

*Lecture Slides for*



INTRODUCTION TO

# *Machine Learning*

ETHEM ALPAYDIN

© The MIT Press, 2004

*alpaydin@boun.edu.tr*

*<http://www.cmpe.boun.edu.tr/~ethem/i2ml>*



CHAPTER 14:

*Assessing and Comparing  
Classification Algorithms*

# Introduction

- Questions:
  - Assessment of the expected error of a learning algorithm: Is the error rate of 1-NN less than 2%?
  - Comparing the expected errors of two algorithms: Is  $k$ -NN more accurate than MLP ?
- Training/validation/test sets
- Resampling methods:  $K$ -fold cross-validation

# Algorithm Preference

- Criteria (Application-dependent):
  - Misclassification error, or risk (loss functions)
  - Training time/space complexity
  - Testing time/space complexity
  - Interpretability
  - Easy programmability
- Cost-sensitive learning

# Resampling and K-Fold Cross-Validation

- The need for multiple training/validation sets  
 $\{X_i, V_i\}_i$ : Training/validation sets of fold  $i$
- $K$ -fold cross-validation: Divide  $X$  into  $k$ ,  $X_i, i=1, \dots, K$

$$V_1 = X_1 \quad T_1 = X_2 \cup X_3 \cup \boxtimes \cup X_K$$

$$V_2 = X_2 \quad T_2 = X_1 \cup X_3 \cup \boxtimes \cup X_K$$

$\boxtimes$

$$V_K = X_K \quad T_K = X_1 \cup X_2 \cup \boxtimes \cup X_{K-1}$$

- $T_i$  share  $K-2$  parts

# 5×2 Cross-Validation

- 5 times 2 fold cross-validation (Dietterich, 1998)

$$T_1 = X_1^{(1)} \quad V_1 = X_1^{(2)}$$

$$T_2 = X_1^{(2)} \quad V_2 = X_1^{(1)}$$

$$T_3 = X_2^{(1)} \quad V_3 = X_2^{(2)}$$

$$T_4 = X_2^{(2)} \quad V_4 = X_2^{(1)}$$

⊠

$$T_9 = X_5^{(1)} \quad V_9 = X_5^{(2)}$$

$$T_{10} = X_5^{(2)} \quad V_{10} = X_5^{(1)}$$

# Bootstrapping

- Draw instances from a dataset *with replacement*
- Prob that we do not pick an instance after N draws

