

PRESENTED BY



strataconf.com
#StrataHadoop

R Day: Validating models in R

- Tuesday, 03/29/2016 11:00 AM – 12:30 PM
Location: LL20C
- Nina Zumel and John Mount, **Win-Vector LLC**.
- We demonstrate a number of techniques, R packages, and R code for validating predictive models, using example code, data, and live demonstrations/exercises. Learn how to determine if there is usable signal in your data, select variables, and choose models using R and R graphics (ggplot2). Increase your statistical efficiency and squeeze more signal out of your data.
- Materials: <https://github.com/WinVector/ValidatingModelsInR>

“Essentially, all models are wrong, but some are useful.”

– George Box

Goals of this Tutorial

- Give you a sophisticated tool box of model quality measures that are:
 - Statistically well founded.
 - Business motivated!
 - Organized by a taxonomy of needs.
 - With ready to go R code and graphs.

Biography

Nina Zumel

Win-Vector LLC

Dr. **Nina Zumel** is a principal consultant and founder at **Win-Vector LLC** a San Francisco data science consultancy and training company. Nina started her advanced education with an EE degree from UC Berkeley and holds a Ph.D. in Robotics from Carnegie Mellon University. Nina has worked as research scientist at SRI and developed revenue optimization platforms. She frequently writes and speaks on statistics and machine learning.

Nina is also the coauthor of the popular book of *Practical Data Science with R* (Manning Publications, 2014).

