Chapter ML:II (continued)

II. Machine Learning Basics

- □ Regression
- □ Concept Learning: Search in Hypothesis Space
- □ Concept Learning: Search in Version Space
- Measuring Performance

True Misclassification Rate

Definition 8 (True Misclassification Rate)

Let X be a feature space with a finite number of elements. Moreover, let C be a set of classes, let $y:X\to C$ be a classifier, and let c be the target concept to be learned. Then the true misclassification rate, denoted as $\mathit{Err}^*(y)$, is defined as follows:

$$\textit{Err}^*(y) = \frac{|\{\mathbf{x} \in X : y(\mathbf{x}) \neq c(\mathbf{x})\}|}{|X|}$$

True Misclassification Rate

Definition 8 (True Misclassification Rate)

Let X be a feature space with a finite number of elements. Moreover, let C be a set of classes, let $y:X\to C$ be a classifier, and let c be the target concept to be learned. Then the true misclassification rate, denoted as $\mathit{Err}^*(y)$, is defined as follows:

$$\textit{Err}^*(y) = \frac{|\{\mathbf{x} \in X : y(\mathbf{x}) \neq c(\mathbf{x})\}|}{|X|}$$

Problem:

 \Box Usually the *total function* c is unknown.

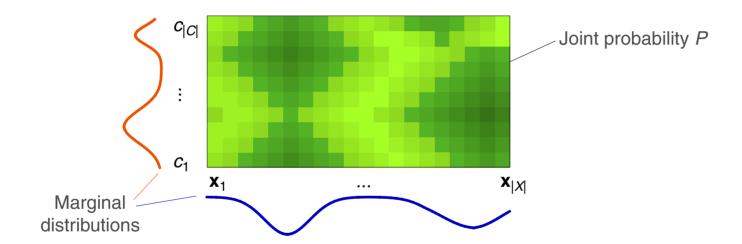
Solution:

□ Estimation of $Err^*(y)$ with $Err(y, D_s)$, i.e., evaluating y on a subset $D_s \subseteq D$ of carefully chosen examples D. Recall that for the feature vectors in D the target concept c is known.

- Instead of the term "true misclassification rate" we may also use the term "true misclassification error" or simply "true error".
- □ The English word "rate" can be used to denote both the mathematical concept of a flow quantity (a change of a quantity per time unit) as well as the mathematical concept of a *portion*, a *percentage*, or a *ratio*, which has a stationary (= time-independent) semantics. Note that the latter semantics is meant here when talking about the misclassification rate.
- □ Unfortunately, the German word "Rate" is often (mis)used to denote the mathematical concept of a portion, a percentage, or a ratio. Taking a precise mathematical standpoint, the correct German words are "Anteil" or "Quote". I.e., a semantically correct translation of misclassification rate is "Missklassifikationsanteil", and not "Missklassifikationsrate".

True Misclassification Rate: Probabilistic Foundation

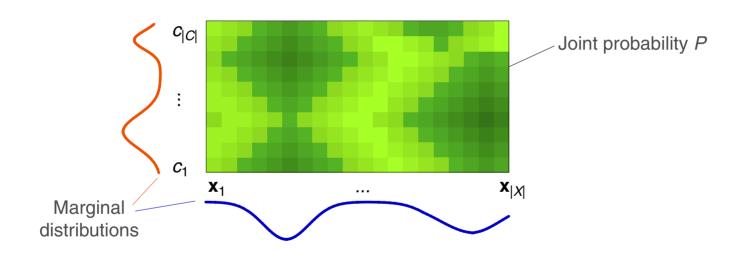
Let X and C be defined as before. Moreover, let P be a probability measure on $X \times C$. Then $P(\mathbf{x}, c)$ (precisely: $P(\mathcal{H} = \mathbf{x}, \mathcal{C} = c)$) denotes the probability that feature vector $\mathbf{x} \in X$ belongs to class $c \in C$. Illustration:



ML:II-100 Basics

True Misclassification Rate: Probabilistic Foundation

Let X and C be defined as before. Moreover, let P be a probability measure on $X \times C$. Then $P(\mathbf{x}, c)$ (precisely: $P(\mathcal{H} = \mathbf{x}, \mathcal{C} = c)$) denotes the probability that feature vector $\mathbf{x} \in X$ belongs to class $c \in C$. Illustration:

