

CS 4260 and CS 5260

Vanderbilt University

Lecture on Evaluation

This lecture assumes that you have

- Read Section 7.1 through 7.2 of *ArtInt* and

ArtInt: Poole and Mackworth, *Artificial Intelligence 2E*

at <http://artint.info/2e/html/ArtInt2e.html>

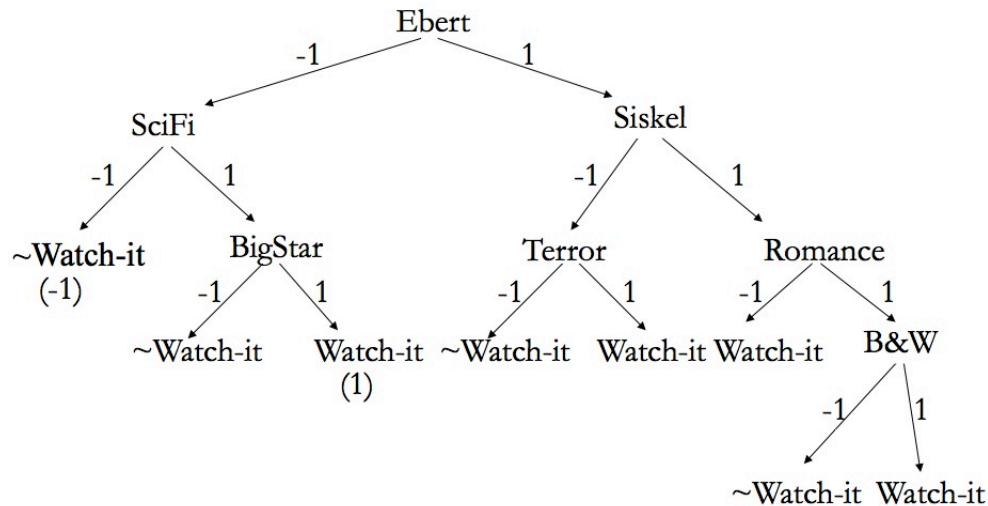
to include slides at <http://artint.info/2e/slides/ch04/lect1.pdf>

We have seen two supervised machine learning strategies

- Naïve Bayesian learning (was optional during ML week; required later)

$$P(C=c1) * P(\text{SciFi} = -1 | c1) * P(\text{Terror} = 1 | c1) * P(\text{Romance} = -1 | c1) * P(\text{Ebert} = 1 | c1) * P(\text{Siskel} = 1 | c1) * \dots]$$

- Decision tree learning



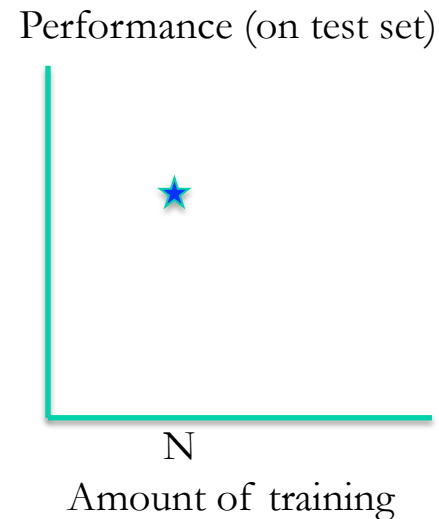
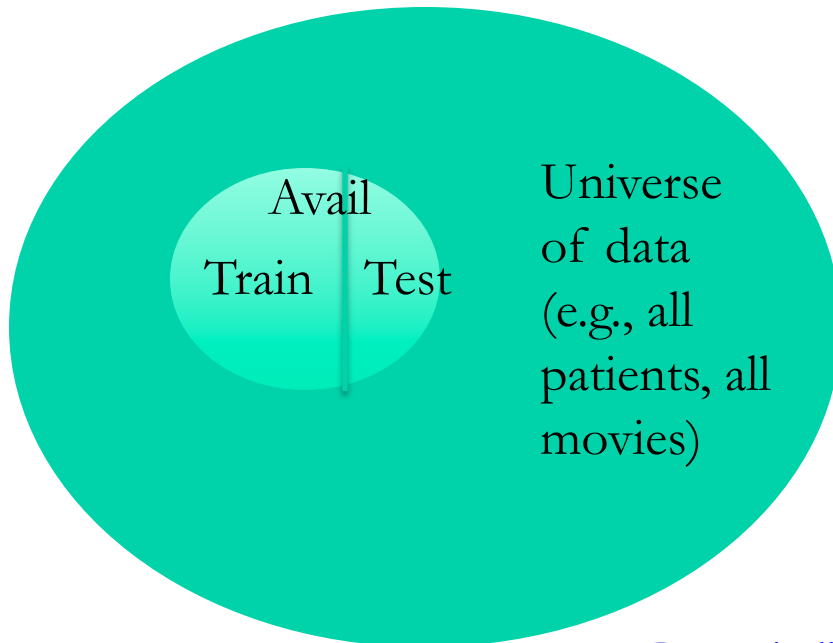
How can we compare them?

