

"Essentially, all models are wrong, but some are useful." - George Box

Own-Wester LLC

Practical Data

O

Science wit



- Data scientists become expert in all these steps.
- · Managers need to be expert in at least model evaluation.

**Biography** 

Win-Vector LLC Dr. Nina Zumet founder at Win-Vector LLC a San Francisco data science ee from UC

#### Win Vector LLC Dr. John Mount is a prin

-0

ant and founde In mould be a prime company, John has worked as a computational scientist in biotechnology and ading algorithm designer and has managed a research team for Shopping.com (now an eBay bed and big obtained devention in mathematics at ULC Barkeley and holds a PDD in from Carnegie Mello

so the coauthor of Practical Data Science with R (Manning Publi Please contact conta WinVectorLLC . m for projects and collabora

What is a Good Model? Performance metrics for classifiers / decision procedures.

### How do you measure model performance?



### Which Metrics Are Appropriate?

| Metric                   | Example   |
|--------------------------|---|
| Precision                | If the test comes back<br>positive, is the patient really<br>diabetic?  |
| Recall<br>Sensitivity    | Do we miss any diabetics through this test?   |
| Precision<br>Specificity | Diagnoses that lead to costly treatment   |
| Recall<br>Sensitivity    | Diagnosing conditions that are costly if untreated  |
| Accuracy                 |   |
|                          | Metric<br>Precision<br>Recall<br>Sensitivity<br>Precision<br>Specificity<br>Recall<br>Sensitivity<br>Accuracy |

### **Technical Metrics**

- · AUC (ROC), deviance, and others.
- Good metrics for data scientists and between data scientists
- Useful proxy measures for comparing candidate models
- · Not always easily translatable to business goals

### ROC/AUC

- Graph of trade-off between true positive and false positive rates as labeling threshold T is varied.
- AUC: area under the curve
- Probability that a randomly chosen positive example will score higher than a randomly chosen negative example (with appropriate tiebreaking).
- Invariant to monotonic transformations of scoring function
- Independent of target class prevalence

### Predictive modeling schematic

- · Define a useful business goal.
- · Choose a convincing performance measure.
- · Collect input ("independent") variables.
- · Build a model.
- · Confirm you have a decisive model.
- · Refine/repeat.

Own-Wester LL

### **Big risks**

- Not being able to produce a good model ("under fit").
- · Not being able to falsify a bad model ("over fit").
- Model depending on variables that are really only available after the outcome is known ("data leakage").

### Example: KDD2009

- KDD conference 2009 contest data.
- Predict from a few hundred features which credit accounts will "churn" or cancel.
- Training data: measurements from past accounts known to have cancelled or not cancelled in a fixed time interval.

O

Own-Wester LLC

# Assume we have our data ready to go

- Data acquisition, documentation, cleaning, and preparation is by far the largest most critical part of real world data science.
- For this demo we are going to assume this is done and move on to the part people always want to hear about: model construction.

### Pre-packaged software

- You don't need a Ph.D. to perform machine learning, because people with Ph.D.s have already implemented and shared very powerful methods:
  - Gradient boosted trees.
  - Random forests.
  - Deep learning.

• For this demonstration we will exhibit R, decision trees, gradient boosting, and h2o deep learning.





### Step back

• We have shown the typical "left to their own devices" data scientist workflow concentrating on improving models "in the lab" only on our training data.

• In a real application

- The biggest modeling improvements come from commissioning new measurements and features.
- High performance on training is not the true end goal, and must be
  viewed with suspicion (issues like over-fit, and data-leaks are important).

• What evidence do we have we actually solved the problem?

(switch) Please stand by....



| Question #1: W<br>the b | Vhich model is<br>est?          |
|-------------------------|---------------------------------|
|                         | Performance on<br>Training Data |
| Decision Tree           | 0.71                            |
| Gradient Boosting       | 0.73                            |
| Neural Net              | 0.76                            |
|                         | O Wester LLC*                   |









|                   | Performance on | Performance on |
|-------------------|----------------|----------------|
|                   |                |                |
| Decision Tree     | 0.71           | 0.70           |
| Gradient Boosting | 0.73           | 0.73           |
| Neural Net        | 0.76           | 0.69           |





# Question #2: How do you tune the modeling algorithm?

- How many trees for gradient boosting? How deep?
- What's the best learning rate for NN?
- How many iterations before you stop updating the NN?



| Which gradient<br>(xgboost) i | boosting model s the best?      |
|-------------------------------|---------------------------------|
| Number of trees               | Performance on<br>Training Data |
| 50                            | 0.93                            |
| 100                           | 0.98                            |
| 200                           | 1.0                             |
|                               | Owener                          |

| Again:<br>Best on train ≠ Best in future |                                 |                             |  |  |
|--|---------------------------------|-----------------------------|--|--|
| Number of Trees                          | Performance on<br>Training Data | Performance on<br>Test Data |  |  |
| 50                                       | 0.93                            | 0.73                        |  |  |
| 100                                      | 0.98                            | 0.72                        |  |  |
| 200                                      | 1.0                             | 0.70                        |  |  |
|  | Overfit!                        |                             |  |  |









Cross-val vs. Holdout

## Data Science : Data-rich

- · Generally, we will prefer train-validation-test split
  - · Lots of data to spare for holdout
  - Large data sets make computational efficiency attractive

· Possible exception: very rare target class

### When to Consider Cross-val

O<sub>We-Vec</sub>

- · Data sets too small for train-validation-test split
- · Lots of modeling parameters
- Rare target or rare features of interest

# Cross validation and existing packages

- Many modeling procedures have cross-validation or validation set use baked in
  - Picking parameters, stopping criteria, etc.
  - gradient boosting with gbm, h2o neural nets....
- Reduces upward bias, but doesn't eliminate it
  - Still need a holdout set

### Holdout: No peeking! (Or not too much)

- In practice: fit model->evaluate->tweak model ...
  Too many iterations and performance estimates are upwardly biased again
  - Especially if the holdout (or validation) set is small
- Recent differential-privacy related results to alleviate this
   http://www.win-vector.com/blog/2015/10/a-simplerexplanation-of-differential-privacy/



### Thank you

All examples: <u>http://winvector.github.io/ModelTesting/</u>