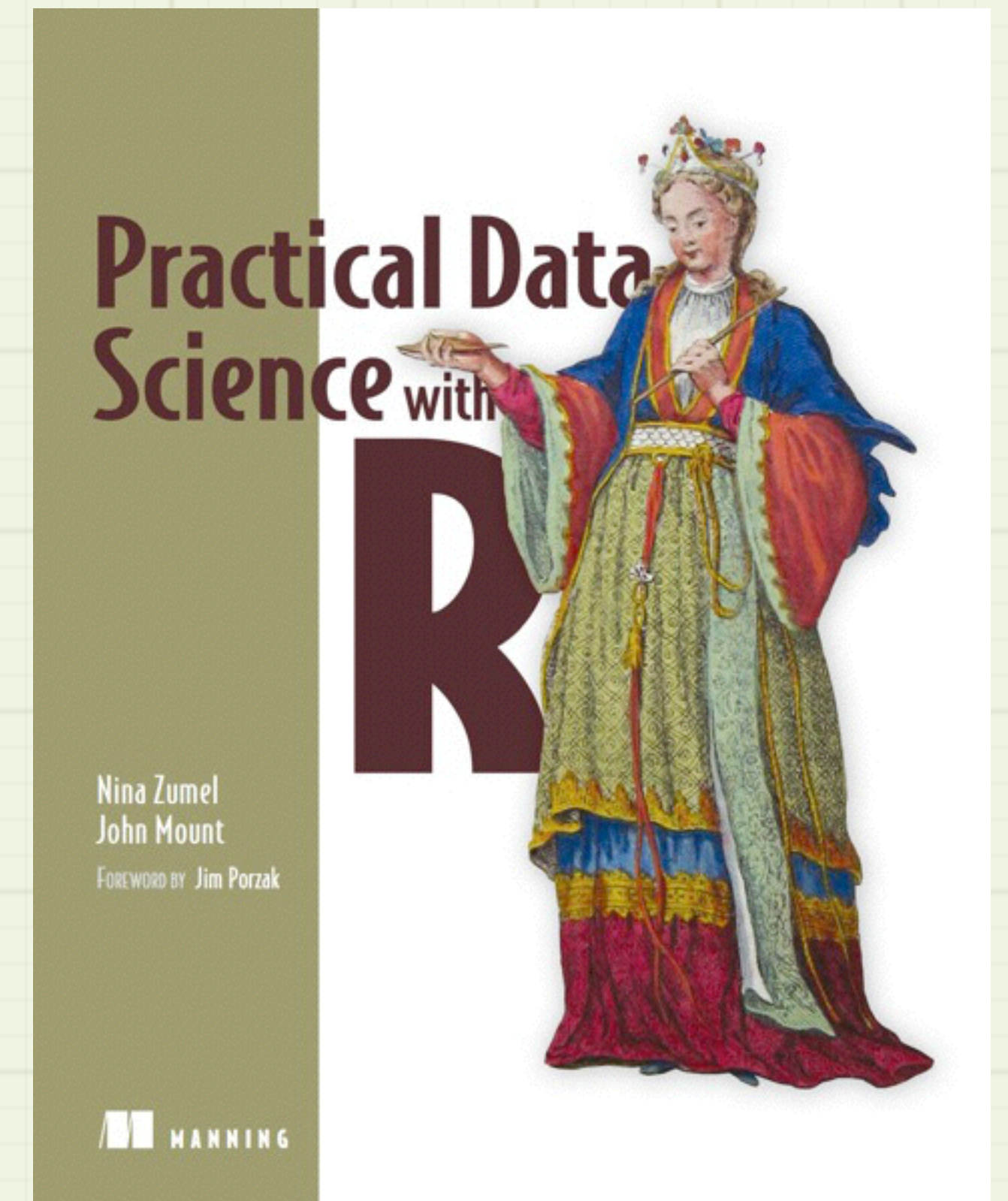
John Mount

Win Vector LLC

Dr. **John Mount** is a principal consultant and founder at **Win-Vector LLC** a San Francisco data science consultancy and training company. John has worked as a computational scientist in biotechnology and a stock-trading algorithm designer and has managed a research team for Shopping.com (now an eBay company). John started his advanced education in mathematics at UC Berkeley and holds a PhD in computer science from Carnegie Mellon.

John is also the coauthor of *Practical Data Science with R* (Manning Publications, 2014).

Please contact contact@win-vector.com for projects and collaborations. <http://win-vector.com/> Twitter: [@WinVectorLLC](https://twitter.com/WinVectorLLC).



Why are we giving a statistics talk during R Day of a data science conference?

- With R classic statistical advice becomes immediately actionable.
- With “big data” you can remove unwanted assumptions and directly estimate model quality.
- If you see something new we can influence how you work. If nothing is new, maybe we can influence how you teach.
- *A great* excuse to use R Markdown, ggplot2, and other packages.

