

Lecture 13: Validation

■ Resampling methods

- Holdout
- Cross Validation
 - Random Subsampling
 - K-Fold Cross-Validation
 - Leave-one-out
- The Bootstrap
 - Bias and variance estimation

■ Three-way data partitioning



Introduction

- **Almost invariably, all the pattern recognition techniques that we have introduced have one or more free parameters**

- The number of neighbors in a kNN classification rule (or the kNN density estimation method)
- The bandwidth of the kernel function in kernel density estimation
- The network size, learning parameters and weights in Multilayer Perceptrons

- **Two questions arise at this point**

- How do we select the “optimal” parameter(s) for a given classification problem?
- Once we have chosen a model, how do we estimate its true error rate?
 - The true error rate is the classifier’s error rate when tested on the ENTIRE POPULATION

- **If we had access to an unlimited number of examples these questions have a straightforward answer**

- Choose the model that provides the lowest error rate on the entire population
 - And, of course, that error rate is the true error rate

- **In real applications we only have access to a finite set of examples, usually smaller than we wanted**

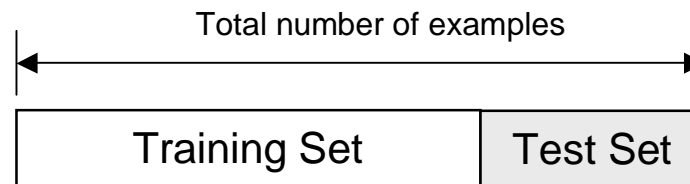
- One approach is to use the entire training data to select our classifier and estimate the error rate
 - This naïve approach has two fundamental problems
 - The final model will normally overfit the training data: it will not be able to generalize to new data
 - The problem of overfitting is more pronounced with models that have a large number of parameters
 - The error rate estimate will be overly optimistic (lower than the true error rate)
 - In fact, it is not uncommon to have 100% correct classification on training data
- A much better approach is to split the training data into disjoint subsets: the holdout method



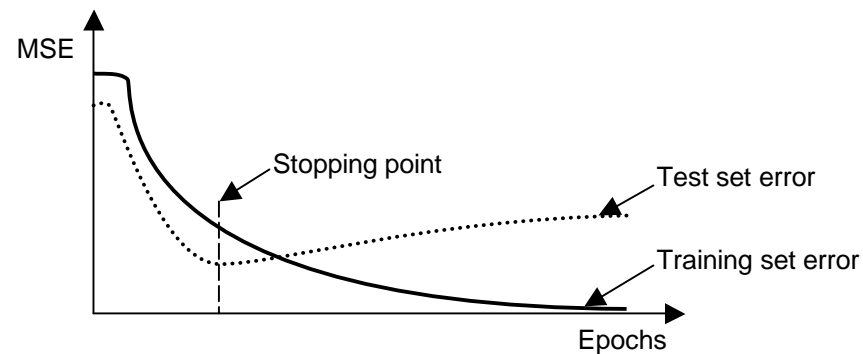
The holdout method (1)

■ Split dataset into two groups

- Training set: used to train the classifier
- Test set: used to estimate the error rate of the trained classifier



■ A typical application the holdout method is determining a stopping point for the back propagation error



The holdout method (2)

- **The holdout method has two basic drawbacks**

- In problems where we have a sparse dataset we may not be able to afford the “luxury” of setting aside a portion of the dataset for testing
- Since it is a single train-and-test experiment, the holdout estimate of error rate will be misleading if we happen to get an “unfortunate” split

- **The limitations of the holdout can be overcome with a family of resampling methods at the expense of higher computational cost**