How can we characterize their learning independent of each other?

How do we parameterize each of them to maximize performance?

Training and Test Data Sets

- Testing a classifier/predictor on data that was used for training is overly optimistic,
 - even if the method doesn't memorize each data per se
- More realistic, in most cases, is to test on previously unseen data
- What does this tell us?
- If there are N training data, then test set accuracy (or error) approximates (to an unknown extent) the performance of classifiers constructed by the learning method on N training data

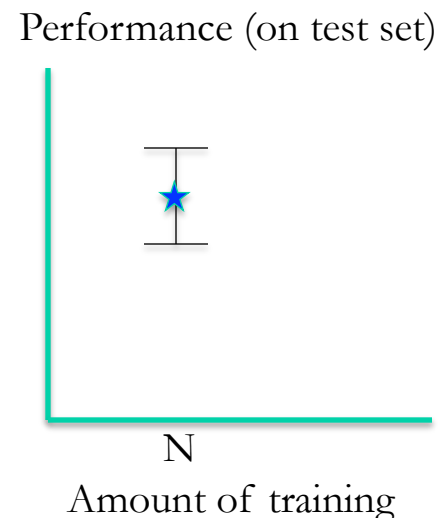


But typically, most interested in performance on universe

Using repeated Train and Test splits

Given: M data available, Avail
Learning Trials, L
Training set size, N
Test set size, M-N

Local: Training Set, Train
Test Set, Test
Classifier
AggregatedPerformance (e.g., Mean, Median, Mode)



Initialize AggregatedPerformance

Repeat L times {

Train \leftarrow Randomly draw N data from Avail, “without replacement”

Test \leftarrow Avail – Train

Classifier \leftarrow Learn(Train)

AggregatedPerformance \leftarrow Performance(Classifier, Test) + AggregatedPerformance

}

Return AggregatedPerformance

This provides approximation of performance (on Universe)
of learning method at training size of N

Generating Learning Curves through repeated Train and Test splits

Given: M data available, Avail
 Learning Trials, L
 Training set sizes, $N_1 \dots N_{\max}$
 Test set size, $M - N_{\max}$

Local: Training Set, Train
 Test Set, Test
 Classifier
 AggregatedPerformanceVector

Initialize AggregatedPerformanceVector

Repeat L times {

Train \leftarrow Randomly draw N_{\max} data from Avail,
 “without replacement”

Test \leftarrow Avail – Train

Partition Train into 1 to max bins, TrainBin₁ through TrainBin_{max}

For k = 1 to max {

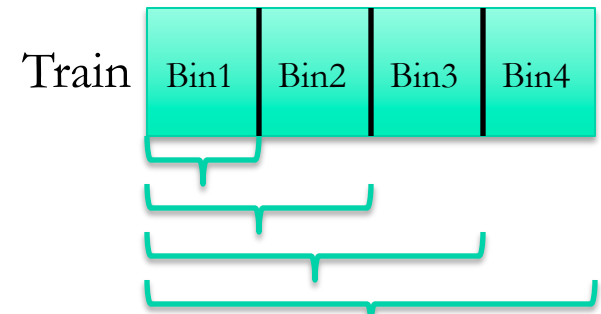
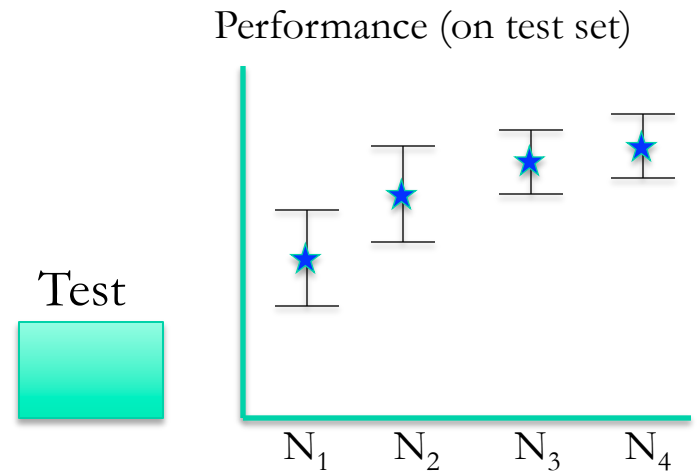
Classifier \leftarrow Learn(Union of TrainBin₁ through TrainBin_k)