$$\left(1 - \frac{1}{N}\right)^N \approx e^{-1} = 0.368$$

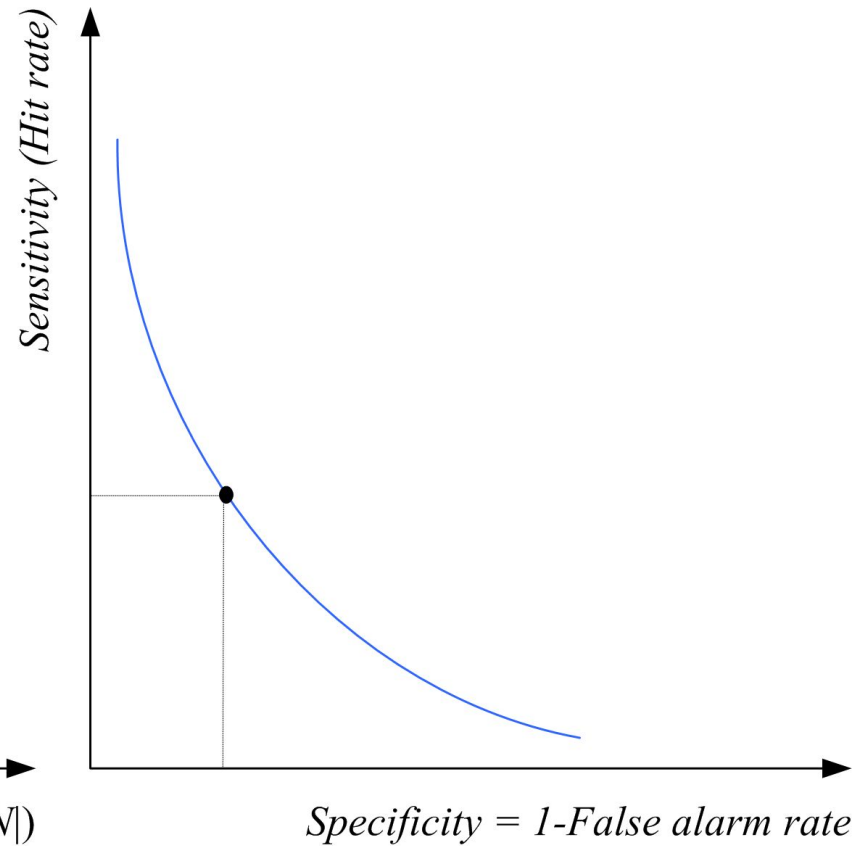
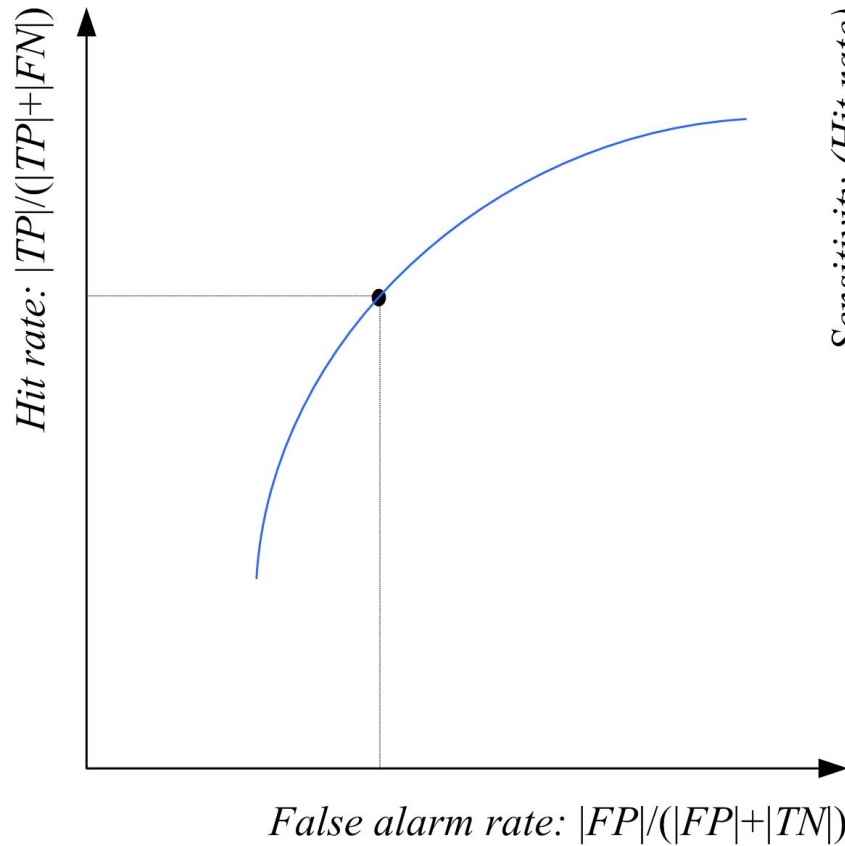
that is, only 36.8% is new!

# Measuring Error

	Predicted class	
True Class	Yes	No
Yes	TP: True Positive	FN: False Negative
No	FP: False Positive	TN: True Negative



# ROC Curve



# Interval Estimation

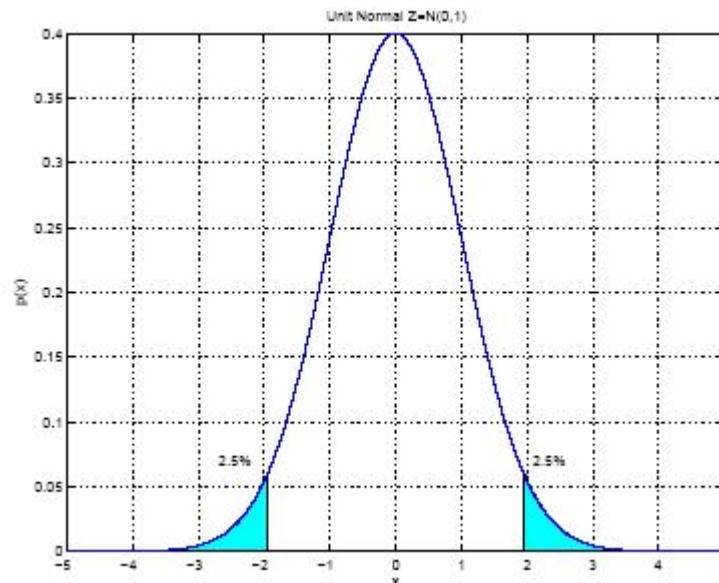
- $X = \{x^t\}_t$  where  $x^t \sim N(\mu, \sigma^2)$
- $m \sim N(\mu, \sigma^2/N)$

$$\sqrt{N} \frac{(m - \mu)}{\sigma} \sim Z$$

$$P\left\{-1.96 < \sqrt{N} \frac{(m - \mu)}{\sigma} < 1.96\right\} = 0.95$$

$$P\left\{m - 1.96 \frac{\sigma}{\sqrt{N}} < \mu < m + 1.96 \frac{\sigma}{\sqrt{N}}\right\} = 0.95$$

$$P\left\{m - z_{\alpha/2} \frac{\sigma}{\sqrt{N}} < \mu < m + z_{\alpha/2} \frac{\sigma}{\sqrt{N}}\right\} = 1 - \alpha$$



