
COMP 551 – Applied Machine Learning

Lecture 6: Performance evaluation. Model assessment and selection.

Instructor: Herke van Hoof (herke.vanhoof@mail.mcgill.ca)

Slides mostly by: Joelle Pineau

Class web page: www.cs.mcgill.ca/~hvanho2/comp551

Unless otherwise noted, all material posted for this course are copyright of the instructors, and cannot be reused or reposted without the instructors' written permission.

Today's quiz (on myCourses)

- Quiz on classification on myCourses

Project questions

- Best place to ask questions: MyCourses forum
- Others can browse questions/answers so everyone can learn from them
- If you have a specific problem, try to visit the office hour of the responsible TA (mentioned on exercise) – they are best placed to help you!

Project 1 hand in

- Original data: Jan 26
- We'll accept submissions until Jan 29, noon (strict deadline)
 - Hardcopy (in box) & code/data (on mycourses)
- Late policy: within 1 week late will be accepted with 30% penalty
- Caution: project 2 will still be available from Jan 26!
- Hand-in box:
 - Opposite 317 in McConnell building

Evaluating performance

- Different objectives:
 - Selecting the right model for a problem.
 - Testing performance of a new algorithm.
 - Evaluating impact on a new application.

Performance metrics for classification

- Not all errors have equal impact!
- There are different types of mistakes, particularly in the **classification setting**.

Example 1

Int J Neural Syst. 2013 Oct;23(5):1350023. doi: 10.1142/S0129065713500238. Epub 2013 Jul 21.

Application of intrinsic time-scale decomposition (ITD) to EEG signals for automated seizure prediction.

Martis RJ¹, Acharya UR, Tan JH, Petznick A, Tong L, Chua CK, Ng EY.

Author information

Abstract

Intrinsic time-scale decomposition (ITD) is a new nonlinear method of time-frequency representation which can decipher the minute changes in the nonlinear EEG signals. In this work, we have automatically classified normal, interictal and ictal EEG signals using the features derived from the ITD representation. The energy, fractal dimension and sample entropy features computed on ITD representation coupled with decision tree classifier has yielded an average classification accuracy of 95.67%, sensitivity and specificity of 99% and 99.5%, respectively using 10-fold cross validation scheme. With application of the nonlinear ITD representation, along with conceptual advancement and improvement of the accuracy, the developed system is clinically ready for mass screening in resource constrained and emerging economy scenarios.

Example 1

Int J Neural Syst. 2013 Oct;23(5):1350023. doi: 10.1142/S0129065713500238. Epub 2013 Jul 21.

Application of intrinsic time-scale decomposition (ITD) to EEG signals for automated seizure prediction.

Martis RJ¹, Acharya UR, Tan JH, Petznick A, Tong L, Chua CK, Ng EY.

Author information

Abstract

Intrinsic time-scale decomposition (ITD) is a new nonlinear method of time-frequency representation which can decipher the minute changes in the nonlinear EEG signals. In this work, we have automatically classified normal, interictal and ictal EEG signals using the features derived from the ITD representation. The energy, fractal dimension and sample entropy features computed on ITD representation coupled with decision tree classifier has yielded an average classification accuracy of 95.67%, sensitivity and specificity of 99% and 99.5%, respectively using 10-fold cross validation scheme. With application of the nonlinear ITD representation, along with conceptual advancement and improvement of the accuracy, the developed system is clinically ready for mass screening in resource constrained and emerging economy scenarios.

Why not just report classification accuracy?

Performance metrics for classification

- Not all errors have equal impact!
- There are different types of mistakes, particularly in the **classification setting**.
 - E.g. Consider the diagnostic of a disease. Two types of mis-diagnostics:
 - Patient does not have disease but received positive diagnostic (**Type I error**);
 - Patient has disease but it was not detected (**Type II error**).

Performance metrics for classification

- Not all errors have equal impact!
- There are different types of mistakes, particularly in the **classification setting**.
 - E.g. Consider the diagnostic of a disease. Two types of mis-diagnostics:
 - Patient does not have disease but received positive diagnostic (**Type I error**);
 - Patient has disease but it was not detected (**Type II error**).
 - E.g. Consider the problem of spam classification:
 - A message that is not spam is assigned to the spam folder (**Type I error**);
 - A message that is spam appears in the regular folder (**Type II error**).

Performance metrics for classification

- Not all errors have equal impact!
- There are different types of mistakes, particularly in the **classification setting**.
 - E.g. Consider the diagnostic of a disease. Two types of mis-diagnostics:
 - Patient does not have disease but received positive diagnostic (**Type I error**);
 - Patient has disease but it was not detected (**Type II error**).
 - E.g. Consider the problem of spam classification:
 - A message that is not spam is assigned to the spam folder (**Type I error**);
 - A message that is spam appears in the regular folder (**Type II error**).
- How many **Type I errors** are you willing to tolerate, for a reasonable rate of **Type II errors** ?

Example 2

Clin Neurophysiol. 2013 Sep;124(9):1745-54. doi: 10.1016/j.clinph.2013.04.006. Epub 2013 May 3.

Seizure prediction in patients with mesial temporal lobe epilepsy using EEG measures of state similarity.

Gadhoumi K¹, Lina JM, Gotman J.

+ Author information

Abstract

OBJECTIVES: In patients with intractable epilepsy, predicting seizures above chance and with clinically acceptable performance has yet to be demonstrated. In this study, an intracranial EEG-based seizure prediction method using measures of similarity with a reference state is proposed.

METHODS: 1565 h of continuous intracranial EEG data from 17 patients with mesial temporal lobe epilepsy were investigated. The recordings included 175 seizures. In each patient the data was split into a training set and a testing set. EEG segments were analyzed using continuous wavelet transform. During training, a reference state was defined in the immediate preictal data and used to derive three features quantifying the discrimination between preictal and interictal states. A classifier was then trained in the feature space. Its performance was assessed using testing set and compared with a random predictor for statistical validation.

