

# EVALUATING HYPOTHESES

- Estimating the accuracy of a hypothesis is relatively straightforward when data is plentiful.
- when we must learn a hypothesis and estimate its future accuracy given only a limited set of data, two key difficulties arise:

**Bias in the estimate:** First, the observed accuracy of the learned hypothesis over the training examples is often a poor estimator of its accuracy over future examples. To obtain an unbiased estimate of future accuracy, we typically test the hypothesis on some set of test examples chosen independently of the training examples and the hypothesis

**Variance in the estimate:** Second, even if the hypothesis accuracy is measured over an unbiased set of test examples independent of the training examples, the measured accuracy can still vary from the true accuracy, depending on the makeup of the particular set of test examples. The smaller the set of test examples, the greater the expected variance.

## ESTIMATING HYPOTHESIS ACCURACY

### Notation

- $X$ : the space of instances
- $D$ : the probability distribution of encountering instances from  $X$
- $f$ : the target function
- $H$ : the hypothesis space
- $h$ : a particular hypothesis in  $H$
- $(x, f(x))$ : a training instance
- $S$ : all training instances

## **Two Questions**

- Given  $h$  constructed from  $n$  examples drawn randomly from  $D$ , what is the best estimate of  $h$  over future instances drawn from  $D$ ?
- What is the probable error in this accuracy estimate?

**sample error:** The error rate of the hypothesis over the sample of data that is available.

**true error :** The error rate of the hypothesis over the entire unknown distribution  $\mathcal{D}$  of examples..

## True error vs. sample error

---

The **true error** of hypothesis  $h$  with respect to target function  $f$  and distribution  $\mathcal{D}$  is the probability that  $h$  will misclassify an instance drawn at random according to  $\mathcal{D}$ .

$$error_{\mathcal{D}}(h) \equiv \Pr_{x \in \mathcal{D}} [f(x) \neq h(x)]$$

The **sample error** of  $h$  with respect to target function  $f$  and data sample  $S$  is the proportion of examples  $h$  misclassifies

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

Where  $\delta(f(x) \neq h(x))$  is 1 if  $f(x) \neq h(x)$ , and 0 otherwise.

Q: How well does  $error_S(h)$  estimate  $error_{\mathcal{D}}(h)$ ?

## Problems Estimating Error

---

1. *Bias*: If  $S$  is training set,  $error_S(h)$  is optimistically biased

$$bias \equiv E[error_S(h)] - error_{\mathcal{D}}(h)$$

For unbiased estimate,  $h$  and  $S$  must be chosen independently

2. *Variance*: Even with unbiased  $S$ ,  $error_S(h)$  may still *vary* from  $error_{\mathcal{D}}(h)$

## Example

---

Hypothesis  $h$  misclassifies 12 of the 40 examples in  $S$

$$error_S(h) = \frac{12}{40} = .30$$

Q: What is  $error_{\mathcal{D}}(h)$ ?

- It is like estimating  $Pr(tail)$  from the results of a series of coin-tossing experiments, i.e.,  $Pr(tail) = \frac{n_T}{N}$ , where  $n_T$  is the number of tail events.
- That is,  $error_S(h)$  is a natural *estimator* for  $error_{\mathcal{D}}(h)$ , but,  $error_S(h)$  will be different for different choice of  $S$  just as  $\frac{n_T}{N}$  has experimental variation.
- We need a statistical measure of the confidence about the estimator  $error_S(h)$ .

# Estimators

---

- Experiment:
  1. choose sample  $S$  of size  $n$  according to distribution  $\mathcal{D}$
  2. measure  $error_S(h)$
- $error_S(h)$  is a random variable (i.e., result of an experiment)
- $error_S(h)$  is an unbiased *estimator* for  $error_{\mathcal{D}}(h)$
- Given observed  $error_S(h)$ , what can we conclude about  $error_{\mathcal{D}}(h)$ ?

