

Evaluating What's been Learned

Metrics

Lecture 4

Marina Santini

Acknowledgements
Slides borrowed and adapted from:
Data Mining by I. H. Witten, E. Frank and M. A. Hall

2016

Lecture 4: Evaluating what's been learned

1

Outline

- Issues: training, testing, tuning
- Predicting performance: confidence limits
- Holdout, cross-validation, bootstrap
- Comparing schemes: the t-test
- Predicting probabilities: loss functions
- Cost-sensitive measures

2016

Lecture 4: Evaluating what's been learned

2

Output: representing structural patterns

- Many different ways of representing patterns
 - ♦ Decision trees, rules, instance-based, ...
- Also called "knowledge" representation
- Representation determines inference method
- Understanding the output is the key to understanding the underlying learning methods
- Different types of output for different learning problems (e.g. classification, regression, ...)

2016

Lecture 4: Evaluating what's been learned

3

Evaluation: the key to success

- How predictive is the model we learned?
- Error on the training data is *not* a good indicator of performance on future data
 - ♦ Otherwise 1-NN would be the optimum classifier!
- Simple solution that can be used if lots of (labeled) data is available:
 - ♦ Split data into training and test set
- However: (labeled) data is usually limited
 - ♦ More sophisticated techniques need to be used

2016

Lecture 4: Evaluating what's been learned

4

Issues in evaluation

- Statistical reliability of estimated differences in performance (→ significance tests)
- Choice of performance measure:
 - ♦ Number of correct classifications
 - ♦ Accuracy of probability estimates
 - ♦ Error in numeric predictions
- Costs assigned to different types of errors
 - ♦ Many practical applications involve costs

2016

Lecture 4: Evaluating what's been learned

5

Training and testing I

- Natural performance measure for classification problems: *error rate*
 - ♦ *Success*: instance's class is predicted correctly
 - ♦ *Error*: instance's class is predicted incorrectly
 - ♦ Error rate: proportion of errors made over the whole set of instances
- *Resubstitution error*: error rate obtained from training data
- Resubstitution error is (hopelessly) optimistic!

2016

Lecture 4: Evaluating what's been learned

6

Training and testing II

- **Test set:** independent instances that have played no part in formation of classifier
- Assumption: both training data and test data are representative samples of the underlying problem
- Test and training data may differ in nature
- Example: classifiers built using customer data from two different towns *A* and *B*
 - To estimate performance of classifier from town *A* in completely new town, test it on data from *B*

2016

Lecture 4: Evaluating what's been learned

7

Note on parameter tuning

- It is important that the test data is not used *in any way* to create the classifier
- Some learning schemes operate in two stages:
 - Stage 1: build the basic structure
 - Stage 2: optimize parameter settings
- The test data can't be used for parameter tuning!
- Proper procedure uses *three sets*: *training data*, *validation (development set) data*, and *test data*
 - Validation data is used to optimize parameters

2016

Lecture 4: Evaluating what's been learned

8

Making the most of the data

- Once evaluation is complete, *all the data* can be used to build the final classifier
- Generally, the larger the training data the better the classifier (but returns diminish)
- The larger the test data the more accurate the error estimate
- **Holdout procedure:** method of splitting original data into training and test set
 - Dilemma: ideally both training set *and* test set should be large!

2016

Lecture 4: Evaluating what's been learned

9

Predicting performance

- Assume the estimated error rate is 25%. How close is this to the true error rate?
 - Depends on the amount of test data
- Prediction is just like tossing a (biased!) coin
 - "Head" is a "success", "tail" is an "error"
- In statistics, a succession of independent events like this is called a *Bernoulli process*
 - Statistical theory provides us with confidence intervals for the true underlying proportion

2016

Lecture 4: Evaluating what's been learned

10

Confidence intervals

- We can say: p lies within a certain specified interval with a certain specified confidence
- Example: $S=750$ successes in $N=1000$ trials
 - Estimated success rate: 75%
 - How close is this to true success rate p ?
 - Answer: with 80% confidence p in [73.2, 76.7]
- Another example: $S=75$ and $N=100$
 - Estimated success rate: 75%
 - With 80% confidence p in [69.1, 80.1]