RESULTS: Better than random prediction performance was achieved in 7 patients. The sensitivity was higher than 85%, the warning rate was less than 0.35/h and the proportion of time under warning was less than 30%.

CONCLUSION: Seizures are predicted above chance in 41% of patients using measures of state similarity.

SIGNIFICANCE: Sensitivity and specificity levels are potentially interesting for closed-loop seizure control applications.

Example 3

Conf Proc IEEE Eng Med Biol Soc. 2012;2012:1061-4. doi: 10.1109/EMBC.2012.6346117.

Low complexity algorithm for seizure prediction using Adaboost.

Ayinala M¹, Parhi KK.

Author information

Abstract

This paper presents a novel low-complexity patient-specific algorithm for seizure prediction. Adaboost algorithm is used in two stages of the algorithm: feature selection and classification. The algorithm extracts spectral power features in 9 different sub-bands from the electroencephalogram (EEG) recordings. We have proposed a new feature ranking method to rank the features. The key (top ranked) features are used to make a prediction on the seizure event. Further, to reduce the complexity of classification stage, a non-linear classifier is built based on the Adaboost algorithm using decision stumps (linear classifier) as the base classifier. The proposed algorithm achieves a sensitivity of 94.375% for a total of 71 seizure events with a low false alarm rate of 0.13 per hour and 6.5% of time spent in false alarms using an average of 5 features for the Freiburg database. The low computational complexity of the proposed algorithm makes it suitable for an implantable device.

PMID: 23366078 [PubMed - indexed for MEDLINE]

Terminology

- Type of classification outputs:
 - True positive (m_{11}): Example of class 1 predicted as class 1.
 - False positive (m_{01}): Example of class 0 predicted as class 1. Type 1 error.
 - True negative (m_{00}): Example of class 0 predicted as class 0.
 - False negative (m_{10}): Example of class 1 predicted as class 0. Type II error.
- Total number of instances: $m = m_{00} + m_{01} + m_{10} + m_{11}$

Terminology

- Type of classification outputs:
 - **True positive** (m11): Example of class 1 predicted as class 1.
 - **False positive** (m01): Example of class 0 predicted as class 1. **Type 1 error.**
 - **True negative** (m00): Example of class 0 predicted as class 0.
 - **False negative** (m10): Example of class 1 predicted as class 0. **Type II error.**
- Total number of instances: $m = m00 + m01 + m10 + m11$
- **Error rate:** $(m01 + m10) / m$
 - If the classes are imbalanced (e.g. 10% from class 1, 90% from class 0), one can achieve low error (e.g. 10%) by classifying everything as coming from class 0!

Confusion matrix

- Many software packages output this matrix.

$$\begin{bmatrix} m_{00} & m_{01} \\ m_{10} & m_{11} \end{bmatrix}$$

Confusion matrix

- Many software packages output this matrix.

$$\begin{bmatrix} m_{00} & m_{01} \\ m_{10} & m_{11} \end{bmatrix}$$

- Be careful! Sometimes the format is slightly different

(E.g. http://en.wikipedia.org/wiki/Precision_and_recall#Definition_.28classification_context.29)

predicted class (expectation)	actual class (observation)	
	tp	fp
	(true positive) Correct result	(false positive) Unexpected result
	fn	tn
	(false negative) Missing result	(true negative) Correct absence of result

Common measures

- **Accuracy** $= (TP + TN) / (TP + FP + FN + TN)$
- **Precision** $= \text{True positives} / \text{Total number of declared positives}$
 $= TP / (TP + FP)$
- **Recall** $= \text{True positives} / \text{Total number of actual positives}$
 $= TP / (TP + FN)$

Common measures

- Accuracy = $(TP + TN) / (TP + FP + FN + TN)$
- Precision = True positives / Total number of declared positives
= $TP / (TP + FP)$
- Recall = True positives / Total number of actual positives
= $TP / (TP + FN)$
- Sensitivity is the same as recall.
- Specificity = True negatives / Total number of actual negatives
= $TN / (FP + TN)$

Text
classification



Medicine



Common measures

- Accuracy = $(TP + TN) / (TP + FP + FN + TN)$
- Precision = True positives / Total number of declared positives
= $TP / (TP + FP)$
- Recall = True positives / Total number of actual positives
= $TP / (TP + FN)$
- Sensitivity is the same as recall.
- Specificity = True negatives / Total number of actual negatives
= $TN / (FP + TN)$
- False positive rate = $FP / (FP + TN)$ (= 1-specificity)

Text
classification



Medicine



Common measures

- Accuracy = $(TP + TN) / (TP + FP + FN + TN)$
- Precision = True positives / Total number of declared positives
= $TP / (TP + FP)$
- Recall = True positives / Total number of actual positives
= $TP / (TP + FN)$
- Sensitivity is the same as recall.
- Specificity = True negatives / Total number of actual negatives
= $TN / (FP + TN)$
- False positive rate = $FP / (FP + TN)$ (= 1-specificity)
- F1 measure
$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Text
classification