The two issues

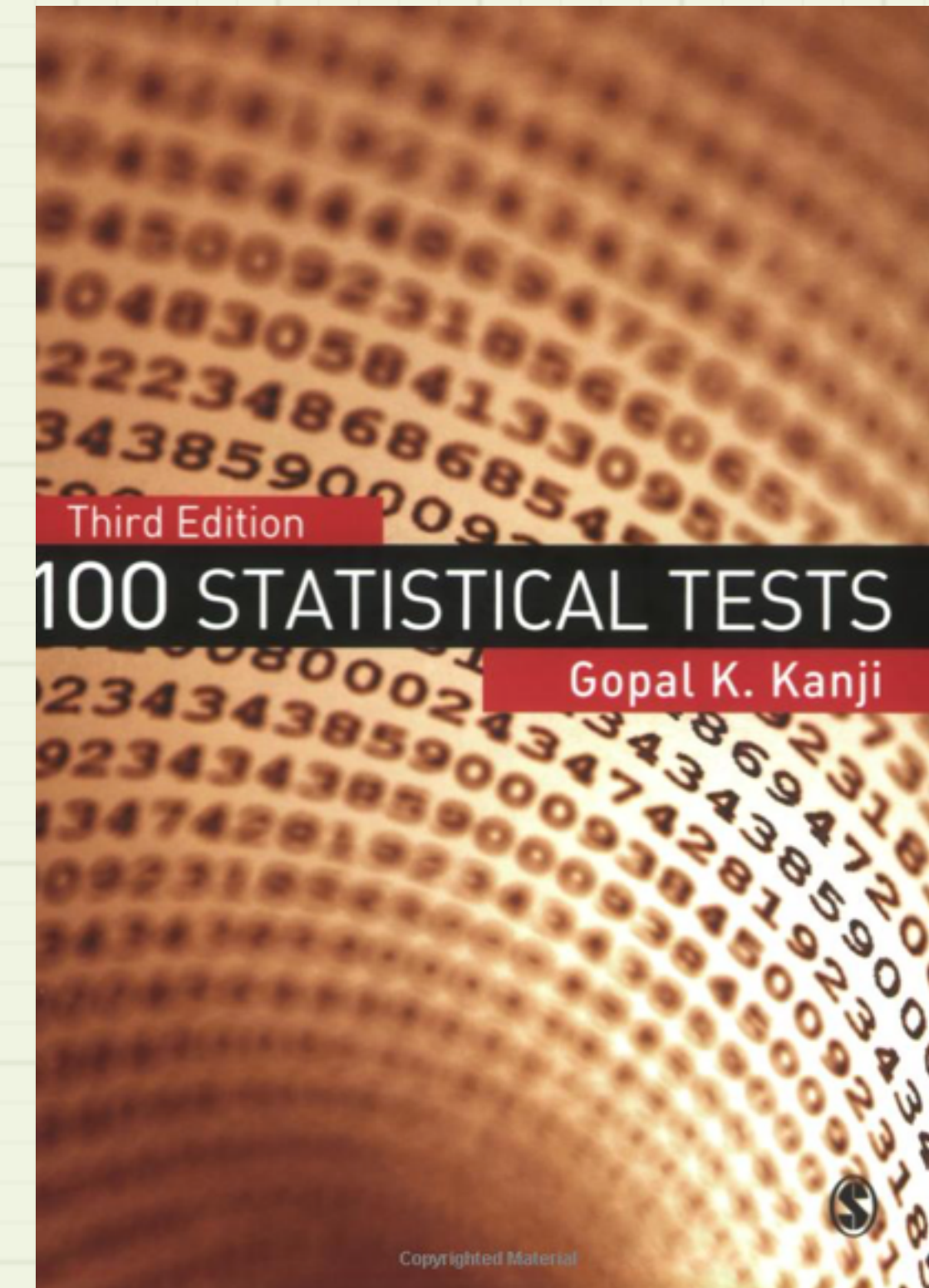
- Choosing a convincing measure.
 - Part 1 of this session.
- Confirming you have a decisive measurement.
 - Part 2 of this session.

Things are simpler when you

Start here

- Separate what to measure from how to confirm significance or estimate posteriors.
- You rely on a programming environment for composable, reusable methods, simulation and visualization.

Not here



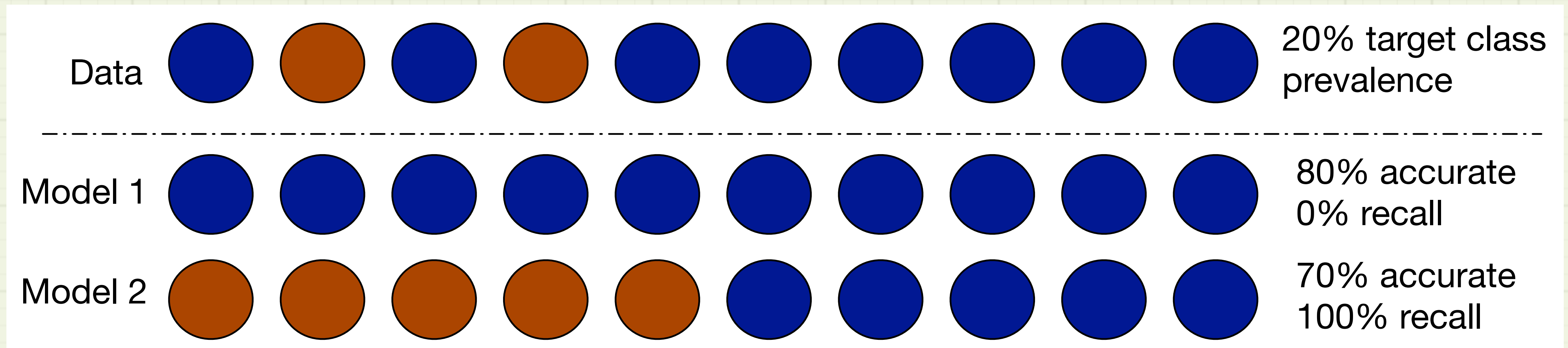
(reserve this to improve solutions later)

What is a Good Model?