2016

Lecture 4: Evaluating what's been learned

11

Mean and variance

- Mean and variance for a Bernoulli trial: $p, p(1-p)$
- Expected success rate $f=S/N$
- Mean and variance for $f: p, p(1-p)/N$
- For large enough N , f follows a Normal distribution
- $c\%$ confidence interval $[-z \leq X \leq z]$ for random variable with 0 mean is given by:
- With a symmetric distribution:

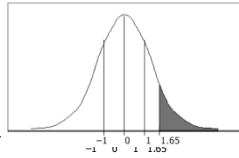
2016

Lecture 4: Evaluating what's been learned

12

Confidence limits

- Confidence limits for the normal distribution with 0 mean and a variance of 1:



$\Pr[X \geq z]$	z
0.1%	3.09
0.5%	2.58
1%	2.33
5%	1.65
10%	1.28
20%	0.84
40%	0.25

- Thus:
- To use this we have to reduce our random variable f to have 0 mean and unit variance

Examples

- $f = 75\%$, $N = 1000$, $c = 80\%$ (so that $z = 1.28$):
- $f = 75\%$, $N = 100$, $c = 80\%$ (so that $z = 1.28$):
- Note that normal distribution assumption is only valid for large N (i.e. $N > 100$)
- $f = 75\%$, $N = 10$, $c = 80\%$ (so that $z = 1.28$):

(should be taken with a grain of salt)

Holdout estimation

- What to do if the amount of data is limited?
- The *holdout* method reserves a certain amount for testing and uses the remainder for training
 - Usually: one third for testing, the rest for training
- Problem: the samples might not be representative
 - Example: class might be missing in the test data
- Advanced version uses *stratification*
 - Ensures that each class is represented with approximately equal proportions in both subsets

Repeated holdout method

- Holdout estimate can be made more reliable by repeating the process with different subsamples
 - In each iteration, a certain proportion is randomly selected for training (possibly with stratification)
 - The error rates on the different iterations are averaged to yield an overall error rate
- This is called the *repeated holdout* method
- Still not optimum: the different test sets overlap
 - Can we prevent overlapping?

Cross-validation

- Cross-validation* avoids overlapping test sets
 - First step: split data into k subsets of equal size
 - Second step: use each subset in turn for testing, the remainder for training
- Called *k-fold cross-validation*
- Often the subsets are stratified before the cross-validation is performed
- The error estimates are averaged to yield an overall error estimate

More on cross-validation

- Standard method for evaluation: stratified ten-fold cross-validation
- Why ten?
 - Extensive experiments have shown that this is the best choice to get an accurate estimate
 - There is also some theoretical evidence for this
- Stratification reduces the estimate's variance
- Even better: repeated stratified cross-validation
 - E.g. ten-fold cross-validation is repeated ten times and results are averaged (reduces the variance)

Leave-One-Out cross-validation

- Leave-One-Out: a particular form of cross-validation:
 - ♦ Set number of folds to number of training instances
 - ♦ I.e., for n training instances, build classifier n times
- Makes best use of the data
- Involves no random subsampling
- Very computationally expensive
 - ♦ (exception: NN)

2016 Lecture 4: Evaluating what's been learned 19

Leave-One-Out-CV and stratification

- Disadvantage of Leave-One-Out-CV: stratification is not possible
 - ♦ It *guarantees* a non-stratified sample because there is only one instance in the test set!
- Extreme example: random dataset split equally into two classes
 - ♦ Best inducer predicts majority class
 - ♦ 50% accuracy on fresh data
 - ♦ Leave-One-Out-CV estimate is 100% error!