100(1-  $\alpha$ ) percent  
confidence  
interval

$$P\left\{\sqrt{N} \frac{(m - \mu)}{\sigma} < 1.64\right\} = 0.95$$

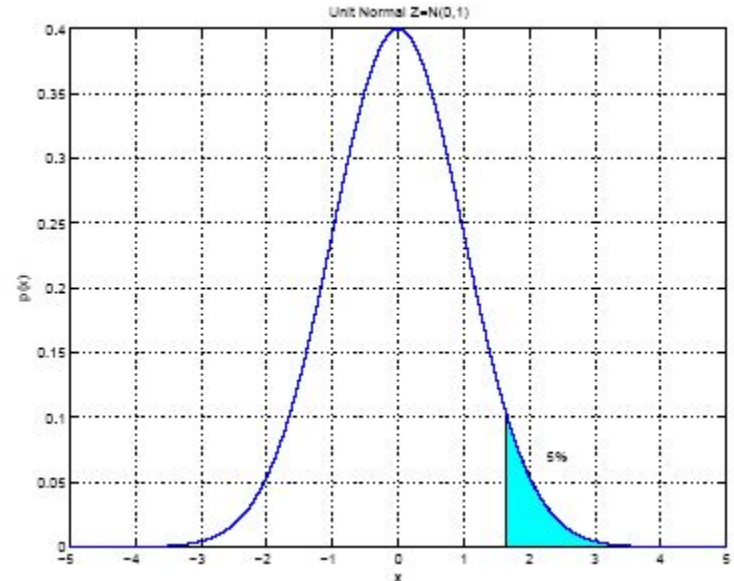
$$P\left\{m - 1.64 \frac{\sigma}{\sqrt{N}} < \mu\right\} = 0.95$$

$$P\left\{m - z_{\alpha} \frac{\sigma}{\sqrt{N}} < \mu\right\} = 1 - \alpha$$

When  $\sigma^2$  is not known:

$$S^2 = \sum_t (x^t - m)^2 / (N - 1) \quad \frac{\sqrt{N} (m - \mu)}{S} \sim t_{N-1}$$

$$P\left\{m - t_{\alpha/2, N-1} \frac{S}{\sqrt{N}} < \mu < m + t_{\alpha/2, N-1} \frac{S}{\sqrt{N}}\right\} = 1 - \alpha$$



# Hypothesis Testing

- Reject a **null hypothesis** if not supported by the sample with enough confidence
- $X = \{x^t\}_t$  where  $x^t \sim N(\mu, \sigma^2)$

$$H_0: \mu = \mu_0 \text{ vs. } H_1: \mu \neq \mu_0$$

Accept  $H_0$  with **level of significance**  $\alpha$  if  $\mu_0$  is in the  $100(1-\alpha)$  confidence interval

$$\frac{\sqrt{N}(m - \mu_0)}{\sigma} \in (-z_{\alpha/2}, z_{\alpha/2})$$

Two-sided test

	Decision	
Truth	Accept	Reject
True	Correct	Type I error
False	Type II error	Correct (Power)

- One-sided test:  $H_0: \mu \leq \mu_0$  vs.  $H_1: \mu > \mu_0$

Accept if

$$\frac{\sqrt{N}(m - \mu_0)}{\sigma} \in (-\infty, z_\alpha)$$

- Variance unknown: Use  $t$ , instead of  $z$

Accept  $H_0: \mu = \mu_0$  if

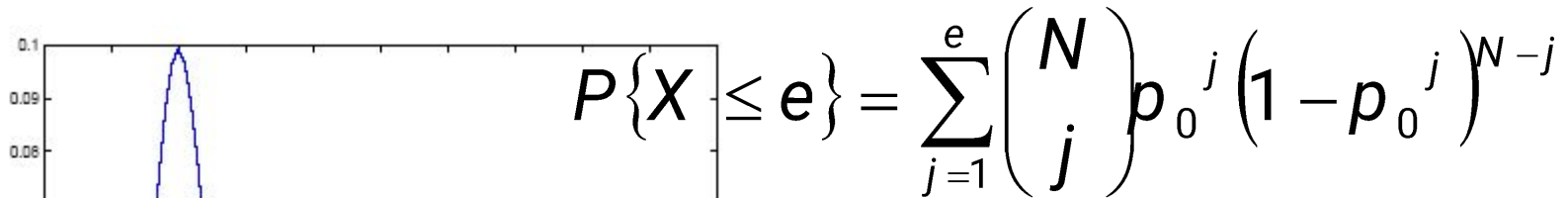
$$\frac{\sqrt{N}(m - \mu_0)}{S} \in (-t_{\alpha/2, N-1}, t_{\alpha/2, N-1})$$

