

Comp 135 Introduction to Machine Learning and Data Mining

Fall 2016

Professor: Roni Khardon

Computer Science
Tufts University

Evaluating machine learning outcomes

- How should we define a good outcome for a machine learning algorithm?
- How do we know when an algorithm is doing well in a particular application?
- How can we compare different classifiers?
- How can we compare different learning algorithms?

comp135

Roni Khardon, Tufts University

Evaluating machine learning outcomes

- More concrete questions:
- What quantity should we measure?
- How can we estimate it?
- Can we get quantitative guarantees for estimates and comparisons?

comp135

Roni Khardon, Tufts University

What to Measure?

- We have so far focused on classification and looked at accuracy
- But this depends on the application and data distribution ...

comp135

Roni Khardon, Tufts University

Confusion Matrix

- + - <-- classified as
- TP FN | true label +
- FP TN | true label -

$$\text{Acc} = \frac{TP + TN}{TP + TN + FP + FN}$$

- In many applications one class (typically the +) has very low frequency
- Running algorithm to optimize Acc often yields bad results. **Why?**

comp135

Roni Khardon, Tufts University

IR Terminology

- Context: search for items with **QueryTerm**
- De-emphasize role of Neg examples
- Aims to measure quality of response

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN}$$

- Tradeoff given by Precision/Recall curve
- And by $F = \frac{2 \cdot R \cdot P}{R + P}$

comp135

Roni Khardon, Tufts University

Medical Community Terminology

- Context: test to identify Med Condition

$$\text{Sensitivity} = \text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

- These measure "accuracy" within each class

comp135

Roni Khardon, Tufts University

Signal Detection Terminology

- Context: identify "signal"

$$\text{TPrate} = \text{Recall} = \frac{TP}{TP + FN}$$

$$\text{FPrate} = 1 - \text{Specificity} = \frac{FP}{TN + FP}$$

- The ROC curve (receiver operator characteristic) plots TPrate vs FPrate

comp135

Roni Khardon, Tufts University

ROC Curves

- The AROC (area under ROC curve) has a nice interpretation:

AROC = probability of correctly ranking a random Pos example above a random Neg example

comp135

Roni Khardon, Tufts University

From Ranking to ROC Curve

- Many learning algorithms (e.g., Naïve Bayes, Perceptron) provide a numerical output that can be used to rank examples in addition to the prediction of +/- label
- This can be used to produce a ROC curve for the algorithm by changing its threshold specifying the transition from Neg to Pos.

comp135

Roni Khardon, Tufts University

How to Measure?

- Validation Set Method:
 - keep aside a portion of the example set
 - Train model on remaining data
 - Measure performance on validation set
- (+) Unbiased estimate of quantity
- (-) Wastes data ...
- (-) Variance in estimate due to choice of validation set. *Can we fix this?*

comp135

Roni Khardon, Tufts University

How to Measure?

- Validation Set Method:
 - keep aside a portion of the example set
 - Train model on remaining data
 - Measure performance on validation set
- Can reduce variance by repeating k times and averaging
- But this introduces bias in estimates because the train/test data in different runs are highly correlated

comp135

Roni Khardon, Tufts University

How to Measure?

- Cross Validation Method:
- Divide data into k portions (called folds)
- Repeat k times
 - Train model on data from all folds except k
 - Measure performance on k 'th fold
- Average the performance
- Test data in different folds is disjoint
- But training sets are not

comp135

Roni Khardon, Tufts University

Stratified Cross Validation

- Another source of variance in cross validation is the fact that the test sets across folds do not have exactly the same class frequencies
- Stratified cross validation removes this variation by splitting into random subsets of the same size per class

comp135

Roni Khardon, Tufts University

Implementing Stratified Cross Validation

- Split data according to class
- Randomly permute the examples in each class
- Partition each class (by permuted index) into k equal size portions (up to ± 1)
- Join portions i from all classes to get fold i
- This gives k portions of the same size and having the same class proportions

comp135

Roni Khardon, Tufts University

Leave One Out Method

- Setting $k=N$ (number of examples) we get the leave one out method
- Often effective but
- High variance in individual estimates
- Expensive if we need to train N times
- But for some algorithms, e.g. kNN, can get away without significant expense.

comp135

Roni Khardon, Tufts University

Model Selection

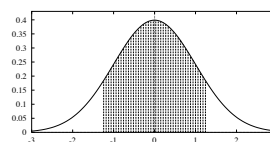
- Cross validation can be applied to configure parameters in many algorithms during the training phase.
 - Pick k for nearest neighbors
 - Pick prune threshold in DT
- In this case we perform cross validation on training set only
- Using the result we pick best parameter value
- Finally: we train again on the entire data

comp135

Roni Khardon, Tufts University

Quantitative Comparisons

Review notions from Normal distributions



Beware to check if using tables for one-sided bounds or two sided bounds. These slides use two sided bounds as in the shaded area. Some other sources use one sided bounds. It is straightforward to extract one from the other.

