



Nebřaška Lincoln	Setting Goals (cont'd)	Nebraska Lincoln	Types of Error
CSCE 478/878 Lecture 4: Experimental Design and Analysis Stephen Scott Introduction Outline Estimating Error Comparing Learning Algorithms Other Performance Measures	 Need to note that, in addition to statistical variations, what we determine is limited to the application that we are studying E.g., if naïve Bayes better than ID3 on spam filtering, that means nothing about face recognition In planning experiments, need to ensure that training data not used for evaluation I.e., don't test on the training set! Will bias the performance estimator Also holds for validation set used to prune DT, tune parameters, etc. Validation set serves as part of training set, but not used for model building 	CSCE 473878 Lecture 4: Experimental Design and Analysis Stephen Scott Introduction Outline Goals Estimating Error Estimating Error Conflored Homy Cargaring Learning Algorithms Other Performance Measures	 For now, focus on straightforerror For hypothesis <i>h</i>, recall the error from Chapter 2: Empirical error (or sam that <i>h</i> gets wrong: error_V(h) ≡ 1/V where δ(C(x) ≠ h(x)) is 1 Generalization error (or new, randomly selected, error_D(h) ≡ where D is probability dis from

Nebraska Lincoln	Types of Error
CSCE 478/878 Lecture 4: Experimental Design and Analysis Stephen Scott	 For now, focus on straightforward, 0/1 classification error For hypothesis <i>h</i>, recall the two types of classification error from Chapter 2: Empirical error (or sample error) is fraction of set V
Introduction	that h gets wrong:
Outline Goals Estimating	$error_{\mathcal{V}}(h) \equiv rac{1}{ \mathcal{V} } \sum_{x \in \mathcal{V}} \delta(C(x) \neq h(x)) \;\;,$
Types of Error Estimating Error Confidence Intervals Comparing	where $\delta(C(x) \neq h(x))$ is 1 if $C(x) \neq h(x)$, and 0 otherwise • Generalization error (or true error) is probability that a new, randomly selected, instance is misclassified by <i>h</i>
Learninġ Algorithms Other	$error_{\mathcal{D}}(h) \equiv \Pr_{x \in \mathcal{D}}[C(x) \neq h(x)]$,
	where T is probability distribution instances are drawn

stribution instances are drawn

Why do we care about error_V(h)?

Nebraska Estimating True Error

• **Bias**: If T is training set, $error_T(h)$ is optimistically biased

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

- For unbiased estimate (bias = 0), h and \mathcal{V} must be chosen independently \Rightarrow **don't test on the training set!**
- (By the way, this is distinct from inductive bias)
- Variance: Even with unbiased V, $error_{V}(h)$ may still vary from $error_{D}(h)$

Nebraska Estimating True Error (cont'd)

8/878 cture 4: erimental

ioals

timating

Experiment:

Choose sample V of size N according to distribution D
Measure *error*_V(h)

 $error_{\mathcal{V}}(h)$ is a random variable (i.e., result of an experiment)

 $error_{\mathcal{V}}(h)$ is an **unbiased estimator** for $error_{\mathcal{D}}(h)$

Given observed $error_{\mathcal{V}}(h)$, what can we conclude about $error_{\mathcal{D}}(h)$?



ntroductior Dutline

stimating

Confidence Intervals

• \mathcal{V} contains *N* examples, drawn independently of *h* and each other

● *N* ≥ 30

e

lf

Then with approximately 95% probability, $error_{\mathcal{D}}(h)$ lies in

$$rror_{\mathcal{V}}(h) \pm 1.96 \sqrt{\frac{error_{\mathcal{V}}(h)(1 - error_{\mathcal{V}}(h))}{N}}$$