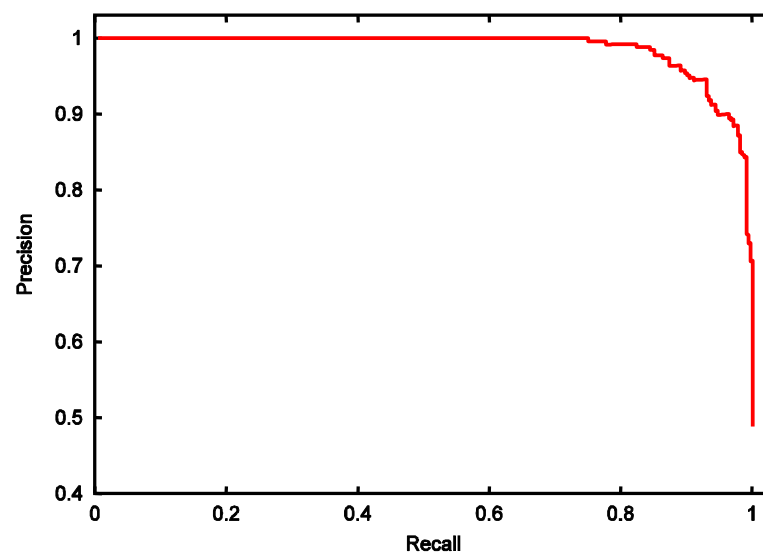
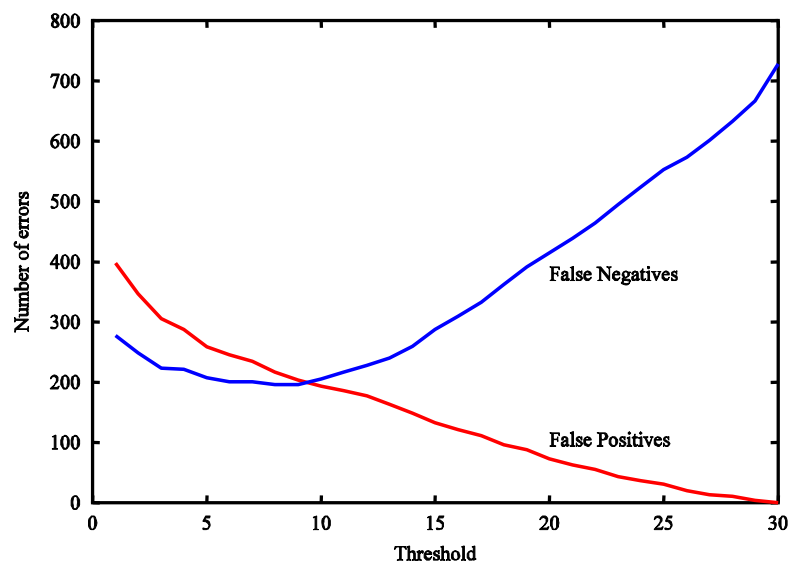
Medicine



Trade-off

- Often have a trade-off between false positives and false negatives.

E.g. Consider 30 different classifiers trained on a class. Classify a new sample as positive if K classifiers output positive. Vary K between 0 and 30.



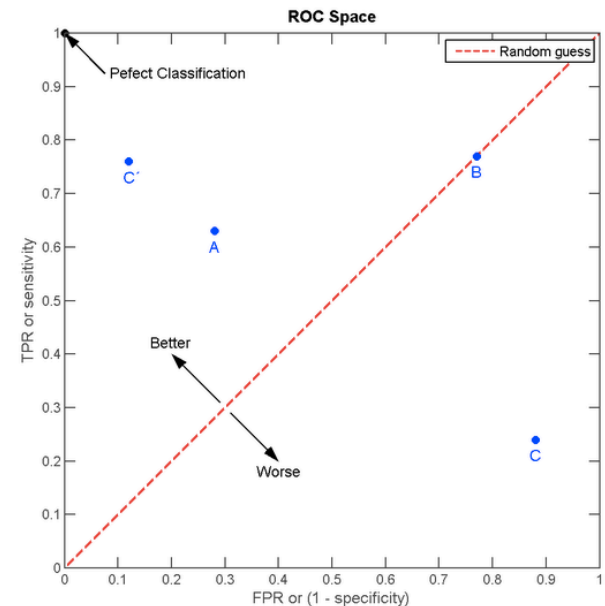
Receiver-operator characteristic (ROC) curve

- Characterizes the performance of a binary classifier over a range of classification thresholds

Data from 4 prediction results:

A			B			C			C'		
TP=63	FP=28	91	TP=77	FP=77	154	TP=24	FP=88	112	TP=76	FP=12	88
FN=37	TN=72	109	FN=23	TN=23	46	FN=76	TN=12	88	FN=24	TN=88	112
100	100	200	100	100	200	100	100	200	100	100	200
TPR = 0.63			TPR = 0.77			TPR = 0.24			TPR = 0.76		
FPR = 0.28			FPR = 0.77			FPR = 0.88			FPR = 0.12		
PPV = 0.69			PPV = 0.50			PPV = 0.21			PPV = 0.86		
F1 = 0.66			F1 = 0.61			F1 = 0.22			F1 = 0.81		
ACC = 0.68			ACC = 0.50			ACC = 0.18			ACC = 0.82		

ROC curve:

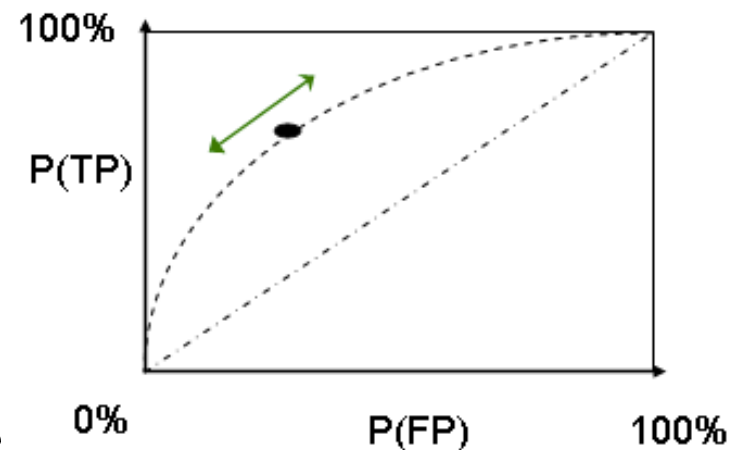
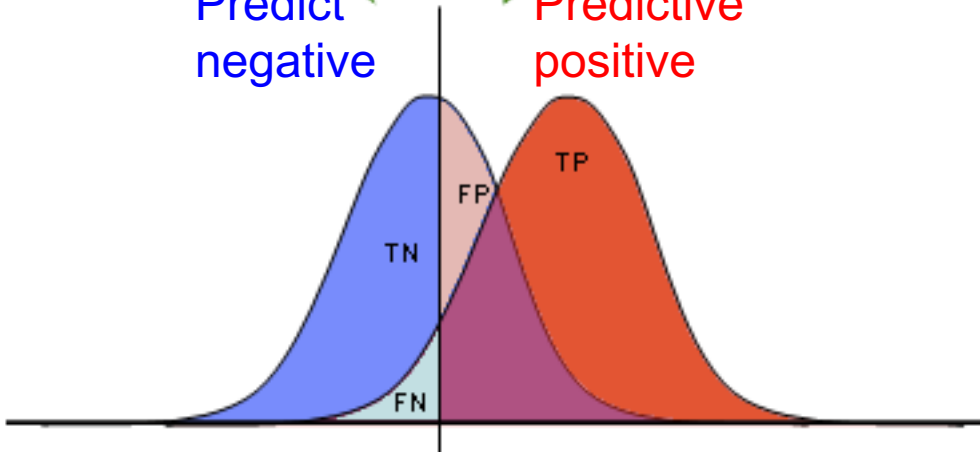


Example from: http://en.wikipedia.org/wiki/Receiver_operating_characteristic

Understanding the ROC curve

- Consider a classification problem where data is generated by 2 Gaussians (blue = negative class; red = positive class).
- Consider the decision boundary (shown as a vertical line on the left figure), where you predict Negative on the left of the boundary and predict Positive on the right of the boundary.
- Changing that boundary defines the ROC curve on the right.

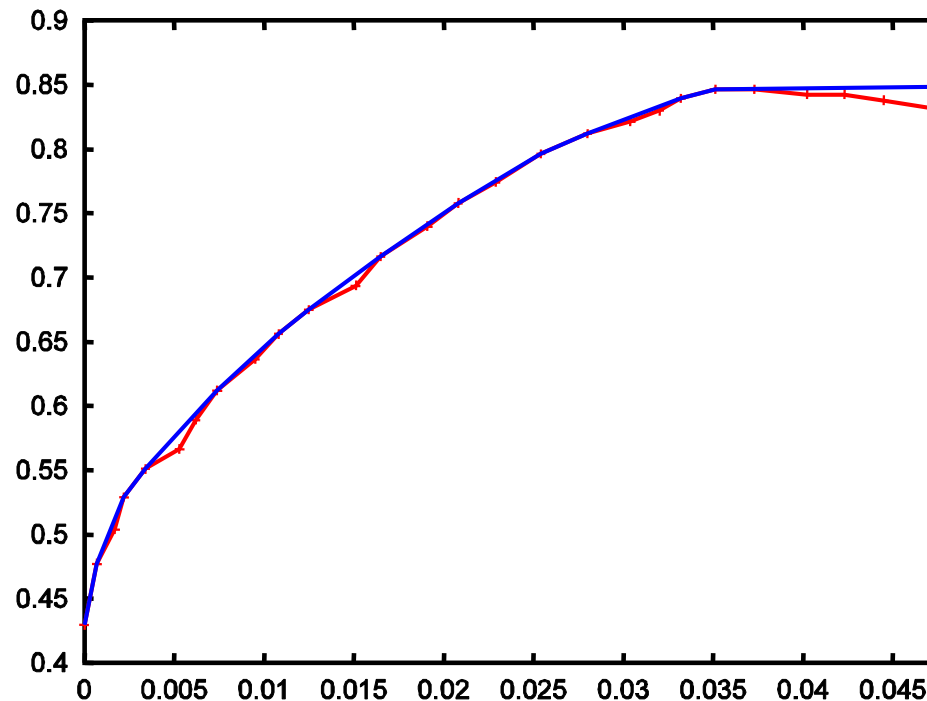
Predict negative ← → Predictive positive



Figures from: http://en.wikipedia.org/wiki/Receiver_operating_characteristic

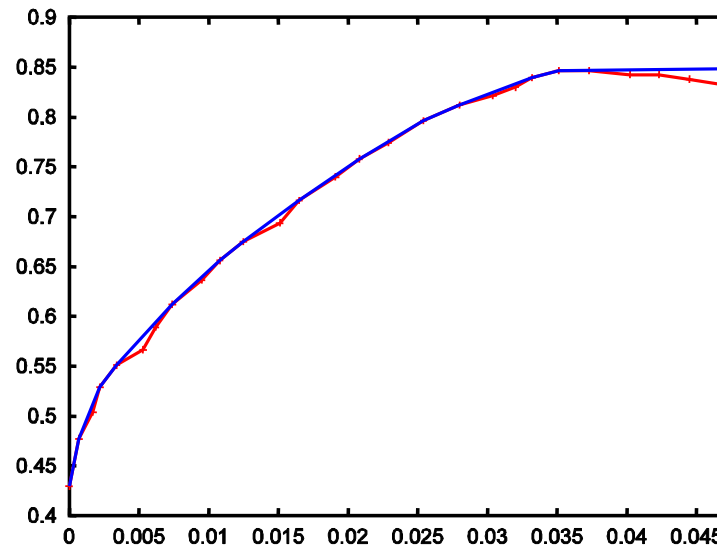
Building the ROC curve

- In many domains, the empirical ROC curve will be non-convex (red line). Take the convex hull of the points (blue line).