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

Image from Mitchell text/slides

comp135

Roni Khardon, Tufts University

Quantitative Comparisons

- From probability to estimates

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

- Probability Statement

[with probability $N\%$] $x \in \mu \pm z_N \sigma$

- Confidence interval

[with confidence $N\%$] $\mu \in x \pm z_N \sigma$

comp135

Roni Khardon, Tufts University

Confidence Interval

[with confidence $N\%$] $\mu \in x \pm z_N \sigma$

- Assume $x \sim \mathcal{N}(\mu, \sigma^2)$

- And we sample to observe value of x

[with confidence 95%] $\mu \in x \pm 1.96\sigma$

comp135

Roni Khardon, Tufts University

Evaluating one classifier

- We measure the performance on test set of n examples
- If the error rate is p then the number of mistakes on the test set is distributed as Binomial(n, p)
- Our estimate is $\hat{p} = \frac{\# \text{ mistakes}}{n}$
- \hat{p} distributed approx as $\hat{p} \sim \mathcal{N}(p, \frac{p(1-p)}{n})$

comp135

Roni Khardon, Tufts University

Evaluating one classifier

- \hat{p} distributed approx as $\hat{p} \sim \mathcal{N}(p, \frac{p(1-p)}{n})$

- Therefore, we can apply

[with confidence $N\%$] $\mu \in x \pm z_N \sigma$

- To get

[with confidence $N\%$] $p \in \hat{p} \pm z_N \sqrt{\frac{p(1-p)}{n}}$

comp135

Roni Khardon, Tufts University

Evaluating one classifier

- **Good news:** we have an interval for p

[with confidence $N\%$] $p \in \hat{p} \pm z_N \sqrt{\frac{p(1-p)}{n}}$

- **But:** we cannot plug in p in variance part

- Solution 1: $p(1-p) \leq 0.25$

[with confidence $N\%$] $p \in \hat{p} \pm z_N \sqrt{\frac{1}{4n}}$

- Solution 2: (not justified)

[with confidence $N\%$] $p \in \hat{p} \pm z_N \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

comp135

Roni Khardon, Tufts University

Evaluating Learning Algorithm

- Cross validation error estimate $\hat{e} = \frac{1}{k} \sum e_i$
- Define e as the expected error rate when running alg (on data of this size)

- Then we have: $e_i \sim \mathcal{N}(e, \sigma_e^2)$

$\hat{e} \sim \mathcal{N}(e, \sigma_e^2/k)$

- And we can use

[with confidence $N\%$] $\mu \in x \pm z_N \sigma$

comp135

Roni Khardon, Tufts University

Evaluating Learning Algorithm

- Cross validation error estimate $\hat{e} = \frac{1}{k} \sum e_i$
- Define e as the expected error rate when running alg (on data of this size)
- Then we have:

$$e_i \sim \mathcal{N}(e, \sigma_e^2)$$

$$\hat{e} \sim \mathcal{N}(e, \sigma_e^2/k)$$

- And we can use

$$[\text{with confidence } N\%] \quad e \in \hat{e} \pm z_N \sqrt{\sigma_e^2/k}$$

comp135

Roni Khardon, Tufts University

Evaluating Learning Algorithm

$$[\text{with confidence } N\%] \quad e \in \hat{e} \pm z_N \sqrt{\sigma_e^2/k}$$

- **But:** we do not know σ_e
- Solution 1: plug in the estimate s instead of σ_e
- Solution 2: a more accurate interval through T random variable

$$s = \sqrt{\frac{1}{k-1} \sum (e_i - \hat{e})^2}$$

Sol2 is more accurate and as easy to implement so avoid Sol1.

comp135

Roni Khardon, Tufts University

Evaluating Learning Algorithm

- It turns out [we skip definitions and details] that $\frac{\hat{e} - e}{(s/\sqrt{k})} \sim T_{k-1}$