2016 Lecture 4: Evaluating what's been learned 20

Comparing data mining schemes

- Frequent question: which of two learning schemes performs better?
- Note: this is domain dependent!
- Obvious way: compare 10-fold CV estimates
- Generally sufficient in applications (we don't lose if the chosen method is not truly better)
- However, what about machine learning research?
 - ♦ Need to show convincingly that a particular method works better

2016 Lecture 4: Evaluating what's been learned 21

Comparing schemes II


- Want to show that scheme A is better than scheme B in a particular domain
 - ♦ For a given amount of training data
 - ♦ On average, across all possible training sets
- Let's assume we have an infinite amount of data from the domain:
 - ♦ Sample infinitely many dataset of specified size
 - ♦ Obtain cross-validation estimate on each dataset for each scheme
 - ♦ Check if mean accuracy for scheme A is better than mean accuracy for scheme B

2016 Lecture 4: Evaluating what's been learned 22

Paired t-test

- In practice we have limited data and a limited number of estimates for computing the mean
- *Student's t-test* tells whether the means of two samples are significantly different
- In our case the samples are cross-validation estimates for different datasets from the domain
- Use a *paired t-test* because the individual samples are paired
 - ♦ The same CV is applied twice

William Gosset
 Born: 1876 in Canterbury; Died: 1937 in Beaconsfield, England
 Obtained a post as a chemist in the Guinness brewery in Dublin in 1899. Invented the t-test to handle small samples for quality control in brewing. Wrote under the name "Student".



2016 Lecture 4: Evaluating what's been learned 23

Student's distribution

- With small samples ($k < 100$) the mean follows *Student's distribution with $k-1$ degrees of freedom*
- Confidence limits:

9 degrees of freedom

Assuming we have 10 estimates

Pr[X ≥ z]	z
0.1%	4.30
0.5%	3.25
1%	2.82
5%	1.83
10%	1.38
20%	0.88

normal distribution

Pr[X ≥ z]	z
0.1%	3.09
0.5%	2.58
1%	2.33
5%	1.65
10%	1.28
20%	0.84

2016 Lecture 4: Evaluating what's been learned 24

Performing the test

- Fix a significance level
 - If a difference is significant at the $\alpha\%$ level, there is a $(100-\alpha)\%$ chance that the true means differ
- Divide the significance level by two because the test is two-tailed
 - I.e. the true difference can be +ve or - ve
- Look up the value for z that corresponds to $\alpha/2$
- If $t \leq -z$ or $t \geq z$ then the difference is significant
 - I.e. the *null hypothesis* (that the difference is zero) can be rejected

2016 Lecture 4: Evaluating what's been learned 25

Aside: statistical significance

Hypothesis

- predicts a difference between 2 classification results

Null hypothesis

- states that there will be *no* difference between the classification results

p-value

- probability value (betw 0 and 1).

Alfa

- is the significance level set for the study

2016 Lecture 4: Evaluating what's been learned 26

Predicting probabilities

- Performance measure so far: success rate
- Also called *0-1 loss function*:

□

- Most classifiers produces class probabilities
- Depending on the application, we might want to check the accuracy of the probability estimates
- 0-1 loss is not the right thing to use in those cases

2016 Lecture 4: Evaluating what's been learned 27

Informational loss function

- The informational loss function is $-\log(p_c)$, where c is the index of the instance's actual class
- Number of bits required to communicate the actual class
- Let $p_1^* \dots p_k^*$ be the true class probabilities
- Then the expected value for the loss function is:

- Justification: minimized when $p_j = p_j^*$
- Difficulty: *zero-frequency problem*

2016 Lecture 4: Evaluating what's been learned 28

Counting the cost

- In practice, different types of classification errors often incur different costs
- Examples:
 - Terrorist profiling
 - "Not a terrorist" correct 99.99% of the time
 - Loan decisions
 - Oil-slick detection
 - Fault diagnosis
 - Promotional mailing

2016 Lecture 4: Evaluating what's been learned 29

Counting the cost

- The *confusion matrix*:

		Predicted class	
		Yes	No
Actual class	Yes	True positive	False negative
	No	False positive	True negative

There are many other types of cost!