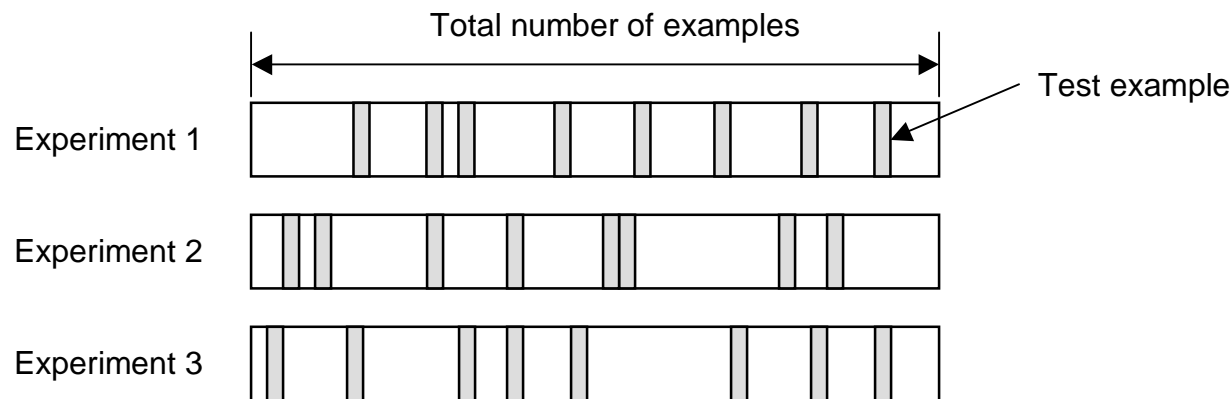
- Cross Validation
 - Random Subsampling
 - K-Fold Cross-Validation
 - Leave-one-out Cross-Validation
- Bootstrap



Random Subsampling

■ Random Subsampling performs K data splits of the entire dataset

- Each data split randomly selects a (fixed) number of examples without replacement
- For each data split we retrain the classifier from scratch with the training examples and then estimate E_i with the test examples



■ The true error estimate is obtained as the average of the separate estimates E_i

- This estimate is significantly better than the holdout estimate

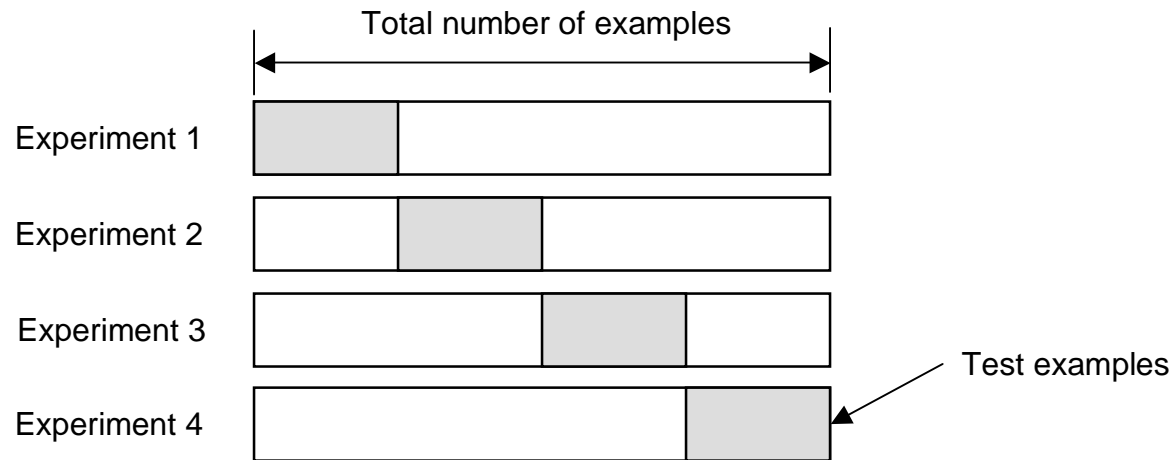
$$E = \frac{1}{K} \sum_{i=1}^K E_i$$



K-Fold Cross-validation

■ Create a K-fold partition of the the dataset

- For each of K experiments, use K-1 folds for training and a different fold for testing
- This procedure is illustrated in the following diagram for K=4