## Confidence Intervals for Discrete-Valued Hypotheses

Here we give an answer to the question “How good an estimate of  $error_{\mathcal{D}}(h)$  is provided by  $error_S(h)$ ?” for the case in which  $h$  is a discrete-valued hypothesis. More specifically, suppose we wish to estimate the true error for some discrete-valued hypothesis  $h$ , based on its observed sample error over a sample  $S$ , where

- the sample  $S$  contains  $n$  examples drawn independent of one another, and independent of  $h$ , according to the probability distribution  $\mathcal{D}$
- $n \geq 30$
- hypothesis  $h$  commits  $r$  errors over these  $n$  examples (i.e.,  $error_S(h) = r/n$ ).



---

## Confidence Intervals

---

If

- $S$  contains  $n$  examples, drawn independently of  $h$  and each other
- $n \geq 30$

Then

- With approximately 95% probability,  $error_{\mathcal{D}}(h)$  lies in interval

$$error_S(h) \pm 1.96 \sqrt{\frac{error_S(h)(1 - error_S(h))}{n}}$$

If

- $S$  contains  $n$  examples, drawn independently of  $h$  and each other
- $n \geq 30$

Then

- With approximately  $N\%$  probability,  $error_{\mathcal{D}}(h)$  lies in interval

$$error_S(h) \pm z_N \sqrt{\frac{error_S(h)(1 - error_S(h))}{n}}$$

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

- A *random variable* can be viewed as the name of an experiment with a probabilistic outcome. Its value is the outcome of the experiment.
- A *probability distribution* for a random variable  $Y$  specifies the probability  $\Pr(Y = y_i)$  that  $Y$  will take on the value  $y_i$ , for each possible value  $y_i$ .
- The *expected value*, or *mean*, of a random variable  $Y$  is  $E[Y] = \sum_i y_i \Pr(Y = y_i)$ . The symbol  $\mu_Y$  is commonly used to represent  $E[Y]$ .
- The *variance* of a random variable is  $\text{Var}(Y) = E[(Y - \mu_Y)^2]$ . The variance characterizes the width or dispersion of the distribution about its mean.
- The *standard deviation* of  $Y$  is  $\sqrt{\text{Var}(Y)}$ . The symbol  $\sigma_Y$  is often used to represent the standard deviation of  $Y$ .
- The *Binomial distribution* gives the probability of observing  $r$  heads in a series of  $n$  independent coin tosses, if the probability of heads in a single toss is  $p$ .
- The *Normal distribution* is a bell-shaped probability distribution that covers many natural phenomena.
- The *Central Limit Theorem* is a theorem stating that the sum of a large number of independent, identically distributed random variables approximately follows a Normal distribution.
- An *estimator* is a random variable  $Y$  used to estimate some parameter  $p$  of an underlying population.
- The *estimation bias* of  $Y$  as an estimator for  $p$  is the quantity  $(E[Y] - p)$ . An unbiased estimator is one for which the bias is zero.
- A  $N\%$  *confidence interval* estimate for parameter  $p$  is an interval that includes  $p$  with probability  $N\%$ .

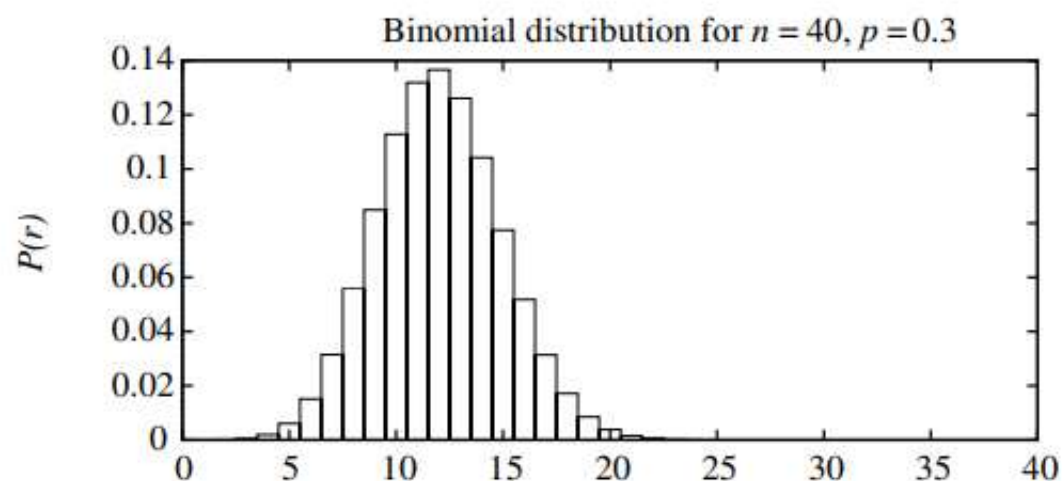
## Error Estimation and Estimating Binomial Proportions

