

Generalization Quiz

1. Generalization is about how well our classifier does on testing data / future data.

2. Overfitting - predictor too flexible
 - can capture every detail of the data
 - Predictor A vs. B.

$$E_{\text{train A}} < E_{\text{train B}}$$

$$E_{\text{test A}} > E_{\text{test B}}$$

underfitting - predictor not flexible enough
 - cannot capture relevant patterns

$$E_{\text{train A}} > E_{\text{train B}}$$

$$E_{\text{test A}} > E_{\text{test B}}$$

3. Complexity control

Naive Bayes - no. of attrs and limits on distribution parameters

Decision Trees - no. of nodes and/or pruning confidence

lin. Regr - degree of polynomial / no. of attributes.

4.
$$E_{\text{train}} = \frac{1}{n} \sum_{i=1}^n \text{error}(f_0(x_i), y_i)$$

\downarrow error fn \downarrow trn. instance \downarrow label
 predictor

5. We can't est. gen err from trn err!

6. gen err - use the estimate to tell us how well we expect to perform on unseen data.

Testing err - est. gen err.

Training err - min. to build best predictor

7. If we use N randomly selected testing instances, and get an error rate of E , then the following are true:

- The error rate on another rand-sel set will be dist. approximately Gaussian w/ mean E and variance $E(1-E)/N$

$$\mu \pm \sqrt{t} \cdot \phi \rightarrow E \pm \sqrt{\frac{E(1-E)}{N}} \cdot \phi$$

- E is our best estimate of the error rate in any set of ran-sel-set.

8. Properties of confidence intervals

- CI varies roughly with the square root of the no. of samples

- CI is fully specified by a confidence level, a mean and an interval either side mean (variance)

- If N is reasonably large, then 95% CI for E_{rate} in random sample of size N is about ± 2 SDs on either side of the mean.

9. We use CI to describe the range of error rates we'd expect to see when testing future unseen sets of instances.
10. Suppose we want to decide whether NB or DT are better for classification tasks, which attr to use (NB), no. of nodes (DT). How would we do it?
- Divide N instances to 3 sets
 - K training \rightarrow train on all attr/nodes
 - L validation \rightarrow test on this, pick best performance.
 - M testing \rightarrow error rate, confidence interval
11. Cross validation
- every instance is used for testing
 - pick subsets in turn, train on other subsets and test on the one we picked, then average.
 - every instance is used for training
 - less likely to get biased testing set
 - each fold has distinct test & training set.
12. Leave one out
- cross validation: all but one of the instances are used for training.
13. Stratification.
- Deals w/ problems in k -fold
 - not useful for leave one out
 - ensures test/trn sets have representative balance of classes.