- E.g.: cost of collecting training data

2016 Lecture 4: Evaluating what's been learned 30

Binary and Multiclass confusion matrices

		Prediction	
		Positive	Negative
Actual	Positive	TP	FN
	Negative	FP	TN

		Prediction				
		Class 1	Class 2	Class 3	...	Class n
Actual	Class 1	Accurate				
	Class 2		Accurate			
	Class 3			Accurate		
	...				Accurate	
	Class n					Accurate

2016

Aside: the kappa statistic

- Two confusion matrices for a 3-class problem: actual predictor (left) vs. random predictor (right)

		Predicted class						Predicted class			
		a	b	c	total			a	b	c	total
Actual class	a	88	10	2	100	Actual class	a	60	30	10	100
	b	14	40	6	60		b	36	18	6	60
	c	18	10	12	40		c	24	12	4	40
total		120	60	20		total		120	60	20	

- Number of successes: sum of entries in diagonal (D)
- Kappa statistic:

measures relative improvement over random predictor

2016

Lecture 4: Evaluating what's been learned

32

Classification with costs

- Two cost matrices:

		Predicted class				Predicted class		
		yes	no			a	b	c
Actual class	yes	0	1	a	0	1	1	
	no	1	0	b	1	0	1	
				c	1	1	0	

- Success rate is replaced by average cost per prediction
 - Cost is given by appropriate entry in the cost matrix

2016

Lecture 4: Evaluating what's been learned

33

Cost-sensitive classification

- Can take costs into account when making predictions
 - Basic idea: only predict high-cost class when very confident about prediction
- Given: predicted class probabilities
 - Normally we just predict the most likely class
 - Here, we should make the prediction that minimizes the expected cost
 - Expected cost: dot product of vector of class probabilities and appropriate column in cost matrix
 - Choose column (class) that minimizes expected cost

2016

Lecture 4: Evaluating what's been learned

34

Lift charts

- In practice, costs are rarely known
- Decisions are usually made by comparing possible scenarios
- Example: promotional mailout to 1,000,000 households
 - Mail to all; 0.1% respond (1000)
 - Data mining tool identifies subset of 100,000 most promising, 0.4% of these respond (400)
40% of responses for 10% of cost may pay off
 - Identify subset of 400,000 most promising, 0.2% respond (800)
- A lift chart allows a visual comparison

2016

Lecture 4: Evaluating what's been learned

35

Generating a lift chart

- Sort instances according to predicted probability of being positive:

	Predicted probability	Actual class
1	0.95	Yes
2	0.93	Yes
3	0.93	No
4	0.88	Yes
...

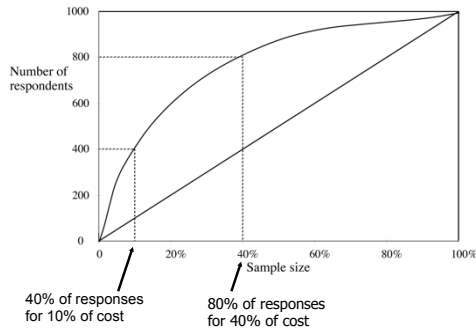
- x axis is sample size
- y axis is number of true positives

2016

Lecture 4: Evaluating what's been learned

36

A hypothetical lift chart



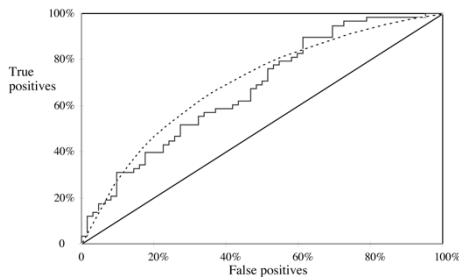
2016 Lecture 4: Evaluating what's been learned 37

ROC curves

- *ROC curves* are similar to lift charts
 - ♦ Stands for “receiver operating characteristic”
 - ♦ Used in signal detection to show tradeoff between hit rate and false alarm rate over noisy channel
- Differences to lift chart:
 - ♦ y axis shows percentage of true positives in sample *rather than absolute number*
 - ♦ x axis shows percentage of false positives in sample *rather than sample size*