Imagine that we were to run  $k$  such random experiments, measuring the random variables  $\text{error}_{s_1}(h)$ ,  $\text{error}_{s_2}(h)$  . . .  $\text{error}_{s_k}(h)$

$error_S(h)$  is a Random Variable

Rerun the experiment with different randomly drawn  $S$  (of size  $n$ )

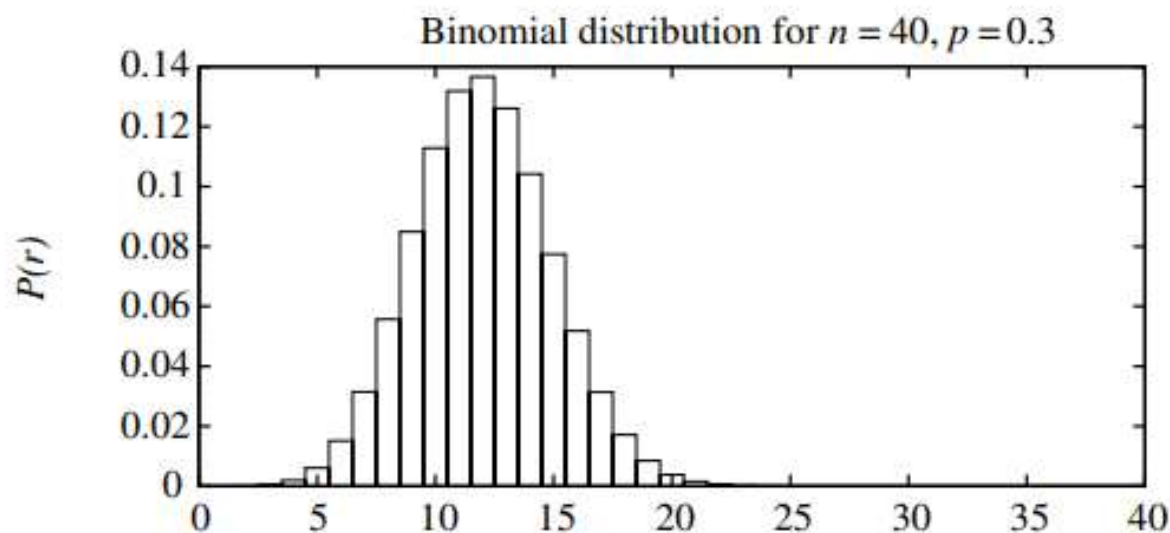
Probability of observing  $r$  misclassified examples:



$$P(r) = \frac{n!}{r!(n-r)!} error_{\mathcal{D}}(h)^r (1 - error_{\mathcal{D}}(h))^{n-r}$$

# Binomial Probability Distribution

---



$$P(r) = \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r}$$

Probability  $P(r)$  of  $r$  heads in  $n$  coin flips, if  $p = \Pr(\text{heads})$

- Expected, or mean value of  $X$ ,  $E[X]$ , is

$$E[X] \equiv \sum_{i=0}^n iP(i) = np$$

- Variance of  $X$  is

$$Var(X) \equiv E[(X - E[X])^2] = np(1 - p)$$

- Standard deviation of  $X$ ,  $\sigma_X$ , is

$$\sigma_X \equiv \sqrt{E[(X - E[X])^2]} = \sqrt{np(1 - p)}$$

**$error_S(h)$  is a Random Variable**

$error_S(h)$  follows a *Binomial* distribution, with

- mean  $\mu_{error_S(h)} = error_D(h)$
- standard deviation  $\sigma_{error_S(h)}$

