

2. Classification Evaluation

References

T. Mitchell. Machine Learning. Chapter 5

2.1 Statistical Evaluation

Performance evaluation in classification based on accuracy or error rate.

Consider a typical classification problem:

$f: X \rightarrow Y$

X instance space

D : prob distribution over X

S : dataset: a sample from X

Consider a hypothesis h , **solution of a learning algorithm** obtained from S and estimate the accuracy:

Errors

True error: is the prob that h will misclassify an instance drawn at random according to D . (true error cannot be computed)

$$\text{error}_D(h) \equiv \Pr_{x \in D}[f(x) \neq h(x)]$$

Sample error: is the proportion of examples h misclassifies (computed only on a small data sample)

$$\text{error}_S(h) \equiv \frac{1}{n} \sum_{x \in S} \delta(f(x) \neq h(x))$$

where $\delta = 1$ if $f(x) \neq h(x)$, 0 otherwise. **Accuracy(h)** = 1 - error(h)

the goal of a learning system is to be accurate in $h(x)$, for every x not in S

If $\text{accuracy}_S(h)$ is very high, but $\text{accuracy}_D(h)$ is poor, then our system would not be very useful.

Probabilities

bias = $E[\text{error}_S(h)] - \text{error}_D(h)$

1. If S is the training set used to compute h , $\text{error}_S(h)$ is optimistically biased
2. For unbiased estimate, h and S must be chosen independently
 $E[\text{error}_S(h)] = \text{error}_D(h)$
3. Even with unbiased S , $\text{error}_S(h)$ may still vary from $\text{error}_D(h)$.
The smaller the set S , the greater the expected variance.

Confidence intervals

With approximately $N\%$ probability, $\text{error}_D(h)$ lies in interval

$$\text{error}_S(h) \pm z_N \sqrt{\frac{\text{error}_S(h)(1 - \text{error}_S(h))}{n}}$$

where

$N\%$:	50%	68%	80%	90%	95%	98%	99%
z_N :	0.67	1.00	1.28	1.64	1.96	2.33	2.58

Estimators

How to compute $\text{error}_S(h)$

1. Partition the data set D ($D = T \cup S$, T Inters. $S = \text{Vuoto}$, $|T| = 2/3|D|$)
2. Compute a hypothesis h using training set T
3. Evaluate $\text{error}_S(h)$

$\text{error}_S(h)$ is an **unbiased estimator** for $\text{error}_D(h)$; since true error is not computable, but we need it to evaluate h , we must use an estimate

Trade off between training and testing

- More training and less testing improve performance: but E_S does not approximate well E_D
- More training and less testing reduces variance

USE: 2/3 TRAINING, 1/3 TESTING

Comparisons

True comparison: $d = E_d(h_1) - E_d(h_2)$

Its estimator: $d^{\wedge} = E_{s1}(h_1) - E_{s2}(h_2)$

d^{\wedge} is unbiased iff all parameters are independent: $E[d^{\wedge}] = d$

Overfitting: h in H **overfits** training data if there is h' in H :

$$E_s(h) < E_s(h') \quad E_d(h) > E_d(h')$$

Rem.: h is solution of learning algorithm L when using training set $T \rightarrow h = L(T)$

K-fold Cross Validation

1. Partition data set D into k disjoint sets S_1, S_2, \dots, S_k ($|S_i| > 30$)

2. For $i = 1, \dots, k$ do

 use S_i as test set, and the remaining data as training set T_i

$$T_i = \{D - S_i\}$$

$$h_i = L(T_i)$$

$$\text{deltai} = E_{S_i}(h_i)$$

3. Return

$$\text{error}(L, D) = 1/k \text{ Sum } \text{deltai}$$

Comparing Algorithm

$$E(S \text{ in } D)[\text{error}_D(L_A(S)) - \text{error}_D(L_B(S))]$$

where $L(S)$ is the hypothesis output by learner L using training set S
i.e., the expected difference in true error between hypotheses output by learners L_A and L_B ; can be approximated by a K-Fold Cross Validation.

Using A and $B \rightarrow$

$$\bar{\delta} \equiv \frac{1}{k} \sum_{i=1}^k \delta_i$$

Note: if $\bar{\delta} < 0$ we can estimate that L_A is better than L_B .

2.2 Performance metrics

Accuracy is not always a good performance metric, especially with unbalanced data.

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

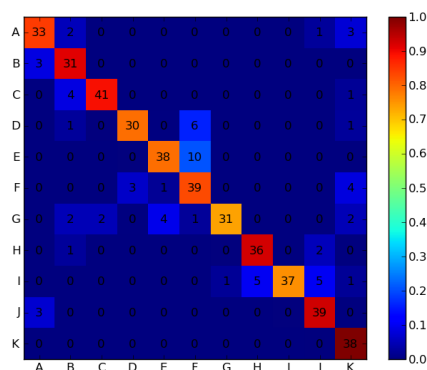
Error rate = $\frac{\text{errors}}{\text{instances}} = \frac{FN + FP}{TP + TN + FP + FN}$

Accuracy = $1 - \text{Error rate} = \frac{TP + TN}{TP + TN + FP + FN}$

Problems when datasets are unbalanced.

Confusion Matrix

In a classification problem with many classes, we can compute how many times an instance of class C_i is classified in class C_j .



Main diagonal contains accuracy;
outside the diagonal the errors.