■ K-Fold Cross validation is similar to Random Subsampling

- The advantage of K-Fold Cross validation is that all the examples in the dataset are eventually used for both training and testing

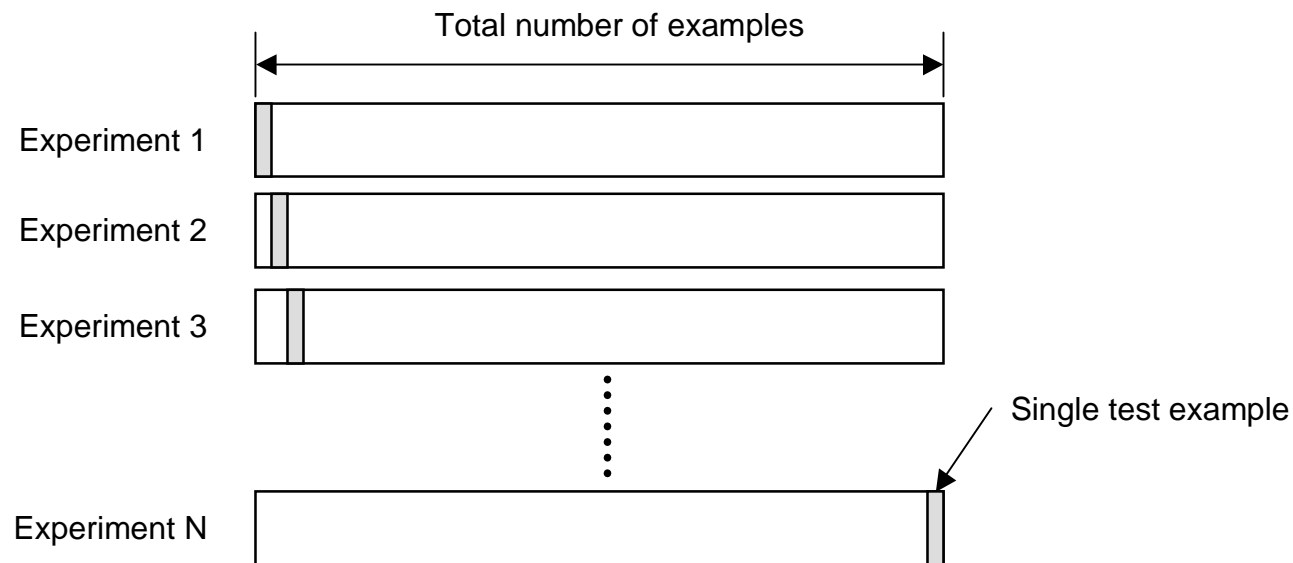
■ As before, the true error is estimated as the average error rate on test examples

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$



Leave-one-out Cross Validation

- Leave-one-out is the degenerate case of K-Fold Cross Validation, where K is chosen as the total number of examples
 - For a dataset with N examples, perform N experiments
 - For each experiment use N-1 examples for training and the remaining example for testing



- As usual, the true error is estimated as the average error rate on test examples

$$E = \frac{1}{N} \sum_{i=1}^N E_i$$



How many folds are needed?

- **With a large number of folds**

- + The bias of the true error rate estimator will be small (the estimator will be very accurate)
- The variance of the true error rate estimator will be large
- The computational time will be very large as well (many experiments)

- **With a small number of folds**

- + The number of experiments and, therefore, computation time are reduced
- + The variance of the estimator will be small
- The bias of the estimator will be large (conservative or smaller than the true error rate)

- **In practice, the choice of the number of folds depends on the size of the dataset**

- For large datasets, even 3-Fold Cross Validation will be quite accurate
- For very sparse datasets, we may have to use leave-one-out in order to train on as many examples as possible