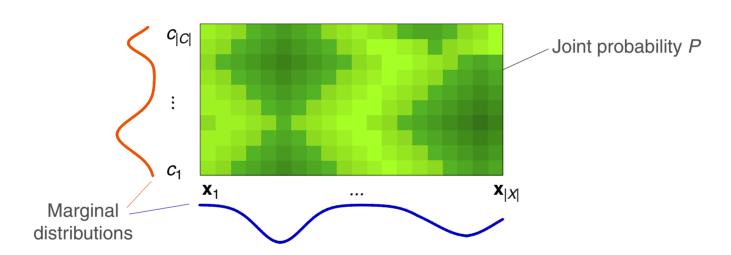


$$\underline{\textit{Err}^*(y)} = \sum_{\mathbf{x} \in X} \sum_{c \in C} P(\mathbf{x}, c) \cdot I(y(\mathbf{x}), c), \quad \text{ with } I(y(\mathbf{x}), c) = \left\{ \begin{array}{l} 0 \ \ \text{if } y(\mathbf{x}) = c \\ 1 \ \ \text{otherwise} \end{array} \right.$$

ML:II-101 Basics ©STEIN/LETTMANN 2005-2018

True Misclassification Rate: Probabilistic Foundation

Let X and C be defined as before. Moreover, let P be a probability measure on $X \times C$. Then $P(\mathbf{x}, c)$ (precisely: $P(\mathcal{H} = \mathbf{x}, \mathcal{C} = c)$) denotes the probability that feature vector $\mathbf{x} \in X$ belongs to class $c \in C$. Illustration:



$$\underline{\textit{Err}^*(y)} = \sum_{\mathbf{x} \in X} \sum_{c \in C} P(\mathbf{x}, c) \cdot I(y(\mathbf{x}), c), \quad \text{with } I(y(\mathbf{x}), c) = \left\{ \begin{array}{l} 0 \quad \text{if } y(\mathbf{x}) = c \\ 1 \quad \text{otherwise} \end{array} \right.$$

 $D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times C$ is a <u>set of examples</u> whose elements are drawn independently and according to the same P.

ML:II-102 Basics

- ☐ See the definition of a probability measure in [ML:IV Probability Basics].
- $\ \square$ $\ \mathcal{H}$ and $\ \mathcal{C}$ are random variables with domains $\ X$ and $\ \mathcal{C}$ respectively.
- Let A and B denote two events, e.g., $A = {}^{"}\mathcal{H} = \mathbf{x}"$ and $B = {}^{"}\mathcal{C} = c"$. Then the following expressions are syntactic variants, i.e., they are semantically equivalent: P(A, B), P(A and B), $P(A \wedge B)$.
- The function $c(\mathbf{x})$ has been modeled as random variable, C, since in the real world the classification of a feature vector \mathbf{x} may not be deterministic but the result of a random (measuring) process. Keyword: label noise.
- □ The elements in D are considered as random variables that are both independent of each other and identically distributed. This property of a set of random variables is abbreviated with "i.i.d."
- If the elements in D or D_s were not chosen according to P, then $Err(y, D_s)$ could not be used as an estimation of $Err^*(y)$. Keyword: sample selection bias

Training Error [True Misclassification Rate]

- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C \text{ is a set of examples.}$
- \Box $D_{tr} = D$ is the training set.
- $\neg y: X \to C$ is a classifier learned on the basis of D_{tr} .

Training error = misclassification rate with respect to D_{tr} :

$$Err(y, D_{tr}) = \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D_{tr} : y(\mathbf{x}) \neq c(\mathbf{x})\}|}{|D_{tr}|}$$

Training Error [True Misclassification Rate]

- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C \text{ is a set of examples.}$
- \Box $D_{tr} = D$ is the training set.
- $\neg y: X \to C$ is a classifier learned on the basis of D_{tr} .

Training error = misclassification rate with respect to D_{tr} :

$$Err(y, D_{tr}) = \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D_{tr} : y(\mathbf{x}) \neq c(\mathbf{x})\}|}{|D_{tr}|}$$

Problems:

- \Box *Err*(y, D_{tr}) is based on examples that are also exploited to learn y.
- \rightarrow *Err*(y, D_{tr}) quantifies memorization but not the generalization capability of y.
- \rightarrow $Err(y, D_{tr})$ is an optimistic estimation, i.e., it is constantly lower compared to the error incurred when applying y in the wild.

2-Fold Cross-Validation (Holdout Estimation) [True Misclassification Rate]

- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C \text{ is a set of examples.}$
- \square $D_{tr} \subset D$ is the training set.
- $\neg y:X\to C$ is a classifier learned on the basis of D_{tr} .
- $D_{ts} \subset D$ with $D_{ts} \cap D_{tr} = \emptyset$ is a test set.

Holdout estimation = misclassification rate with respect to D_{ts} :

$$Err(y, D_{ts}) = \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D_{ts} : y(\mathbf{x}) \neq c(\mathbf{x})\}|}{|D_{ts}|}$$

2-Fold Cross-Validation (Holdout Estimation) [True Misclassification Rate]

- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C \text{ is a set of examples.}$
- \square $D_{tr} \subset D$ is the training set.
- $\neg y:X\to C$ is a classifier learned on the basis of D_{tr} .
- $D_{ts} \subset D$ with $D_{ts} \cap D_{tr} = \emptyset$ is a test set.