E.g. hypothesis *h* misclassifies 12 of the 40 examples in test set \mathcal{V} :

$$error_{\mathcal{V}}(h) = \frac{12}{40} = 0.30$$

Then with approx. 95% confidence, $error_{\mathcal{D}}(h) \in [0.158, 0.442]$

Nebraska Lincoln	Confidence Intervals (cont'd)								
CSCE 478/878 Lecture 4: Experimental Design and Analysis Stephen Scott Introduction Outline Goals Estimating Error Types of Error Estimating Error	If • \mathcal{V} contains N examples, drawn independently of h and each other • $N \ge 30$ Then with approximately c% probability, $error_{\mathcal{D}}(h)$ lies in $error_{\mathcal{V}}(h) \pm z_c \sqrt{\frac{error_{\mathcal{V}}(h)(1 - error_{\mathcal{V}}(h))}{N}}$								
Comparing Learning Algorithms Other Performance Measures	N%: 50% 68% 80% 90% 95% 98% 99% z_c: 0.67 1.00 1.28 1.64 1.96 2.33 2.58								



Nebraska

$error_{\mathcal{V}}(h)$ is a Random Variable

Repeatedly run the experiment, each with different randomly drawn \mathcal{V} (each of size N) Probability of observing *r* misclassified examples:





I.e., let $error_{\mathcal{D}}(h)$ be probability of heads in biased coin, then P(r) = prob. of getting r heads out of $\underset{i=1}{N} \text{flips}_{i=1}^{N}$, the second s



Nebiaska Nebraska Approximate Binomial Dist. with Normal Normal Probability Distribution $error_{\mathcal{V}}(h) = r/N$ is binomially distributed, with tandard deviation 1 0.35 0.3 • mean $\mu_{error_{\mathcal{V}}(h)} = error_{\mathcal{D}}(h)$ (i.e., unbiased est.) 0.25 • standard deviation $\sigma_{error_{\mathcal{V}}(h)}$ 0.2 0.15 $\sigma_{error_{\mathcal{V}}(h)} = \sqrt{\frac{error_{\mathcal{D}}(h)(1 - error_{\mathcal{D}}(h))}{(1 - error_{\mathcal{D}}(h))}}$ 0.1 (increasing N decreases variance) $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$ Want to compute confidence interval = interval centered at $error_{\mathcal{D}}(h)$ containing c% of the weight under the distribution • The probability that X will fall into the interval (a, b) is Approximate binomial by normal (Gaussian) dist: • mean $\mu_{error_{\mathcal{V}}(h)} = error_{\mathcal{D}}(h)$ given by $\int_a^b p(x) dx$ • standard deviation $\sigma_{error_{\mathcal{V}}(h)}$ • Expected, or mean value of X, E[X], is $E[X] = \mu$ • Variance is $Var(X) = \sigma^2$, standard deviation is $\sigma_X = \sigma$ $\sigma_{error_{\mathcal{V}}(h)} \approx \sqrt{\frac{error_{\mathcal{V}}(h)(1 - error_{\mathcal{V}}(h))}{(N_{\mathcal{V}} + \mathcal{O} + \mathcal$

Ν





10) (B) (2) (2) (2) (2) (2)



ioals

Jebraska Lincoln	Central Limit Theorem
CSCE 478/878	How can we justify approximation?
Lecture 4: xperimental Design and Analysis tephen Scott	Consider set of iid random variables Y_1, \ldots, Y_N , all from arbitrary probability distribution with mean μ and finite variance σ^2 . Define sample mean $\overline{Y} \equiv (1/N) \sum_{i=1}^n Y_i$
troduction utline pals	\bar{Y} is itself a random variable, i.e., result of an experiment (e.g., $error_S(h) = r/N$)
stimating ror pes of Error stimating Error onfidence Intervals	Central Limit Theorem : As $N \to \infty$, the distribution governing \bar{Y} approaches normal distribution with mean μ and variance σ^2/N
omparing earning gorithms ther	Thus the distribution of $error_{S}(h)$ is approximately normal for large N , and its expected value is $error_{D}(h)$
easures	(Rule of thumb: $N \ge 30$ when estimator's distribution is binomial; might need to be larger for other distributions)





Lincoin	
CSCE 478/878 Lecture 4: Experimental Design and Analysis Stephen Scott	
Introduction	l
Outline	l
Goals	l
Estimating Error	l
Comparing Learning Algorithms K-Fold CV	l
Student's r Distribution	l
Other Performance Measures	
23/35	