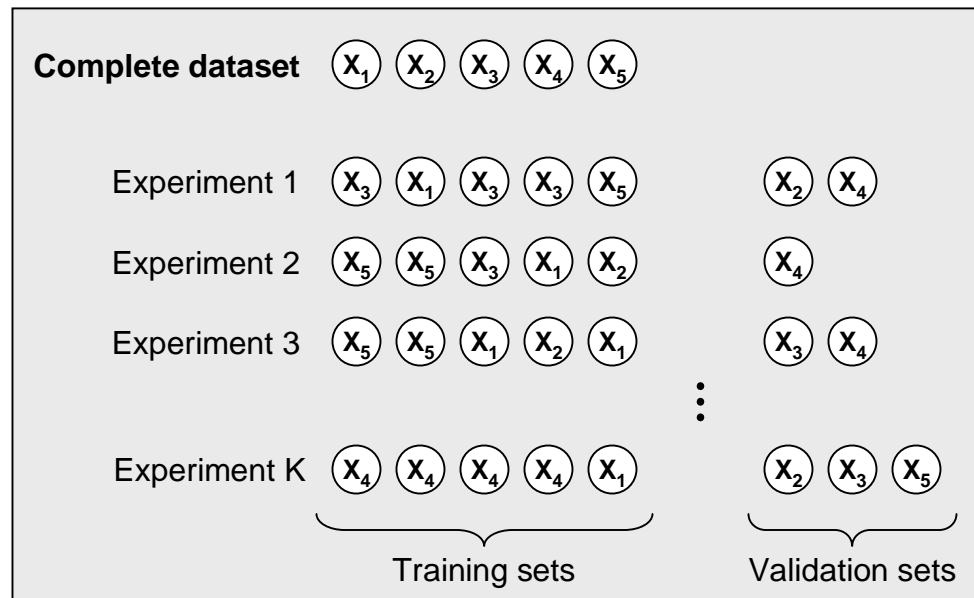
- **A common choice for K-Fold Cross Validation is K=10**



The bootstrap (1)

■ The bootstrap is a resampling technique with replacement

- From a dataset with N examples
 - Randomly select (with replacement) N examples and use this set for training
 - The remaining examples that were not selected for training are used for testing
 - This value is likely to change from fold to fold
- Repeat this process for a specified number of folds (K)



- As usual, the true error is estimated as the average error rate on test examples

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$



The bootstrap (2)

- **Compared to basic cross-validation, the bootstrap increases the variance that can occur in each fold [Efron and Tibshirani, 1993]**
 - This is a desirable property since it is a more realistic simulation of the real-life experiment from which we obtained our dataset
- **Consider a classification problem with C classes, a total of N examples and N_i examples for each class ω_i**
 - The a priori probability of choosing an example from class ω_i is N_i/N
 - Once we choose an example from class ω_i , if we do not replace it for the next selection, then the a priori probabilities will have changed since the probability of choosing an example from class ω_i will now be $(N_i-1)/N$
 - Sampling with replacement preserves the a priori probabilities of the classes throughout the random selection process
- **An additional benefit of the bootstrap is its ability to obtain accurate measures of BOTH the bias and variance of the true error estimate**



Bias and variance of a statistical estimate

- Consider the problem of estimating a parameter α of n unknown distribution G

- To emphasize the fact that α concerns G we will refer to it as $\alpha(G)$

- We collect N examples $X=\{x_1, x_2, \dots, x_N\}$ from the distribution G

- These examples define a discrete distribution G' with mass $1/N$ at each of the examples
- We compute the statistic $\alpha'=\alpha(G')$ as an estimator of $\alpha(G)$
 - In the context of this lecture, $\alpha(G')$ is the estimate of the true error rate for our classifier

- How good is this estimator?

