

Experimentation and Evaluation

Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Eric Eaton (UPenn), David Kauchak (Pomona), and the many others who made their course materials freely available online.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

Experimental Procedure

Learning Goals

- Describe how cross-validation (k-fold, leaveone-out) is used to evaluate model and optimize hyperparameters
- Describe how to compare models statistically using the *t*-test
- Describe how bootstrapping is used to evaluate test performance

Proper Evaluation?

Current plan

- Learn algorithm on training data (subset of full data)
- Evaluate on test data (subset of full data)
- Repeat until happy with results

Is this okay?

- No! Although we are not explicitly looking at test data, we are still "cheating" by biasing our algorithm to test data
- Once you look at / use test data, it is no longer test data!

So, how can we evaluate our algorithm during development?

Based on slide by David Kauchak













Based on slide by David Kauchak

Statistical tests

Setup

- assume some default hypothesis about data that you would like to *disprove*, called the null hypothesis
- e.g. model 1 and model 2 are not statistically different in performance

Test

- calculate test statistic from data (often assuming something about data)
- calculate *p*-value from test statistic
 - *p*-value = probability of seeing test statistic at least as extreme as one actually observed given null hypothesis is true
- compare p-value to threshold α (significance level)
- reject null hypothesis if $p < \alpha$
- note that statistically significant difference is not necessarily a large-magnitude difference

t-test Determines whether two samples come from same underlying distribution or not Null hypothesis - model 1 and model 2 accuracies are no different, i.e. come from same distribution **Assumptions** - there are a number that often are not completely true, but we are often not too far off Our formulation do "paired t-test" _ values can be thought of as pairs, calculated under same conditions (in our • case, same train/test split) gives more power than unpaired *t*-test (we have more information) • for almost all experiments, do "two-tailed" version

• no *a priori* knowledge of which model is better

Based on slide by David Kauchak

Comparing Models

Is model 2 better than model 1?

Sample 1			Sample 2			Sample 3				Sample 4				
split	M1	M2		split	M1	M2	split	M1	M2		split	M1	M2	
1	87	88		1	87	87	1	84	87		1	80	82	
2	85	84		2	92	88	2	83	86		2	84	87	
3	83	84		3	74	79	3	78	82		3	89	90	
4	80	79		4	75	86	4	80	86		4	78	82	
5	88	89		5	82	84	5	82	84		5	90	91	
6	85	85		6	79	87	6	79	87		6	81	83	
7	83	81		7	83	81	7	83	84		7	80	80	
8	87	86		8	83	92	8	83	86		8	88	89	
9	88	89		9	88	81	9	85	83		9	76	77	
10	84	85		10	77	85	10	83	85		10	86	88	
avg	85	85		avg	82	85	avg	82	85		avg	83	85	
				sdev	5.9	3.9	sdev	2.3	1.7		sdev	4.9	4.7	
<i>p</i> = 1				p = 0.15			<i>p</i> = 0.007				<i>p</i> = 0.001			



Summary : Cross-Validation

- Cross-validation generates an approximate estimate of how well the classifier will do on "unseen" data
 - as $k \rightarrow n$, model becomes more accurate (more training data)
 - ... but, CV becomes more computationally expensive (have to train k models)
 - choosing k < n is a compromise
- It is an even better idea to do CV repeatedly!



Comparing Multiple Classifiers

1) Loop for t trials:







Experimentation Good Practices Never look at your test data! During development - compare different models / hyperparameters on development data - use cross-validation to get more consistent results - if you want to be confident with results, use t-test For final evaluation, use bootstrap resampling combined with t-test to compare

Based on slide by David Kauchak

Avoiding Pitfalls

- Is my held-aside test data really representative of going out to collect new data?
- Did I repeat my entire data processing procedure on every fold of cross-validation, using only training data for that fold?
- Have I modified my algorithm so many times, or tried so many approaches, on this same data set that I (human) am overfitting it?

The Short Way (that Many People Actually Use)

- Split into only training data + validation data
- Train on training data, evaluate on validation data
- Report cross-validation performance

 possibly also training performance
- Why is this used?
 - might not be enough data to create held-out test set
 - you cannot trust that authors did not peek at test data anyway =P