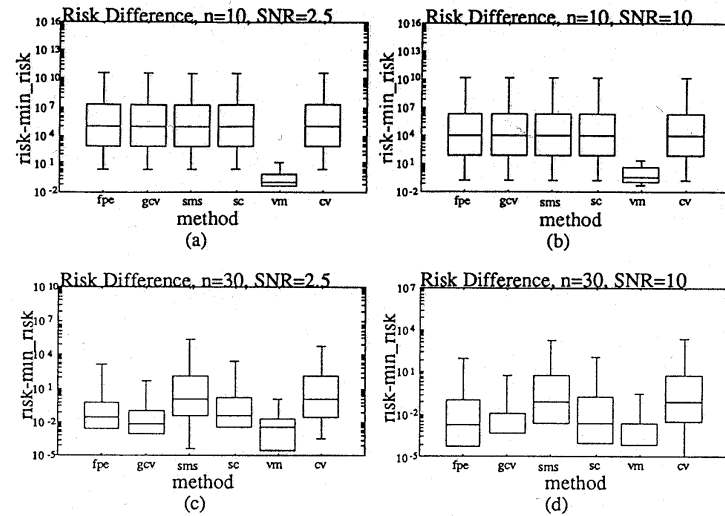


Figure 4.11 Mean squared error for each model selection approach. Training set size and noise level are indicated.





Case Study



Case Study

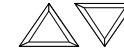
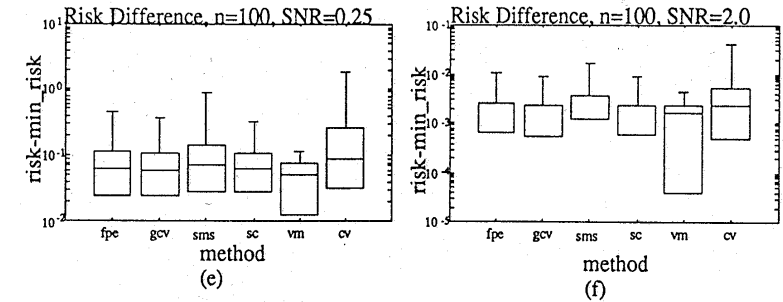
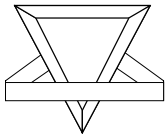


Figure 4.13 Risk difference for each model selection approach. Training set size and noise level are indicated.



4.6 Summary

- SLT framework can be used in three ways:
 - For the interpretation and critical evaluation of empirical learning methods developed in statistics and neural networks.
 - For developing new constructive learning procedures based on SLT.
 - For developing new inductive principles, such as transductive inference and local risk minimization.
- Comments on SLT framework
 - SLT sometimes doesn't seem to conform with real and complex problems and we cannot expect SLT to provide immediate and clear solutions to practical problems. With all these difficulties, all learning methods must be consistent in SLT senses in order to be a reliable one.