AggregatedPerformanceVector[k]

\leftarrow Performance(Classifier, Test) + AggregatedPerformanceVector[k]

}

Return AggregatedPerformanceVector



Comparing Learning Curves of Different Supervised Machine Learning Methods

Given: ...
Local: ...

Initialize *AggregatedPerformanceVectorMethod1* (e.g., TDIDT)

Initialize *AggregatedPerformanceVectorMethod2* (e.g., NBC)

Repeat L times {

Train \leftarrow Randomly draw N_{\max} data from Avail,
“without replacement”

Test \leftarrow Avail – Train

Partition Train into 1-to-max bins, TrainBin₁ through TrainBin_{max}

For k = 1 to max {

Classifier1 \leftarrow Learn1(Union of TrainBin₁ through TrainBin_k)

AggregatedPerformanceVectorMethod1[k]

\leftarrow Performance(Classifier1, Test) + AggregatedPerformanceVector1[k]

Classifier2 \leftarrow Learn2(Union of TrainBin₁ through TrainBin_k)

AggregatedPerformanceVectorMethod2[k]

\leftarrow Performance(Classifier2, Test) + AggregatedPerformanceVector2[k]

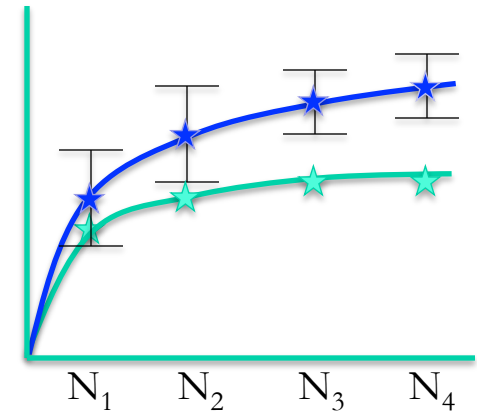
}

Return *AggregatedPerformanceVectorMethod1*, *AggregatedPerformanceVectorMethod2*

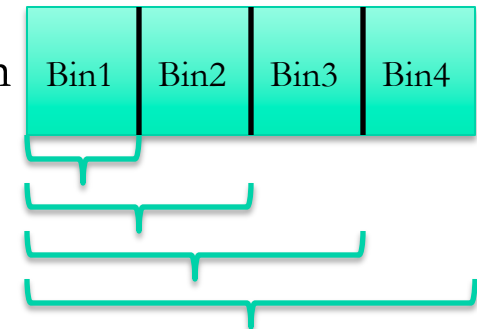
Test



Performance (on test set)