# Assessing Error:

$$H_0: p \leq p_0 \text{ vs. } H_1: p > p_0$$

- Single training/validation set: **Binomial Test**

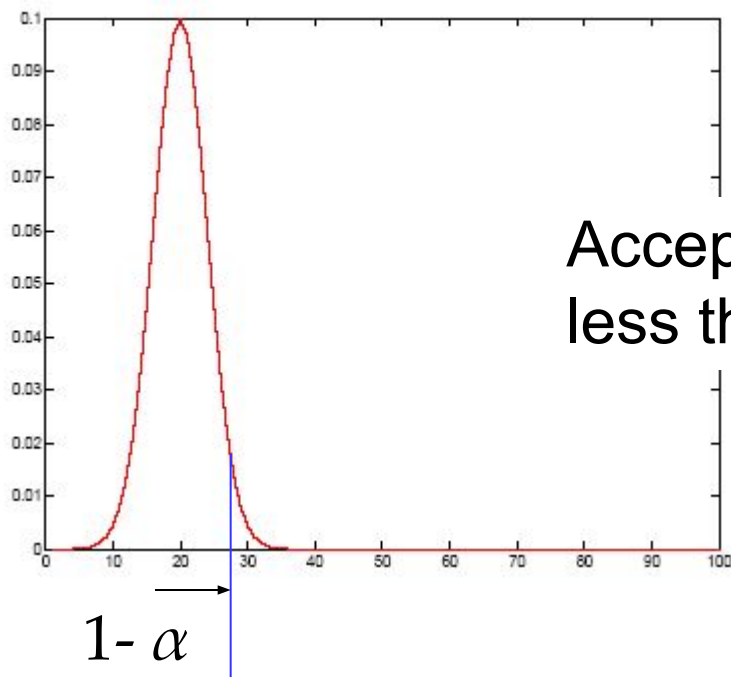
If error prob is  $p_0$ , prob that there are  $e$  errors or less in  $N$  validation trials is



Accept if this prob is less than  $1 - \alpha$

# Normal Approximation to the Binomial

- Number of errors  $X$  is approx  $N$  with mean  $Np_0$  and var  $Np_0(1-p_0)$



$$\frac{X - Np_0}{\sqrt{Np_0(1-p_0)}} \sim Z$$

Accept if this prob for  $X = e$  is less than  $z_{1-\alpha}$

# Paired $t$ Test

- Multiple training/validation sets
- $x_i^t = 1$  if instance  $t$  misclassified on fold  $i$

- Error rate of fold  $i$ :

$$p_i = \frac{\sum_{t=1}^N x_i^t}{N}$$

- With  $m$  and  $s^2$  average and var of  $p_i$  we accept  $p_0$  or less error if

$$\frac{\sqrt{K}(m - p_0)}{S} \sim t_{K-1}$$

is less than  $t_{\alpha, K-1}$



# Comparing Classifiers:

$$H_0: \mu_0 = \mu_1 \text{ vs. } H_1: \mu_0 \neq \mu_1$$

- Single training/validation set: McNemar's Test

$e_{00}$ : Number of examples misclassified by both	$e_{01}$ : Number of examples misclassified by 1 but not 2
$e_{10}$ : Number of examples misclassified by 2 but not 1	$e_{11}$ : Number of examples correctly classified by both

- Under  $H_0$ , we expect  $e_{01} = e_{10} = (e_{01} + e_{10})/2$

$$\frac{(|e_{01} - e_{10}| - 1)^2}{e_{01} + e_{10}} \sim \chi_1^2$$

Accept if  $< \chi_{\alpha,1}^2$

# *K-Fold CV Paired t Test*

- Use  $K$ -fold cv to get  $K$  training/validation folds
- $p_i^1, p_i^2$ : Errors of classifiers 1 and 2 on fold  $i$
- $p_i = p_i^1 - p_i^2$ : Paired difference on fold  $i$
- The null hypothesis is whether  $p_i$  has mean 0

$$H_0 : \mu = 0 \text{ vs. } H_0 : \mu \neq 0$$

$$m = \frac{\sum_{i=1}^K p_i}{K} \quad s^2 = \frac{\sum_{i=1}^K (p_i - m)^2}{K - 1}$$

$$\frac{\sqrt{K}(m - 0)}{s} = \frac{\sqrt{K} \cdot m}{s} \sim t_{K-1} \text{ Accept if in } (-t_{\alpha/2, K-1}, t_{\alpha/2, K-1})$$