Nebiaska

Student's t Distribution (One-Sided Test)

df	0.600	0.700	0.800	0.900	0.950	0.975	0.990	0.995
1	0.325	0.727	1.376	3.078	6.314	12.706	31.821	63.657
2	0.289	0.617	1.061	1.886	2.920	4.303	6.965	9.925
3	0.277	0.584	0.978	1.638	2.353	3.182	4.541	5.841
4	0.271	0.569	0.941	1.533	2.132	2.776	3.747	4.604
5	0.267	0.559	0.920	1.476	2.015	2.571	3.365	4.032
6	0.265	0.553	0.906	1.440	1.943	2.447	3.143	3.707
7	0.263	0.549	0.896	1.415	1.895	2.365	2.998	3.499
8	0.262	0.546	0.889	1.397	1.860	2.306	2.896	3.355
9	0.261	0.543	0.883	1.383	1.833	2.262	2.821	3.250
10	0.260	0.542	0.879	1.372	1.812	2.228	2.764	3.169
11	0.260	0.540	0.876	1.363	1.796	2.201	2.718	3.106
12	0.259	0.539	0.873	1.356	1.782	2.179	2.681	3.055
13	0.259	0.538	0.870	1.350	1.771	2.160	2.650	3.012

If $p + t_{c,K-1} s_p < 0$ our assertion that L^1 has less error than L^2 is supported with confidence c

So if K-fold CV used, compute p, look up $t_{c,K-1}$ and check if $p < -t_{c,K-1} s_p$

One-sided test; says nothing about L^2 over L^1



Nebraska More Specific Performance Measures

Collegation of the second seco

So far, we've looked at a single error rate to compare hypotheses/learning algorithms/etc.

- This may not tell the whole story:
 - 1000 test examples: 20 positive, 980 negative
 - h¹ gets 2/20 pos correct, 965/980 neg correct, for
 - accuracy of (2+965)/(20+980)=0.967 $\bullet\,$ Pretty impressive, except that always predicting
 - negative yields accuracy = 0.980
 Would we rather have h², which gets 19/20 pos correct and 930/980 neg, for accuracy = 0.949?

10, 10, 10, 12, 12, 12, 10, 00, 00

 Depends on how important the positives are, i.e., frequency in practice and/or cost (e.g., cancer diagnosis)

Nebraska Lincoln Confusion Matrices

478/878 Lecture 4:

outline

ioals

Break down error into type: true positive, etc.

	Predicted Class				
True Class	Positive	Negative	Total		
Positive	tp : true positive	fn : false negative	p		
Negative	<i>fp</i> : false positive	tn : true negative	n		
Total	p'	<i>n'</i>	N		

Generalizes to multiple classes

Allows one to quickly assess which classes are missed the most, and into what other class

10+10+15+15+15+15-000

CSCE 478/878 Lecture 4: Experimental Design and Analysis Stephen Scott	
Introduction Outline	
Goals	l
Estimating Error	l
Comparing Learning Algorithms	
Other Performance Measures Confusion Matrices	
ROC Curves Precision-Recall	

Nebraska Lincoln ROC Curves

- Consider an ANN or SVM
- Normally threshold at 0, but what if we changed it?
- Keeping weight vector constant while changing threshold = holding hyperplane's slope fixed while moving along its normal vector



• I.e., get a set of classifiers, one per labeling of test set

• Similar situation with any classifier with confidence value, e.g., probability-based

Nebraska Lincol Plotting *tp* versus *fp*

SCE 8/878

- Consider the "always –" hyp. What is *fp*? What is *tp*? What about the "always +" hyp?
- In between the extremes, we plot TP versus FP by sorting the test examples by the confidence values

Ex	Confidence	label	Ex	Confidence	label
x_1	169.752	+	<i>x</i> ₆	-12.640	-
x_2	109.200	+	x7	-29.124	-
<i>x</i> ₃	19.210	_	<i>x</i> ₈	-83.222	-
<i>x</i> ₄	1.905	+	<i>x</i> 9	-91.554	+
<i>x</i> 5	-2.75	+	<i>x</i> ₁₀	-128.212	-











Precision-Recall Curves (cont'd)



As with ROC, can vary threshold to trade off precision against recall

Can compare curves based on containment

Use $F_\beta\text{-measure to combine at a specific point, where <math display="inline">\beta$ weights precision vs recall:

$$F_{\beta} \equiv (1+\beta^2) \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}$$