

# **Cross Validation**

# **Generally: Cross Validation (CV)**

Set of **validation techniques** that use the training dataset itself to validate model

- Allows maximum allocation of training data from original dataset
- Efficient due to advances in processing power

Cross validation is used to test the effectiveness of any model or its modified forms.



# **Validation Goal**

- Estimate Expected Prediction Error
- Best Fit model
- Make sure that the model does not Overfit



Hastie et al. "Elements of Statistical Learning."

# **HoldOut Validation**

#### Dataset



# **HoldOut Validation**

#### **Training Sample**

#### **Testing Sample**



# **HoldOut Validation**

#### **Training Sample**

**Testing Sample** 

Advantage: Traditional and Easy Disadvantage: Varying Error based on how to sample testing





Often used in practice with *k*=5 or *k*=10.

Create equally sized *k* partitions, or **folds**, of training data

For each fold:

- Treat the *k-1* other folds as training data.
- Test on the chosen fold.

The average of these errors is the validation error



#### Dataset

# Suppose K = 10, 10-Fold CV



| Training Sample |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| Training Sample | Training Sample | Training Sample | Training Sample | Testing Sample  |



| Training Sample |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| Training Sample | Training Sample | Training Sample | Training Sample | Testing Sample  |

#### Calculate RMSE = rmse1



| Training Sample |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| Training Sample | Training Sample | Training Sample | Testing Sample  | Training Sample |



| Training Sample |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| Training Sample | Training Sample | Training Sample | Testing Sample  | Training Sample |

#### Calculate RMSE = rmse2



| Training Sample |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| Training Sample | Training Sample | Testing Sample  | Training Sample | Training Sample |



| Training Sample |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| Training Sample | Training Sample | Testing Sample  | Training Sample | Training Sample |

#### Calculate RMSE = rmse3



# And so on



Testing Sample	Training Sample	Training Sample	Training Sample	Training Sample
Training Sample				

#### Calculate RMSE = rmse10



Testing Sample	Training Sample	Training Sample	Training Sample	Training Sample
Training Sample				

#### RMSE = Avg(rmse1...10)



#### Less matters how we divide up

# Selection bias not present







#### Dataset



#### **Training Sample**



# What just happened?



#### **Training Sample**



**Testing Sample** 

# **Leave-P-Out Validation**



For each data point:

- Leave out p data points and train learner on the rest of the data.
- Compute the test error for the p data points.

Define average of these <sub>n</sub>C<sub>p</sub> error values as validation error





## **Leave-P-Out Validation**

A really exhaustive and thorough way to validate

High Computation Time



# **Question:**

# What's the difference between 10-fold and leave-5-out given dataset n=50?





Your problem set: Final Project

Next week: Ensemble







# **Meta-Learning**

# **Layers of Learning**

Gilberto Titericz Junior (top-ranked user on <u>Kaggle.com</u>) used this setup to win the \$10,000 Otto Group Product Classification Challenge.







# **Introduction: Ensemble Averaging**

Basic ensemble composed of a **committee** of learning algorithms.

Results from each algorithm are averaged into a final result, reducing variance.





Logit	SVM	KNN	Majority Voting	Actual
А	A	В		А
В	A	В		В
A	A	А		Α
A	В	В		В



Logit	SVM	KNN	Majority Voting	Actual
А	А	В	А	A
В	A	В		В
A	A	A		A
A	В	В		В



Logit	SVM	KNN	Majority Voting	Actual
А	А	В	А	А
В	А	В	В	В
A	A	A		A
A	В	В		В



Logit	SVM	KNN	Majority Voting	Actual
А	A	В	А	А
В	А	В	В	В
А	A	А	А	A
A	В	В		В



Logit	SVM	KNN	Majority Voting	Actual
А	A	В	А	A
В	А	В	В	В
А	А	А	A	A
А	В	В	В	В



# Ensemble

Meta-Learning

# **Ensembles and Hypotheses**

- Recall the definition of "hypothesis."
- Machine learning algorithms search the **hypothesis space** for hypotheses.
  - Set of mathematical functions on real numbers
  - Set of possible classification boundaries in feature space
- More searchers → more likely to find a "good" hypothesis



# **General Definition**

**One Hypothesis** 

**One Hypothesis** 

**One Hypothesis** 



One Strong Hypothesis

**One Hypothesis** 



# Why so many models?

A single model on its own is often prone to bias and/or variance.

- **Bias** Systematic or "consistent" error. Associated with underfitting.
- **Variance** Random or "deviating" error. Associated with overfitting.

A tradeoff exists. We want to minimize both as much as we can.





# **Three Main Types**



# **Three Main Types**



# Bagging

Short for **b**ootstrap **<u>agg</u>**regat<u>ing</u>.

A **parallel ensemble**. Models are applied without knowledge of each other.

- Apply each model on a random subset of data.
- Combine the output by averaging (for regression) or by majority vote (for classification)
- A more sophisticated version of ensemble averaging.



Bagging decreases variance and prevents overfitting.

# **Random Forests**

Designed to improve accuracy over CART. Much more difficult to overfit

- Works by building a large number of CART trees
  - Each tree in the forest "votes" on outcome
  - Outcome with the most votes becomes our prediction





# Boosting

A **sequential ensemble**. Models are applied one-by-one based on how previous models have done.

- Apply a model on a subset of data.
- Check to see where the model has badly classified data.
- Apply another model on a new subset of data, giving preference to data badly classified by the model.

Boosting decreases bias and prevents underfitting.



# **Weak Learners**

Important concept in boosting.

Weak learners do only slightly better than the baseline for a given dataset. In isolation, they are not very useful.

While boosting, we improve these learners sequentially to create hyper-powered models.





Short for **<u>ada</u>**ptive **<u>boost</u>**ing. Sequentially generates weak learners, adjusting newer learners based on mistakes of older learners

Combines output of all learners into weighted sum





 $H(x) = sign(\alpha_1 h_1(x) + \alpha_2 h_2(x) + \alpha_3 h_3(x))$ 













D3



# XGBoost

Short for e<u>X</u>treme <u>G</u>radient <u>Boost</u>ing. Sequentially generates weak learners like AdaBoost

- Updates model by computing cost function
  - Computes gradient of cost function
  - Direction of greatest decrease = negative of gradient
  - Creates new learner with parameters adjusted in this direction





# **Stacking**



# Stacking



# **Stacking**



# Stacking pt. 1

Assumption: can improve performance by taking a **weighted average** of the predictions of models.

- Take a bunch of machine learning models.
- Apply these models on subsets of your data (how you choose them is up to you).
- Obtain predictions from each of the models.





# Stacking pt. 2

Once we have predictions from each individual model...

- Perform Top-Layer-ML on the predictions.
  - This gives you the coefficients of the weighted average.
- Result: a massive blend of potentially hundreds of models.





# **CDS Core Team Example: Stacking**

CDS Kaggle Team (2017 March Madness Kaggle competition)

- Each member of the Kaggle team made a logistic regression model based on different features
- Combined these using a stacked model







Your problem set: Final Project

Next week: Thank you all!





# **Random Forest Parameters**

- Minimum number of observations in a branch
  - o min\_samples\_split parameter
  - Smaller the node size, more branches, longer the computation
- Number of trees
  - o n\_estimators parameter
  - Fewer trees means *less accurate* prediction
  - More trees means *longer computation* time
  - Diminishing returns after a couple hundred trees