## 5×2 cv Paired t Test

- Use 5×2 cv to get 2 folds of 5 tra/val replications (Dietterich, 1998)
- $p_i^{(j)}$  : difference btw errors of 1 and 2 on fold  $j=1, 2$  of replication  $i=1, \dots, 5$

$$\bar{p}_i = (p_i^{(1)} + p_i^{(2)}) / 2 \quad s_i^2 = (p_i^{(1)} - \bar{p}_i)^2 + (p_i^{(2)} - \bar{p}_i)^2$$

$$\frac{p_1^{(1)}}{\sqrt{\sum_{i=1}^5 s_i^2 / 5}} \sim t_5$$

Two-sided test: Accept  $H_0: \mu_0 = \mu_1$  if in  $(-t_{\alpha/2,5}, t_{\alpha/2,5})$

One-sided test: Accept  $H_0: \mu_0 \leq \mu_1$  if  $< t_{\alpha,5}$

## 5×2 cv Paired F Test

$$\frac{\sum_{i=1}^5 \sum_{j=1}^2 (p_i^{(j)})^2}{2 \sum_{i=1}^5 s_i^2} \sim F_{10,5}$$

Two-sided test: Accept  $H_0: \mu_0 = \mu_1$  if  $< F_{\alpha,10,5}$

# Comparing $L > 2$ Algorithms: Analysis of Variance (Anova)

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_L$$

- Errors of  $L$  algorithms on  $K$  folds

$$X_{ij} \sim N(\mu_j, \sigma^2), j = 1, \dots, L, i = 1, \dots, K$$

- We construct two estimators to  $\sigma^2$ .

One is valid if  $H_0$  is true, the other is always valid.

We reject  $H_0$  if the two estimators disagree.

If  $H_0$  is true :

$$m_j = \sum_{i=1}^K \frac{X_{ij}}{K} \sim N(\mu, \sigma^2 / K)$$

$$m = \frac{\sum_{j=1}^L m_j}{L} \quad S^2 = \frac{\sum_j (m_j - m)^2}{L - 1}$$

Thus an estimator of  $\sigma^2$  is  $K \cdot S^2$ , namely,

$$\hat{\sigma}^2 = K \sum_{j=1}^L \frac{(m_j - m)^2}{L - 1}$$

$$\sum_j \frac{(m_j - m)^2}{\sigma^2 / K} \sim \chi_{L-1}^2 \quad SSb \equiv K \sum_j (m_j - m)^2$$

So when  $H_0$  is true, we have

$$\frac{SSb}{\sigma^2} \sim \chi_{L-1}^2$$

Regardless of  $H_0$  our second estimator to  $\sigma^2$  is the average of group variances  $S_j^2$  :

$$S_j^2 = \frac{\sum_{i=1}^K (X_{ij} - m_j)^2}{K-1} \quad \hat{\sigma}^2 = \sum_{j=1}^L \frac{S_j^2}{L} = \sum_j \sum_i \frac{(X_{ij} - m_j)^2}{L(K-1)}$$

$$SSW \equiv \sum_j \sum_i (X_{ij} - m_j)^2$$

$$(K-1) \frac{S_j^2}{\sigma^2} \sim \chi_{K-1}^2 \quad \frac{SSW}{\sigma^2} \sim \chi_{L(K-1)}^2$$

$$\left( \frac{SSb / \sigma^2}{L-1} \right) / \left( \frac{SSW / \sigma^2}{L(K-1)} \right) = \frac{SSb / (L-1)}{SSW / (L(K-1))} \sim F_{L-1, L(K-1)}$$

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_L \text{ if } < F_{\alpha, L-1, L(K-1)}$$

# Other Tests

- Range test (Newman-Keuls): 145 23
- Nonparametric tests (Sign test, Kruskal-Wallis)
- Contrasts: Check if 1 and 2 differ from 3,4, and 5
- Multiple comparisons require **Bonferroni correction** If there are  $m$  tests, to have an overall significance of  $\alpha$ , each test should have a significance of  $\alpha/m$ .
- Regression: CLT states that the sum of iid variables from *any* distribution is approximately normal and the preceding methods can be used.
- Other loss functions ?