2016 Lecture 4: Evaluating what's been learned 38

A sample ROC curve



- Jagged curve—one set of test data
- Smooth curve—use cross-validation

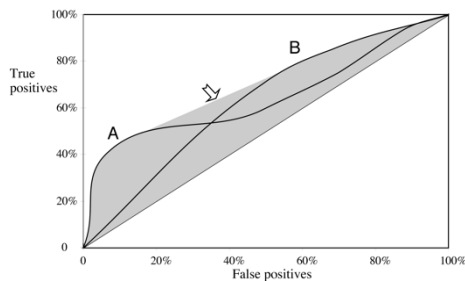
2016 Lecture 4: Evaluating what's been learned 39

Cross-validation and ROC curves

- Simple method of getting a ROC curve using cross-validation:
 - ♦ Collect probabilities for instances in test folds
 - ♦ Sort instances according to probabilities
- This method is implemented in WEKA
- However, this is just one possibility
 - ♦ Another possibility is to generate an ROC curve for each fold and average them

2016 Lecture 4: Evaluating what's been learned 40

ROC curves for two schemes



- For a small, focused sample, use method A
- For a larger one, use method B
- In between, choose between A and B with appropriate probabilities

2016 Lecture 4: Evaluating what's been learned 41

The convex hull

- Given two learning schemes we can achieve any point on the convex hull!
- TP and FP rates for scheme 1: t_1 and f_1
- TP and FP rates for scheme 2: t_2 and f_2
- If scheme 1 is used to predict $100 \times q$ % of the cases and scheme 2 for the rest, then
 - TP rate for combined scheme:
 $q \times t_1 + (1-q) \times t_2$
 - FP rate for combined scheme:
 $q \times f_1 + (1-q) \times f_2$

2016 Lecture 4: Evaluating what's been learned 42

Recall-Precision Curve

- Precision/recall curves have hyperbolic shape
- Summary measures: average precision at 20%, 50% and 80% recall (*three-point average recall*)

2016 Lecture 4: Evaluating what's been learned 43

Some measures and domains

	Domain	Plot	Explanation
Lift chart	Marketing	TP Subset size	$\frac{TP}{(TP+FP)/(TP+FP+TN+FN)}$
ROC curve	Communications	TP rate FP rate	$\frac{TP}{(TP+FN)}$ $\frac{FP}{(FP+TN)}$
Recall-precision curve	Information retrieval	Recall Precision	$\frac{TP}{(TP+FN)}$ $\frac{TP}{(TP+FP)}$

2016 Lecture 4: Evaluating what's been learned 44

Summary


- Percentage of retrieved documents that are relevant:
precision = $TP / (TP + FP)$
- Percentage of relevant documents that are returned:
recall = $TP / (TP + FN)$
- Precision/recall curves have hyperbolic shape
- Summary measures: average precision at 20%, 50% and 80% recall (*three-point average recall*)
- *F-measure* = $(2 \times \text{recall} \times \text{precision}) / (\text{recall} + \text{precision})$
- *sensitivity* \times *specificity* = $(TP / (TP + FN)) \times (TN / (FP + TN))$
- Area under the ROC curve (*AUC*): probability that randomly chosen positive instance is ranked above randomly chosen negative one

2016 Lecture 4: Evaluating what's been learned 45

Model selection criteria

- Model selection criteria attempt to find a good compromise between:
 - The complexity of a model
 - Its prediction accuracy on the training data
- Reasoning: a good model is a simple model that achieves high accuracy on the given data
- Also known as *Occam's Razor*: the best theory is the smallest one that describes all the facts

William of Ockham, born in the village of Ockham in Surrey (England) about 1285, was the most influential philosopher of the 14th century and a controversial theologian.



2016 Lecture 4: Evaluating what's been learned 46

Elegance vs. errors

- Theory 1: very simple, elegant theory that explains the data almost perfectly
- Theory 2: significantly more complex theory that reproduces the data without mistakes
- Theory 1 is probably preferable

2016 Lecture 4: Evaluating what's been learned 47

Simplicity first

- Simple algorithms often work very well!
- There are many kinds of simple structure, eg:
 - One attribute does all the work
 - All attributes contribute equally & independently
 - A weighted linear combination might do
 - Instance-based: use a few prototypes
 - Use simple logical rules
- Success of method depends on the domain

2016 Lecture 4: Evaluating what's been learned 48

The end