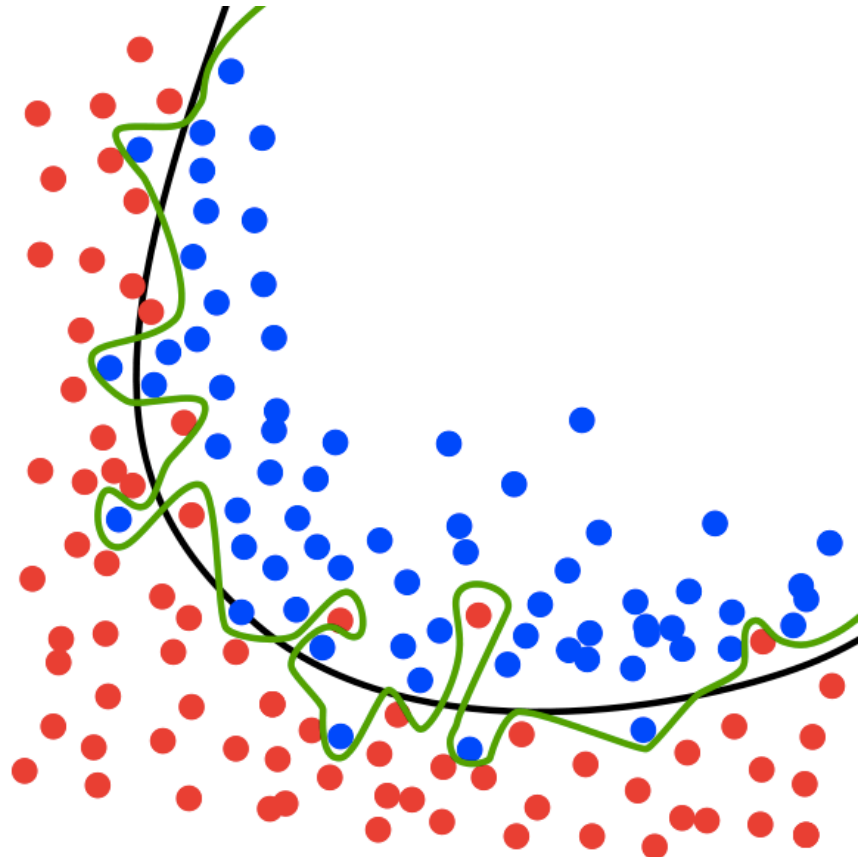
Using the ROC curve

- To compare 2 algorithms over a range of classification thresholds, consider the Area Under the Curve (AUC).
 - A perfect algorithm has $AUC=1$.
 - A random algorithm has $AUC=0.5$.
 - Higher AUC doesn't mean all performance measures are better.