$$\sigma_{error_S(h)} = \sqrt{\frac{error_D(h)(1 - error_D(h))}{n}}$$



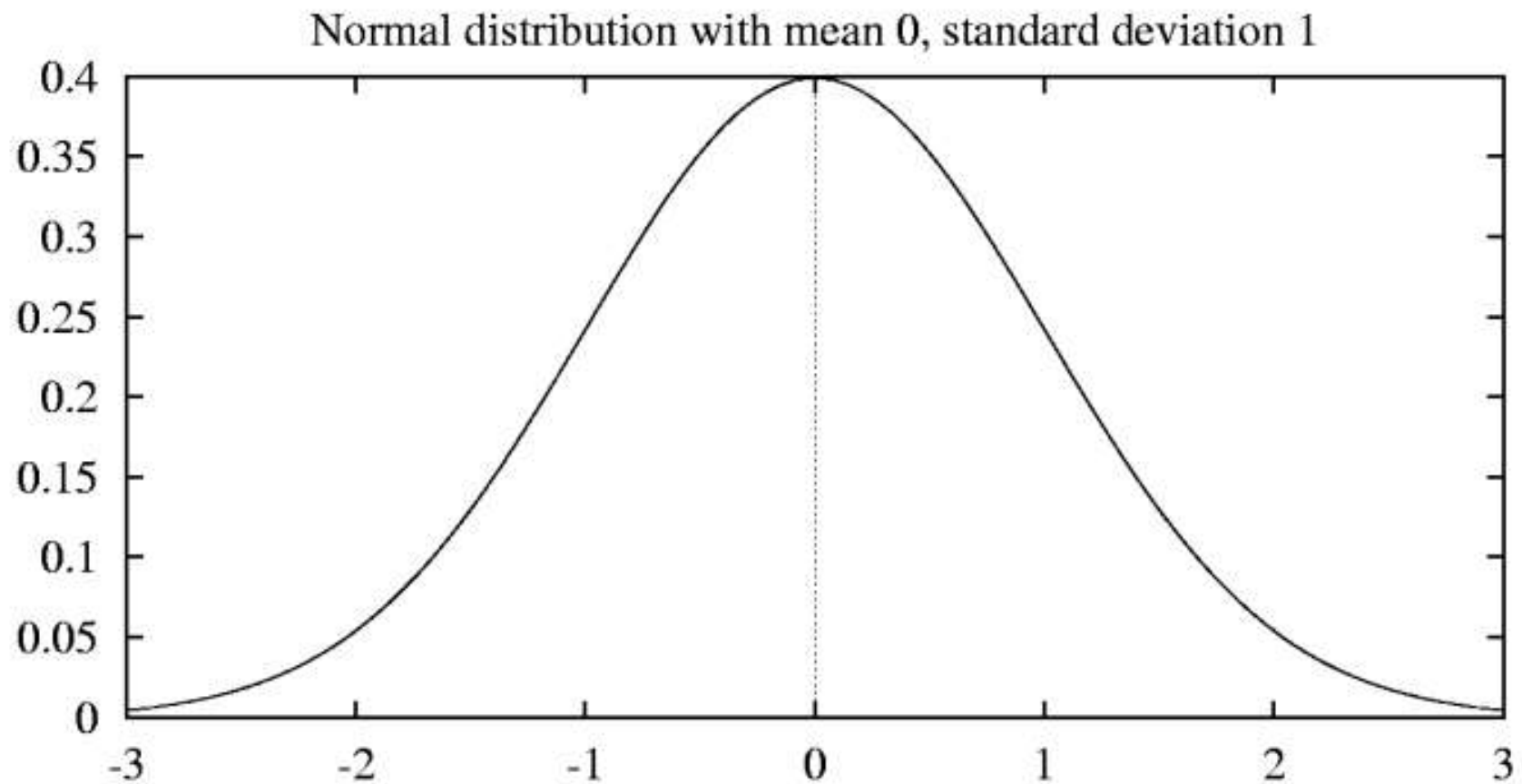
Approximate this by a *Normal* distribution with