- The “goodness” of a statistical estimator is measured by

- BIAS: How much it deviates from the true value

$$\text{Bias} = E_G[\alpha'(G)] - \alpha(G)$$

$$\text{where } E_G[X] = \int_{-\infty}^{+\infty} x g(x) dx$$

- VARIANCE: How much variability it shows for different samples $X=\{x_1, x_2, \dots, x_N\}$ of the population G

$$\text{Var} = E_G[(\alpha' - E_G[\alpha'])^2]$$

- If we are trying to estimate the mean of the population with the sample mean

- The bias of the sample mean is known to be ZERO
- The standard deviation of the sample mean is, from elementary statistics, equal to

$$\text{std}(\bar{x}) = \sqrt{\frac{1}{N(N-1)} \sum_{i=1}^N (x_i - \bar{x})^2}$$

- This term is also known in statistics as the STANDARD ERROR
- Unfortunately, there is no such a neat algebraic formula for almost any estimate other than the sample mean



Bias and variance estimates with the bootstrap

- The bootstrap, with its elegant simplicity, allows us to estimate bias and variance for practically any statistical estimate, be it a scalar or vector (matrix)
 - Here we will only describe the estimation procedure
 - Timothy Masters has an excellent introduction to the bootstrap in his textbook “*Advanced algorithms for neural networks*” if you are interested in the details
- The bootstrap estimate of bias and variance
 - Consider a dataset of N examples $X=\{x_1, x_2, \dots, x_N\}$ from the distribution G
 - This dataset defines a discrete distribution G'
 - Compute $\alpha'=\alpha(G')$ as our initial estimate of $\alpha(G)$
 - Let $\{x_1^*, x_2^*, \dots, x_N^*\}$ be a bootstrap dataset drawn from $X=\{x_1, x_2, \dots, x_N\}$
 - Estimate the parameter α using this bootstrap dataset $\alpha^*(G^*)$
 - Generate K bootstrap datasets and obtain K estimates $\{\alpha^{*1}(G^*), \alpha^{*2}(G^*), \dots, \alpha^{*K}(G^*)\}$
 - The rationale in the bootstrap method is that the effect of generating a bootstrap dataset from the distribution G' is similar to the effect of obtaining the dataset $X=\{x_1, x_2, \dots, x_N\}$ from the original distribution
 - In other words, the distribution $\{\alpha^{*1}(G^*), \alpha^{*2}(G^*), \dots, \alpha^{*K}(G^*)\}$ is related to the initial estimate α' in the same fashion as multiple estimates α' are related to the true value α , so the bias and variance estimates of α' are

$$\text{Bias}(\alpha') = \frac{1}{K} \sum_{i=1}^K \alpha^{*i} - \alpha'$$
$$\text{Var}(\alpha') = \sqrt{\frac{1}{K-1} \sum_{i=1}^K \left(\alpha^{*i} - \frac{1}{K} \sum_{i=1}^K \alpha^{*i} \right)^2}$$



Three-way data splits (1)

- **If model selection and true error estimates are to be computed simultaneously, the data needs to be divided into three disjoint sets [Ripley, 1996]**
 - **Training set:** a set of examples used for learning: to fit the parameters of the classifier
 - In the MLP case, we would use the training set to find the “optimal” weights with the back-prop rule
 - **Validation set:** a set of examples used to tune the parameters of of a classifier
 - In the MLP case, we would use the validation set to find the “optimal” number of hidden units or determine a stopping point for the back propagation algorithm
 - **Test set:** a set of examples used only to assess the performance of a fully-trained classifier
 - In the MLP case, we would use the test to estimate the error rate after we have chosen the final model (MLP size and actual weights)
 - After assessing the final model with the test set, YOU MUST NOT further tune the model
- **Why separate test and validation sets?**
 - The error rate estimate of the final model on validation data will be biased (smaller than the true error rate) since the validation set is used to select the final model
 - After assessing the final model with the test set, YOU MUST NOT tune the model any further
- **Procedure outline**

1. Divide the available data into training, validation and test set
2. Select architecture and training parameters
3. Train the model using the training set
4. Evaluate the model using the validation set
5. Repeat steps 2 through 4 using different architectures and training parameters
6. Select the best model and train it using data from the training and validation set
7. Assess this final model using the test set

- This outline assumes a holdout method
 - If CV or Bootstrap are used, steps 3 and 4 have to be repeated for each of the K folds



Three-way data splits (2)