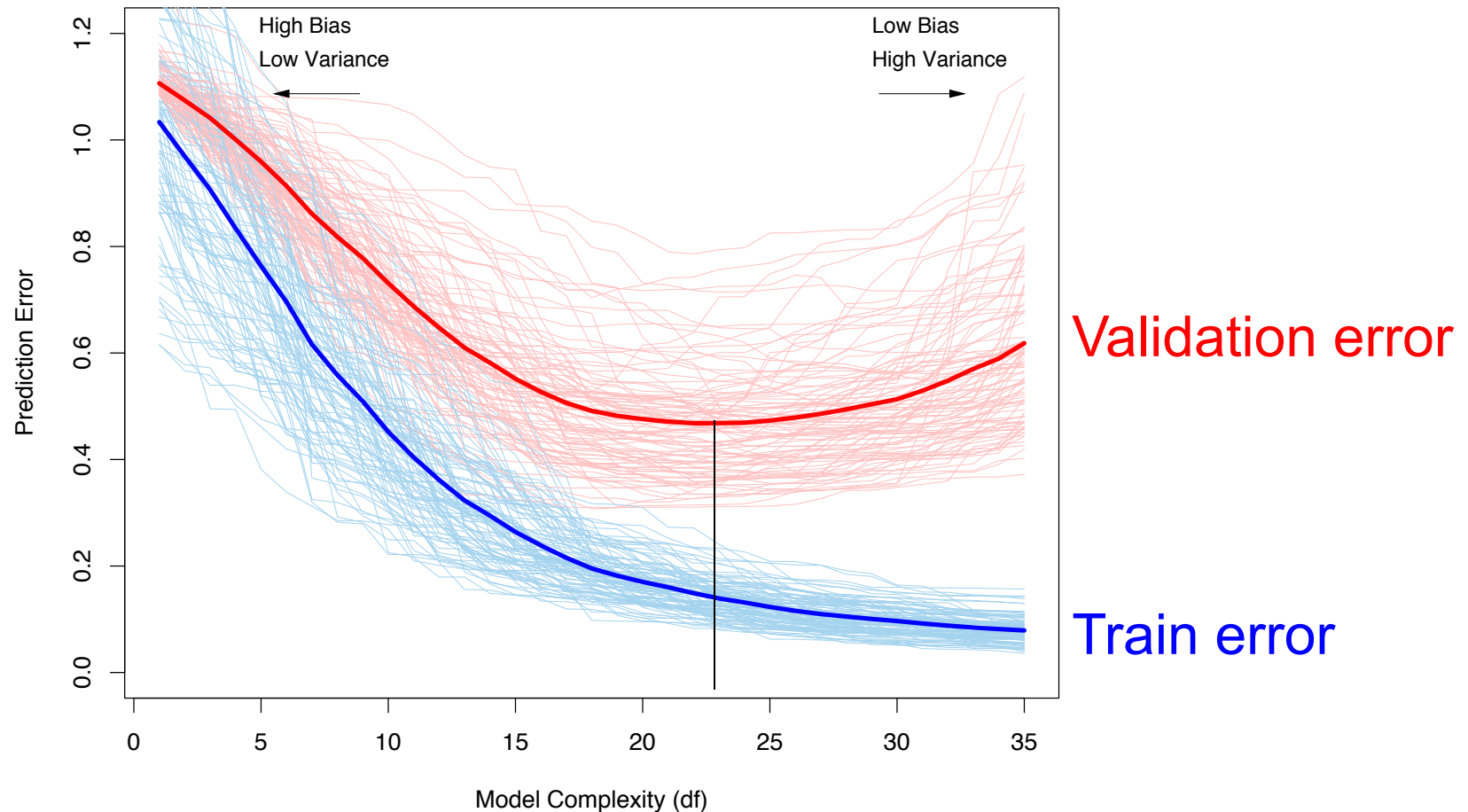
Overfitting

- We have seen that adding more degrees of freedom (more features) always seems to improve the solution!



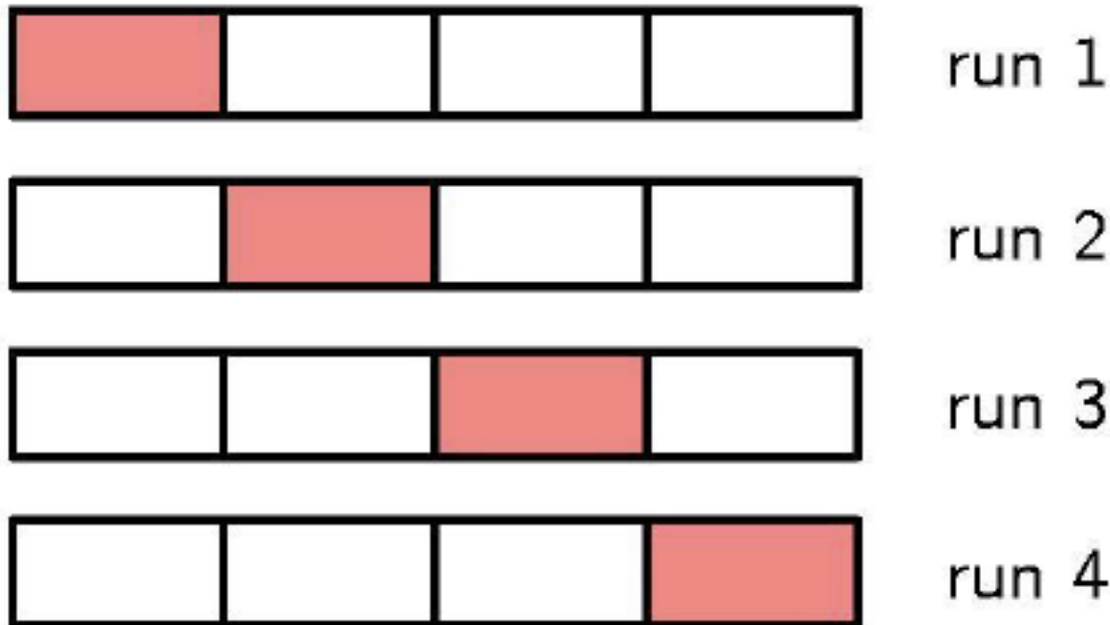
Minimizing the error

- Find the low point in the validation error:



K-fold cross-validation

- **Single test-train split:** Estimation test error with high variance.
- **4-fold test-train splits:** Better estimation of the test error, because it is averaged over four different test-train splits.



K-fold cross-validation

- $K=2$: High variance estimate of $\text{Err}()$.
Fast to compute.
- $K>2$: Improved estimate of $\text{Err}()$; wastes $1/K$ of the data.
 K times more expensive to compute.

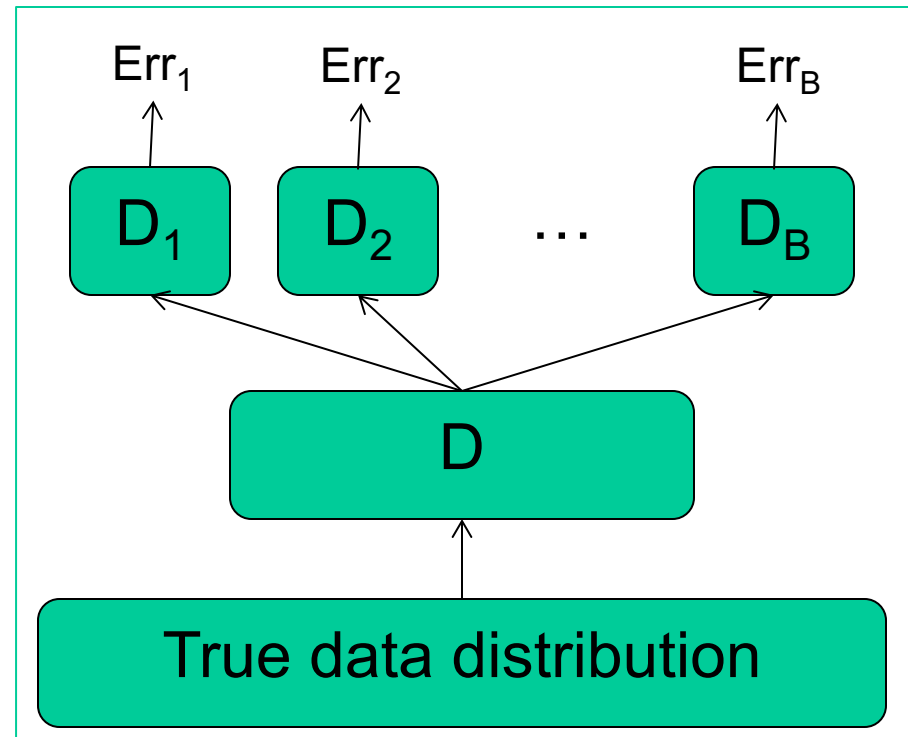
K-fold cross-validation

- $K=2$: High variance estimate of $\text{Err}()$.
Fast to compute.
- $K>2$: Improved estimate of $\text{Err}()$; wastes $1/K$ of the data.
 K times more expensive to compute.
- $K=N$: Lowest variance estimate of $\text{Err}()$. Doesn't waste data.
 N times slower to compute than single train/validate split.