- mean  $\mu_{error_S(h)} = error_{\mathcal{D}}(h)$
- standard deviation  $\sigma_{error_S(h)}$

$$\sigma_{error_S(h)} \approx \sqrt{\frac{error_S(h)(1 - error_S(h))}{n}}$$

# Normal Probability Distribution

---



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

The probability that  $X$  will fall into the interval  $(a, b)$  is given by

$$\int_a^b p(x) dx$$

- Expected, or mean value of  $X$ ,  $E[X]$ , is

$$E[X] = \mu$$

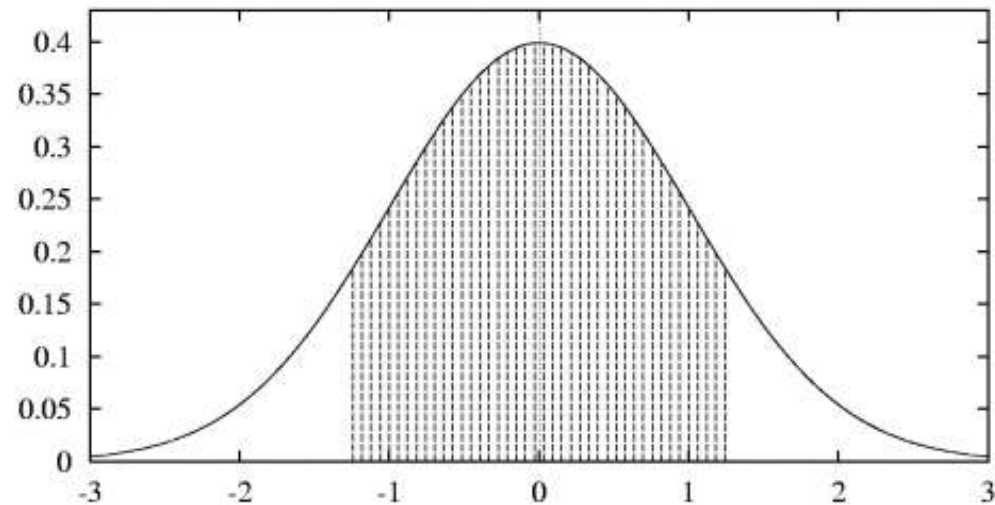
- Variance of  $X$  is

$$Var(X) = \sigma^2$$

- Standard deviation of  $X$ ,  $\sigma_X$ , is

$$\sigma_X = \sigma$$

## Normal Probability Distribution



- 80% of area (probability) lies in  $\mu \pm 1.28\sigma$
- N% of area (probability) lies in  $\mu \pm z_N\sigma$

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

## **$error_S(h)$ is a Random Variable**

---

Approximate  $P\{error_S\}$  by a *Normal* distribution with

- mean  $\mu_{error_S(h)} = error_{\mathcal{D}}(h)$
- standard deviation  $\sigma_{error_S(h)} \approx \sqrt{\frac{error_S(h)(1-error_S(h))}{n}}$

## Confidence Intervals, More Correctly

---

- If  $S$  contains  $n$  examples, drawn independently of  $h$  and each other, and if  $n \geq 30$ , then, with approximately 95% probability,  $error_S(h)$  lies in interval

$$error_{\mathcal{D}}(h) \pm 1.96 \sqrt{\frac{error_{\mathcal{D}}(h)(1 - error_{\mathcal{D}}(h))}{n}}$$

- Equivalently,  $error_{\mathcal{D}}(h)$  lies in interval

$$error_S(h) \pm 1.96 \sqrt{\frac{error_{\mathcal{D}}(h)(1 - error_{\mathcal{D}}(h))}{n}}$$

which is approximately

$$error_S(h) \pm 1.96 \sqrt{\frac{error_S(h)(1 - error_S(h))}{n}}$$

## Central Limit Theorem

---

Consider a set of independent, identically distributed random variables  $Y_1 \dots Y_n$ , all governed by an arbitrary probability distribution with mean  $\mu$  and finite variance  $\sigma^2$ . Define the sample mean,

$$\bar{Y} \equiv \frac{1}{n} \sum_{i=1}^n Y_i$$

**Central Limit Theorem.** As  $n \rightarrow \infty$ , the distribution governing  $\bar{Y}$  approaches a Normal distribution, with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$ .

## Central Limit Theorem

The sum of a large number of independent, identically distributed random variables follows a distribution that is approximately Normal.