Performance Metrics

How do you measure model performance?

- Hint: accuracy is not usually the right way



How are you going to use your model?

- Decision Procedure
 - Predict customer churn; predict home sales price
 - Correct answers are important
- Sorting or prioritization
 - Target at-risk customers for intervention; identify most valuable homes
 - Correct comparisons are important

Metrics for Classifiers

Technical Metrics

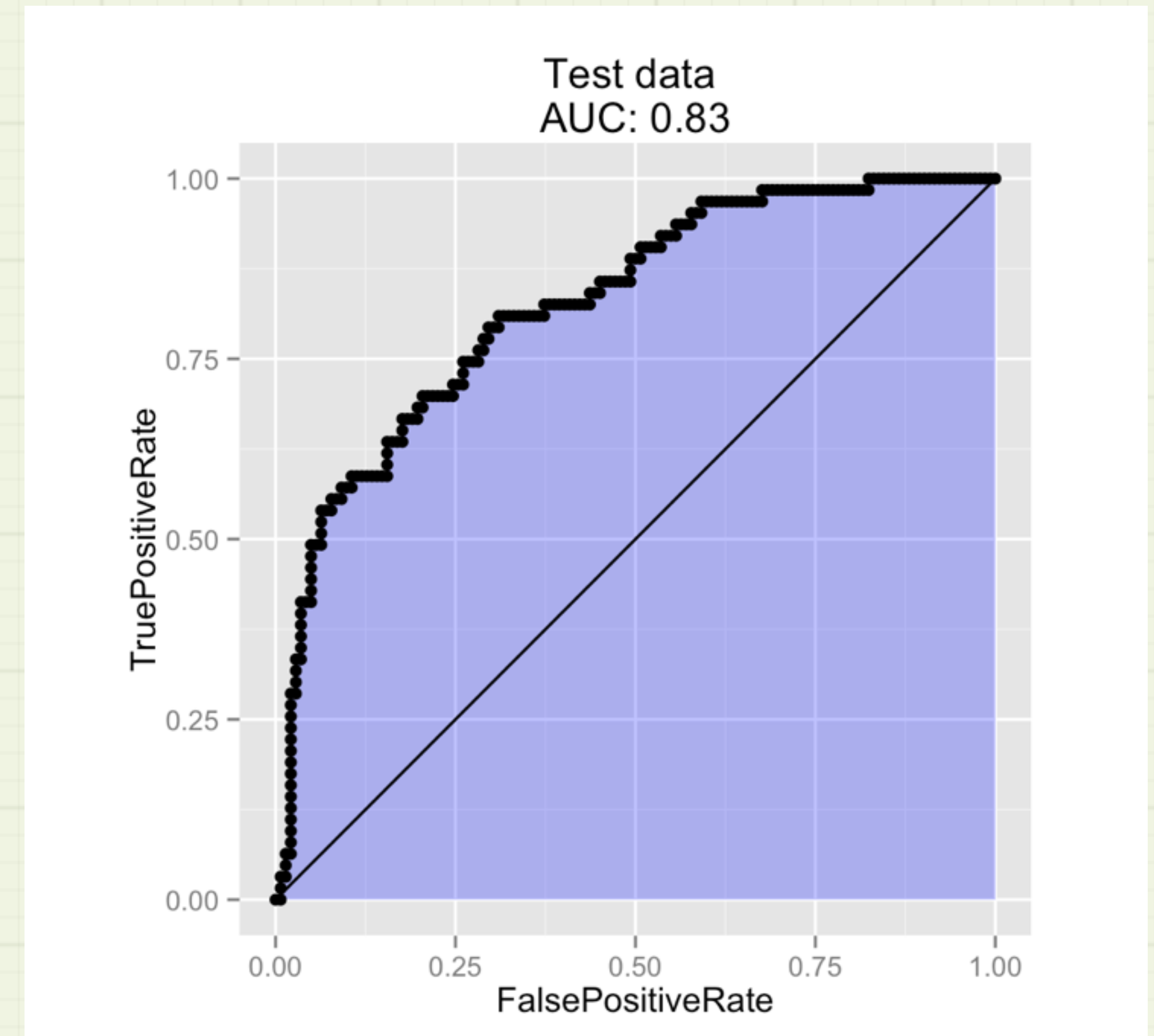
- AUC (ROC), deviance
- Good metrics for data scientists and between data scientists
- Useful proxy measures for comparing candidate models
- Not always easily translatable to business goals