Brief aside: Bootstrapping

- **Basic idea:** Given a dataset D with N examples.
 - Randomly draw (with replacement) B datasets of size N from D .
 - Estimate the measure of interest on each of the B datasets.
 - Take the mean of the estimates.

Is this a good measure
for estimating the error?



Bootstrapping the error

- Use a dataset b to fit a hypothesis f^b . Use the original dataset D to evaluate the error. Average over all bootstrap sets b in B .

$$\widehat{\text{Err}}_{\text{boot}} = \frac{1}{B} \frac{1}{N} \sum_{b=1}^B \sum_{i=1}^N L(y_i, \hat{f}^{*b}(x_i)).$$

- **Problem:** Some of the same samples are used for training the learning and validation.

Bootstrapping the error

- Use a dataset b to fit a hypothesis f^b . Use the original dataset D to evaluate the error. Average over all bootstrap sets b in B .

$$\widehat{\text{Err}}_{\text{boot}} = \frac{1}{B} \frac{1}{N} \sum_{b=1}^B \sum_{i=1}^N L(y_i, \hat{f}^{*b}(x_i)).$$

- **Problem:** Some of the same samples are used for training the learning and validation.
- **Better idea:** Include the error of a data sample i only over classifiers trained with those bootstrap sets b in which i isn't included (denoted C^{-i}).
$$\widehat{\text{Err}}^{(1)} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{*b}(x_i)).$$

(Note: Bootstrapping is a very general ideal, which can be applied for empirically estimating many different quantities.)

Strategy #1

Consider a classification problem with a large number of features, greater than the number of examples ($m \gg n$). Consider the following strategies to avoid over-fitting in such a problem.

Strategy 1:

1. Check for correlation between each feature (individually) and the output. Keep a small set of features showing strong correlation.
2. Divide the examples into k groups at random.
3. Using the features from step 1 and the examples from $k-1$ groups from step 2, build a classifier.
4. Use this classifier to predict the output for the examples in group k and measure the error.
5. Repeat steps 3-4 for each group to produce the cross-validation estimate of the error.

Strategy #2

Consider a classification problem with a large number of features, greater than the number of examples ($m \gg n$). Consider the following strategies to avoid over-fitting in such a problem.

Strategy 2:

1. Divide the examples into k groups at random.
2. For each group, find a small set of features showing strong correlation with the output.
3. Using the features and examples from $k-1$ groups from step 1, build a classifier.
4. Use this classifier to predict the output for the examples in group k and measure the error.
5. Repeat 2-4 for each group to produce the cross-validation estimate of the error.

Strategy #3

Consider a classification problem with a large number of features, greater than the number of examples ($m \gg n$). Consider the following strategies to avoid over-fitting in such a problem.

Strategy 3:

1. Randomly sample n' examples.
2. For the sampled data, find a small set of features showing strong correlation with the output
3. Using the examples from step 1 and features from step 2, build a classifier.
4. Use this classifier to predict the output for those examples in the dataset that are not in n' and measure the error.
5. Repeat steps 1-4 k times to produce the cross-validation estimate of the error.

Summary of 3 strategies

Strategy 1:

1. Check for correlation between each feature (individually) and the output. Keep a small set of features showing strong correlation.
2. Divide the examples into k groups at random.
3. Using the features from step 1 and the examples from $k-1$ groups from step 2, build a classifier.
4. Use this classifier to predict the output for the examples in group k and measure the error.
5. Repeat steps 3-4 for each group to produce the cross-validation estimate of the error.

Strategy 2:

1. Divide the examples into k groups at random.
2. For each group, find a small set of features showing strong correlation with the output.
3. Using the features and examples from $k-1$ groups from step 1, build a classifier.
4. Use this classifier to predict the output for the examples in group k and measure the error.
5. Repeat 2-4 for each group to produce the cross-validation estimate of the error.

Strategy 3:

1. Randomly sample n' examples.
 2. For the sampled data, find a small set of features showing strong correlation with the output.
 3. Using the examples from step 1 and features from step 2, build a classifier.
 4. Use this classifier to predict the output for those examples in the dataset that are not in n' and measure the error.
 5. Repeat steps 1-4 k times to produce the cross-validation estimate of the error.
-

Discussion

- **Strategy 1** is prone to overfitting, because the full dataset is considered in step 1, to select the features. Thus we do not get an unbiased estimate of the generalization error in step 5.
- **Strategy 2** is closest to standard k-fold cross-validation. One can view the joint procedure of selecting the features and building the classifier as the training step, to be applied (separately) on each training fold.
- **Strategy 3** is closer to a **bootstrap** estimate. It can give a good estimate of the generalization error, but the estimate will possibly have higher variance than the one obtained using Strategy 2.

What can we use validation set for?