Holdout estimation = misclassification rate with respect to D_{ts} :

$$Err(y, D_{ts}) = \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D_{ts} : y(\mathbf{x}) \neq c(\mathbf{x})\}|}{|D_{ts}|}$$

Requirements:

- \Box D_{tr} and D_{ts} must be goverend by the same distribution.
- \Box D_{tr} and D_{ts} should have similar sizes.

- \Box A typical value for splitting D into training set D_{tr} and test set D_{ts} is 2:1.
- \Box When splitting D into D_{tr} and D_{ts} one has to ensure that the <u>underlying distribution</u> is maintained. Keywords: stratification, sample selection bias

k-Fold Cross-Validation [Holdout Estimation]

- \neg Form k test sets by splitting D into disjoint sets D_1, \ldots, D_k of similar size.
- \Box For $i = 1, \ldots, k$ do:
 - 1. $y_i: X \to C$ is a classifier learned on the basis of $D \setminus D_i$

2.
$$Err(y_i, D_i) = \frac{|\{(\mathbf{x}, c(\mathbf{x}) \in D_i : y_i(\mathbf{x}) \neq c(\mathbf{x})\}|}{|D_i|}$$

Cross-validated misclassification rate:

$$\textit{Err}_{cv}(y,D) = rac{1}{k} \sum_{i=1}^{k} \textit{Err}(y_i,D_i)$$

n-Fold Cross-Validation (Leave One Out)

Special case with k = n:

 \Box Determine the cross-validated misclassification rate for $D \setminus D_i$ where

$$D_i = \{(\mathbf{x}_i, c(\mathbf{x}_i))\}, i \in \{1, \dots, n\}$$
.

n-Fold Cross-Validation (Leave One Out)

Special case with k = n:

Determine the cross-validated misclassification rate for $D \setminus D_i$ where $D_i = \{(\mathbf{x}_i, c(\mathbf{x}_i))\}, i \in \{1, \dots, n\}$.

Problems:

- \Box High computational effort if D is large.
- \Box Singleton test sets ($|D_i|=1$) are never stratified since they contain a single class only.

- \Box For large k the set $D \setminus D_i$ is of similar size as D. Hence $Err(y_i, D_i)$ is close to Err(y, D), where y is the classifier learned on the basis of the entire set D.
- \Box n-fold cross-validation is a special case of exhaustive cross-validation methods, which learn and test on all possible ways to divide the original sample into a training and a validation set. [Wikipedia]

Bootstrapping [Holdout Estimation]

Resampling the example set D:

- \Box For $j=1,\ldots,l$ do:
 - 1. Form training set D_i by drawing m examples from D with replacement.
 - 2. $y_j: X \to C$ is a classifier learned on the basis of D_j

3.
$$Err(y_j, D \setminus D_j) = \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D \setminus D_j : y_j(\mathbf{x}) \neq c(\mathbf{x})\}|}{|D \setminus D_j|}$$

Bootstrapped misclassification rate:

$$\textit{Err}_{bt}(y, D) = \frac{1}{l} \sum_{j=1}^{l} \textit{Err}(y_j, D \setminus D_j)$$

- Let |D| = n. The probability that an example is not considered is $(1 1/n)^m$. Similarly, the probability that an example is considered at least once is $1 (1 1/n)^m$.
- □ If m gets closer to n, then $1 (1 1/n)^m \approx 1 1/e \approx 0.632$. I.e., each training set contains about 63.2% of the examples in D.
- \Box The classifiers y_1, \ldots, y_l can be used in a combined fashion, called *ensemble*, where the class is determined by means of a majority decision:

$$y(\mathbf{x}) = \operatorname*{argmax}_{c \in C} |\{j \in \{1, \dots, l\} : y_j(\mathbf{x}) = c\}|$$

Misclassification Costs [Holdout Estimation]

Use of a cost measure for the misclassification of a feature vector \mathbf{x} in class c' instead of in class c:

$$cost(c' \mid c)$$
 $\begin{cases} \geq 0 & \text{if } c' \neq c \\ = 0 & \text{otherwise} \end{cases}$

Estimation of $\mathit{Err}^*_{\mathit{cost}}(y)$ based on a sample $D_s \subseteq D$:

$$\textit{Err}_{\textit{cost}}(y, D_s) = \frac{1}{|D_s|} \cdot \sum_{(\mathbf{x}, c(\mathbf{x})) \in D_s} \textit{cost}(y(\mathbf{x}) \mid c(\mathbf{x}))$$

□ The misclassification rate *Err* is a special case of Err_{cost} with $cost(c' \mid c) = 1$ for $c' \neq c$.

ML:II-116 Basics © STEIN/LETTMANN 2005-2018