Domain Metrics

- Precision, Recall, Sensitivity, Specificity
- Good metrics for business

ROC/AUC

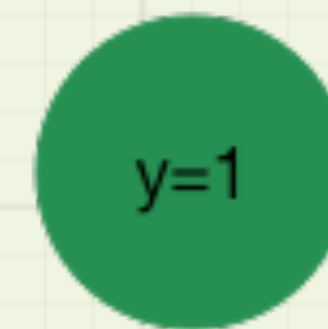
- Trade-off between true positive and false positive rates as labeling threshold T is varied
- AUC: area under the curve
 - Probability that a randomly chosen positive example will score higher than a randomly chosen negative example (with appropriate tie-breaking).
- Invariant to monotonic transformations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



Deviance

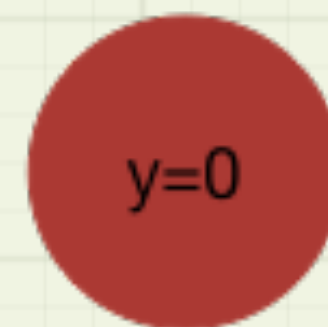
- Deviance penalizes “mismatches” between prediction and true class label
- Analogous to variance
 - smaller is better
- Dependent on dataset size
 - Don't compare unnormalized deviance across evaluation sets of different size

$$\text{deviance} = -2 * [\text{sum}(y * \log(p_y) + (1 - y) * \log(1 - p_y))]$$



and $P(\text{class} = 1) \rightarrow 1$

contribution $\rightarrow 1 * \log(1) \rightarrow 0$



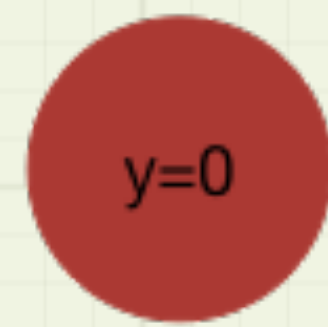
and $P(\text{class} = 1) \rightarrow 0$

contribution $\rightarrow (1 - 0) * \log(1 - 0) \rightarrow 0$



and $P(\text{class} = 1) \rightarrow 0$

contribution $\rightarrow 1 * \log(0) \rightarrow \text{large negative value}$



and $P(\text{class} = 1) \rightarrow 1$

contribution $\rightarrow (1 - 0) * \log(1 - 1) \rightarrow \text{large negative value}$

When Correct Answers are Important

- Confusion matrix
- Recall/Precision
- Sensitivity/Specificity
- Pricing of False Positives, False Negatives

Confusion Matrix

Prediction			
	FALSE	TRUE	
Diabetic FALSE	434	66	500
TRUE	110	158	268
	544	224	768