- **Selecting model class (e.g. number of features, type of features: Exp? Log? Polynomial? Fourier basis?)**
- **Selecting the algorithm (e.g. logistic regression vs naïve Bayes vs LDA)**
- **Selecting hyper-parameters**
 - We often call weights w (or other unknowns in the model) parameters. These are found by algorithm
 - Hyper-parameters are tunable values of the algorithm itself (learning rate, stopping criteria, algorithm-dependent params)
 - Also: regularization parameter λ

A word of caution

- Intensive use of cross-validation can overfit!
- E.g. Given a dataset with 50 examples and 100 features.
 - Consider using any subset of features – 2^{100} possible models!
 - The best of these models will look very good!
 - But it would have looked good even if the output was random!
 - no guarantee it has captures any real pattern in data
 - So no guarantee that it will generalize

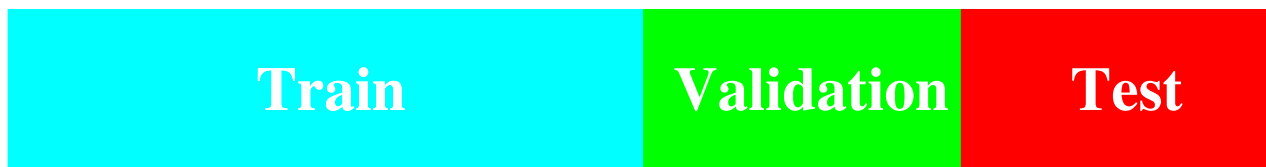
What should we do about this?

Remember from lecture 3

- After adapting the weights to minimize the error on the **train set**, the weights could be exploiting particularities in the **train set**:
 - have to use the **validation set** as proxy for true error
- After choosing the hypothesis class (or other properties, e.g. λ) to minimize error on the **validation set**, the hypothesis class (or other properties) could be adapted to some particularities in the **validation set**
 - **Validation set** is no longer a good proxy for the true error!

To avoid overfitting to the validation set

- When you need to optimize **many** parameters of your model or learning algorithm.
- Use three datasets:
 - The **training set** is used to estimate the parameters of the model.
 - The **validation set** is used to estimate the prediction error for the given model.
 - The **test set** is used to estimate the generalization error once the model is fixed.



What error is measured?

- Scenario: Model selection with validation set. Final evaluation with test set
- Validation error is unbiased error for the current model class
- $\text{Min}(\text{validation error})$ is not an unbiased error for the best model
 - Consequence of using same error to select and evaluate model
- Test error is an unbiased estimate for the chosen model

What can we use test set for?

- Test set should tell us how well the model performs on unseen instances
- If we use test set for any selection purposes, the selection could be based on accidental properties of test set
- **Even if we're just 'taking a peak' during development**
- **The only way to get an unbiased estimate of true loss if is the test set is only used to measure performance of the final model!**


What can we use test set for?

- To prevent overfitting some machine learning competitions limit number of test evaluations
- **Imagenet cheating scandal:** multiple accounts to try more hyperparameters / models on held out test set
- **Not just a theoretical possibility!**

Validation, test, cross validation


- In principle, could ‘cross-validate’ to get estimate of generalization (test-set error)
- In practice, not done so much
 - When designing model, one wants to look at data. This would lead to “strategy 1” from before
 - Having two “cross validation” loops inside each other would make running this type of evaluation very costly
- **So typically:**
 - Test set held out from very beginning. Shouldn’t even look at it
 - Validation:
 - cross validation if we can afford it
 - Hold out validation set from training data if we have plenty of data, or method too expensive for cross validation

Kaggle

 Customer Solutions ▾ Competitions Community ▾ [Sign Up](#) [Login](#)


Welcome to Kaggle, the leading platform for predictive modeling competitions. Here's how to jump into competing on Kaggle —

New to Data Science? [Visit our Wiki »](#)
[Learn about hosting a competition »](#)
[in-Class & Research competitions »](#)




Enter

Find a competition & download the training data. You don't need new software/skills to submit.



Build

Build a model using whatever methods you prefer and upload your predictions to Kaggle.






...Win!

Kaggle scores your solution in real time and you'll see your place on the live leaderboard.

Active Competitions

All Competitions

Active Competitions

	Flight Quest 2: Flight Optimization Optimize flight routes based on current weather and traffic.	10 hours 127 teams \$250,000
	Belkin Energy Disaggregation Competition Disaggregate household energy consumption into individual appliances	35 days 102 teams \$25,000
	Personalize Expedia Hotel Searches - ICDM 2013 Learning to rank hotels to maximize purchases	40 days 138 teams \$25,000

<http://www.kaggle.com/competitions>

Lessons for evaluating ML algorithms

- Error measures are tricky! Always compare to a simple baseline:
 - In classification:
 - Classify all samples as the majority class.
 - Classify with a threshold on a single variable.
 - In regression:
 - Predict the average of the output for all samples.
 - Compare to a simple linear regression.
- Use K-fold cross validation to properly estimate the error. If necessary, use a validation set to estimate hyper-parameters.
- Consider appropriate measures for fully characterizing the performance: Accuracy, Precision, Recall, F1, AUC.

Machine learning that matters

- What can our algorithms do?
- Help make money? Save lives? Protect the environment?
- Accuracy (etc) does not guarantee our algorithm is useful

- How can we develop algorithms and applications that matter?

K. Wagstaff, “Machine Learning that Matters”, ICML 2012.

<http://www.wkiri.com/research/papers/wagstaff-MLmatters-12.pdf>

What you should know

- Understand the concepts of loss, error function, bias, variance.
- Commit to correctly applying cross-validation.
- Understand the common measures of performance.
- Know how to produce and read ROC curves.
- Understand the use of bootstrapping.
- Be concerned about good practices for machine learning!

Read this paper today!

K. Wagstaff, “Machine Learning that Matters”, ICML 2012.

<http://www.wkiri.com/research/papers/wagstaff-MLmatters-12.pdf>