- And that this implies

$$[\text{with confidence } N\%] \quad e \in \hat{e} \pm t_{N,k-1} \frac{s}{\sqrt{k}}$$

- And more concretely

$$[\text{with confidence } N\%] \quad e \in \hat{e} \pm t_{N,k-1} \sqrt{\frac{\sum (e_i - \hat{e})^2}{(k-1)k}}$$

comp135

Roni Khardon, Tufts University

Degrees of freedom	p - 0.1	p - 0.05	p - 0.02	p - 0.01	p - 0.002	p - 0.001
1	6.314	12.706	31.821	63.657	318.310	636.620
2	2.920	4.303	6.965	9.925	22.327	31.598
3	2.353	3.182	4.541	5.841	16.274	12.924
4	2.132	2.776	3.747	4.604	11.715	8.610
5	2.015	2.571	3.365	4.032	9.893	6.965
6	1.943	2.447	3.143	3.707	9.208	5.959
7	1.895	2.365	2.998	3.499	8.785	5.408
8	1.860	2.306	2.896	3.355	8.451	5.041
9	1.833	2.262	2.821	3.250	8.207	4.761
10	1.812	2.228	2.764	3.169	8.044	4.587
11	1.796	2.201	2.718	3.106	7.895	4.437
12	1.782	2.179	2.681	3.055	7.759	4.318
13	1.771	2.160	2.650	3.012	7.632	4.221
14	1.761	2.145	2.624	2.977	7.517	4.140
15	1.753	2.131	2.602	2.947	7.433	4.073
16	1.746	2.120	2.583	2.921	7.366	4.015
17	1.740	2.110	2.567	2.898	7.305	3.965
18	1.734	2.101	2.552	2.878	7.250	3.922
19	1.729	2.093	2.539	2.861	7.200	3.883
20	1.725	2.086	2.528	2.845	7.155	3.850
21	1.721	2.080	2.518	2.831	7.114	3.819
22	1.717	2.074	2.508	2.818	7.076	3.792
23	1.714	2.069	2.500	2.807	7.041	3.767
24	1.711	2.064	2.492	2.797	7.007	3.745
25	1.708	2.060	2.485	2.787	6.975	3.725
26	1.706	2.056	2.479	2.779	6.945	3.707
27	1.703	2.052	2.473	2.771	6.917	3.690
28	1.701	2.048	2.467	2.763	6.890	3.674
29	1.699	2.045	2.462	2.756	6.865	3.659
30	1.697	2.042	2.457	2.750	6.841	3.646
40	1.684	2.021	2.423	2.704	6.576	3.551
60	1.671	2.000	2.380	2.660	6.389	3.490
120	1.658	1.980	2.358	2.617	6.160	3.373
∞	1.645	1.960	2.326	2.576	5.999	3.291

Table from
<http://www.saburchill.com/>

Comparing Two Algorithms

- We can calculate an interval for the error of each and say that one is better if the intervals do not overlap
- But we can do much better when variance σ_e is large
- Alg1 Alg2 Diff
- 60 68 8
- 70 77 7
- 80 88 8

comp135

Roni Khardon, Tufts University

Comparing Two Algorithms

- We can calculate an interval for the error of each and say that one is better if the intervals do not overlap
- But we can do much better when variance σ_e is large
- Compute a confidence interval the difference in performance.
- Matched tests: estimates come from running algorithms on the same folds. This is valid and reduces variance.

comp135

Roni Khardon, Tufts University

Running many experiments and tests

- We develop 100 new algorithms
- Calculate interval with N=95% confidence for diff over baseline alg
- What is the probability of wrongly claiming a significant difference in at least one of these?

$$1 - 0.95^{100} = 1 - 0.059 = 0.9941$$

comp135

Roni Khardon, Tufts University

Running many experiments and tests

- We develop 100 new algorithms
- Calculate interval with N=95% confidence for diff over baseline alg
- Simple correction (tests are not required to be independent)
- To get 0.95 confidence for all the tests together we need each individual run with N=99.95%

$$p(\text{fail}) \leq 100(1 - 0.9995) = 100 * 0.0005 = 0.05$$

$$p(\text{not fail}) \geq 0.95$$

comp135

Roni Khardon, Tufts University

Evaluating machine learning outcomes

- You should now have some answers to these questions:
- What quantity should we measure?
- How can we estimate it?
- Can we get quantitative guarantees for estimates and comparisons?

comp135

Roni Khardon, Tufts University