## Estimating Confidence Intervals In General

---

1. Pick parameter  $p$  to estimate
  - $error_{\mathcal{D}}(h)$
2. Choose an estimator
  - $error_S(h)$
3. Determine probability distribution that governs estimator
  - $error_S(h)$  governed by Binomial distribution, approximated by Normal when  $n \geq 30$
4. Find interval  $(L, U)$  such that  $N\%$  of probability mass falls in the interval
  - Use table of  $z_N$  values

## Difference Between Hypotheses

---

Test  $h_1$  on sample  $S_1$ , test  $h_2$  on  $S_2$

1. Pick parameter to estimate

$$d \equiv error_{\mathcal{D}}(h_1) - error_{\mathcal{D}}(h_2)$$

2. Choose an estimator

$$\hat{d} \equiv error_{S_1}(h_1) - error_{S_2}(h_2)$$

3. Determine probability distribution that governs estimator

$$\sigma_{\hat{d}} \approx \sqrt{\frac{error_{S_1}(h_1)(1 - error_{S_1}(h_1))}{n_1} + \frac{error_{S_2}(h_2)(1 - error_{S_2}(h_2))}{n_2}}$$

4. Find interval  $(L, U)$  such that  $N\%$  of probability mass falls in the interval

$$\hat{d} \pm z_N \sqrt{\frac{\text{error}_{s_1}(h_1)(1 - \text{error}_{s_1}(h_1))}{n_1} + \frac{\text{error}_{s_2}(h_2)(1 - \text{error}_{s_2}(h_2))}{n_2}}$$

## Paired $t$ test to compare $h_A, h_B$

---

1. Partition data into  $k$  disjoint test sets  $T_1, T_2, \dots, T_k$  of equal size, where this size is at least 30.

2. For  $i$  from 1 to  $k$ , do

$$\delta_i \leftarrow \text{error}_{T_i}(h_A) - \text{error}_{T_i}(h_B)$$

3. Return the value  $\bar{\delta}$ , where

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

---

$N\%$  confidence interval estimate for  $d$ :

$$\bar{\delta} \pm t_{N,k-1} s_{\bar{\delta}}$$

$$s_{\bar{\delta}} \equiv \sqrt{\frac{1}{k(k-1)} \sum_{i=1}^k (\delta_i - \bar{\delta})^2}$$

*Note  $\delta_i$  approximately Normally distributed*

## Comparing learning algorithms $L_A$ and $L_B$

---

What we'd like to estimate:

$$E_{S \subset \mathcal{D}}[\text{error}_{\mathcal{D}}(L_A(S)) - \text{error}_{\mathcal{D}}(L_B(S))]$$

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  $L_B$ , when trained using randomly selected training sets  $S$  drawn according to distribution  $\mathcal{D}$ .

But, given limited data  $D_0$ , what is a good estimator?

- could partition  $D_0$  into training set  $S$  and training set  $T_0$ , and measure

$$error_{T_0}(L_A(S_0)) - error_{T_0}(L_B(S_0))$$

- even better, repeat this many times and average the results (next slide)

## Comparing learning algorithms $L_A$ and $L_B$

---

1. Partition data  $D_0$  into  $k$  disjoint test sets  $T_1, T_2, \dots, T_k$  of equal size, where this size is at least 30.
2. For  $i$  from 1 to  $k$ , do
  - use  $T_i$  for the test set, and the remaining data for training set  $S_i$*
  - $S_i \leftarrow \{D_0 - T_i\}$
  - $h_A \leftarrow L_A(S_i)$
  - $h_B \leftarrow L_B(S_i)$
  - $\delta_i \leftarrow \text{error}_{T_i}(h_A) - \text{error}_{T_i}(h_B)$
3. Return the value  $\bar{\delta}$ , where

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



## Comparing learning algorithms $L_A$ and $L_B$

---

Notice we'd like to use the paired  $t$  test on  $\bar{\delta}$  to obtain a confidence interval

but not really correct, because the training sets in this algorithm are not independent (they overlap!)

more correct to view algorithm as producing an estimate of

$$E_{S \subset D_0}[\text{error}_{\mathcal{D}}(L_A(S)) - \text{error}_{\mathcal{D}}(L_B(S))]$$

instead of

$$E_{S \subset \mathcal{D}}[\text{error}_{\mathcal{D}}(L_A(S)) - \text{error}_{\mathcal{D}}(L_B(S))]$$

but even this approximation is better than no comparison