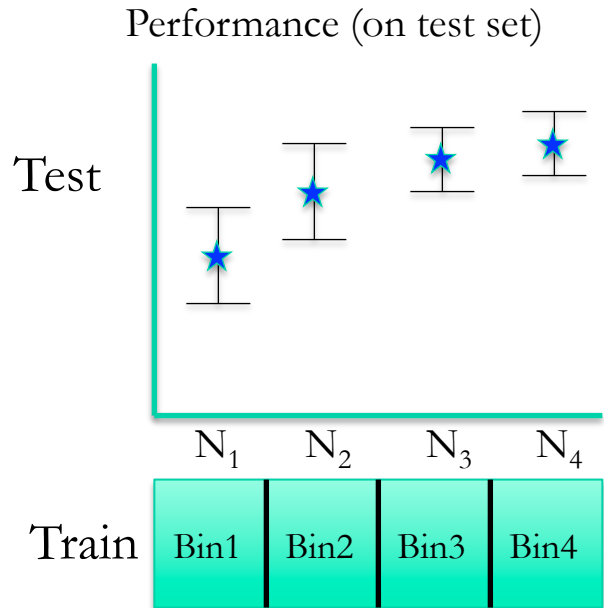
Train



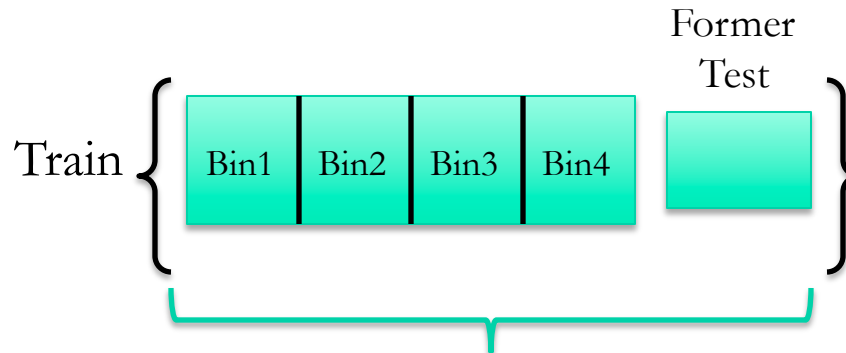
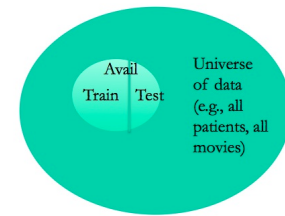
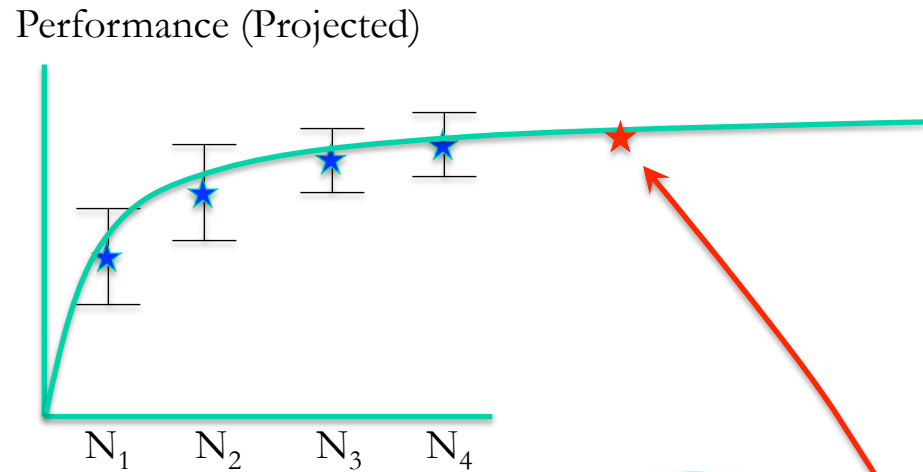
This provides approximations of performance (on Universe)
of each method at different training N

How Might we use in Real Setting

1. Use Training/Test splits to plot performance



2. Curve fit learning behavior to project performance (on universe) at larger N



and deploy classifier

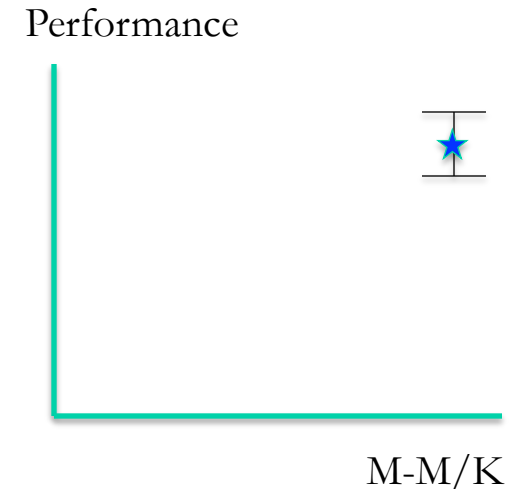
All available data

K-Fold Cross Validation

1. Randomize order of available M data
2. divide available data into K (e.g., 10) equal size bins or folds
3. For $I = 1$ to K {
 - Train on union of all folds, except fold $_I$
 - Test on fold $_I$
4. Average results



All available data



M-Fold Cross Validation (or leave-one-out cross validation)

Divide a data set of size M into M singleton folds, and follow algorithm above (e.g., for each datum, train on $M-1$ other data and test on the datum)

This is often regarded as the best way to leverage the existing data and get as close as one can to estimating performance on a final deployed classifier trained on all data