False Positives

Outcome sums

Prediction sums

False Negatives

Detailed description: The confusion matrix is a 2x2 grid with an additional column for outcome sums. The rows represent the actual 'Diabetic' status (FALSE, TRUE) and the columns represent the predicted status (FALSE, TRUE). The values in the cells are: (Diabetic FALSE, Predicted FALSE) = 434, (Diabetic FALSE, Predicted TRUE) = 66, (Diabetic TRUE, Predicted FALSE) = 110, and (Diabetic TRUE, Predicted TRUE) = 158. The outcome sums are 500 for Diabetic FALSE and 268 for Diabetic TRUE. The prediction sums are 544 for Predicted FALSE and 224 for Predicted TRUE. The total sum is 768. Annotations include: 'False Positives' pointing to the 66 count, 'False Negatives' pointing to the 110 count, 'Outcome sums' pointing to the rightmost column, and 'Prediction sums' pointing to the bottom row.

	Predicted FALSE	Predicted TRUE	Outcome sums
Actual FALSE	434	66	500
Actual TRUE	110	158	268
Prediction sums	544	224	768

Confusion matrix is paramount

- Most common classifier score follow from:
- The confusion matrix
- Or scaled summaries of it such as
 - tpr, fpr, tnr, fnr

False Positive Rate, False Negative Rate:

False Positive Rate

False Positives

All Negatives

$$\frac{66}{(434+66)} = 0.132$$

“The fraction of non-diabetics misdiagnosed as diabetic.”

Prediction	FALSE	TRUE
Diabetic	434	66
Non-Diabetic	110	158

False Positives

All Negatives

False Positive Rate, False Negative Rate: False Negative Rate

False Negatives

All Positives

$$110 / (110 + 158) = 0.41$$

“The fraction of diabetics misdiagnosed as non-diabetic.”

Prediction	FALSE	TRUE
Diabetic	434	66
	110	158

False
Negatives

All Positives

True Positive Rate, True Negative Rate: True Positive Rate

True Positives

All Positives

$$\frac{158}{158+110} = 0.589$$

“The fraction of diabetics correctly identified as such.”

Prediction	FALSE	TRUE
Diabetic	434	66
	110	158

True Positives

All Positives

True Positive Rate, True Negative Rate: True Negative Rate

True Negative

All Negatives

$$\frac{434}{434+66} = 0.868$$

“The fraction of non-diabetics correctly identified as such.”

Prediction	FALSE	TRUE
Diabetic	434	66
	110	158

True
Negatives

All Negatives

What about accuracy?

- Almost all business partners will ask for “accuracy.”
- This is only because it is likely the only score that has been explained to them in any detail.

Accuracy

$$\frac{\sum \text{diagonals}}{\sum \text{entries}}$$
$$\frac{(434 + 158)}{768}$$
$$= 0.77$$

“The fraction of patients correctly diagnosed.”

Prediction	FALSE	TRUE
Diabetic		
FALSE	434	66
TRUE	110	158

Common derived measures

Sensitivity, Specificity: Sensitivity

True Positives

All Positives

$$\frac{158}{158+110} = 0.589$$

“The fraction of diabetics correctly identified as such.”

Prediction	FALSE	TRUE
Diabetic	434	66
	110	158

True Positives

All Positives

Sensitivity, Specificity: Specificity

True Negatives

All Negatives

$$\frac{434}{434+66} = 0.868$$

“The fraction of
non-diabetics
correctly identified
as such.”
(Or 1-FPR)

Prediction	FALSE	TRUE
Diabetic	434	66
	110	158

True
Negatives

All Negatives

Precision, Recall: Precision

True Positives
—
Predicted Positives

$$158 / (158 + 66) = 0.705$$

“The fraction of patients diagnosed as diabetic who really are.”

Prediction	FALSE	TRUE
Diabetic		
FALSE	434	66
TRUE	110	158

True Positives

Predicted Positives

Precision, Recall:

Recall

True Positives

All Positives

$$\frac{158}{158+110} = 0.589$$

“The fraction of diabetics correctly identified as such.”

Prediction	FALSE	TRUE
Diabetic	434	66
	110	158

True Positives

All Positives

labs/ Lab01 ScoringClassifiers

Bringing it all together

Which Metrics Are Appropriate?

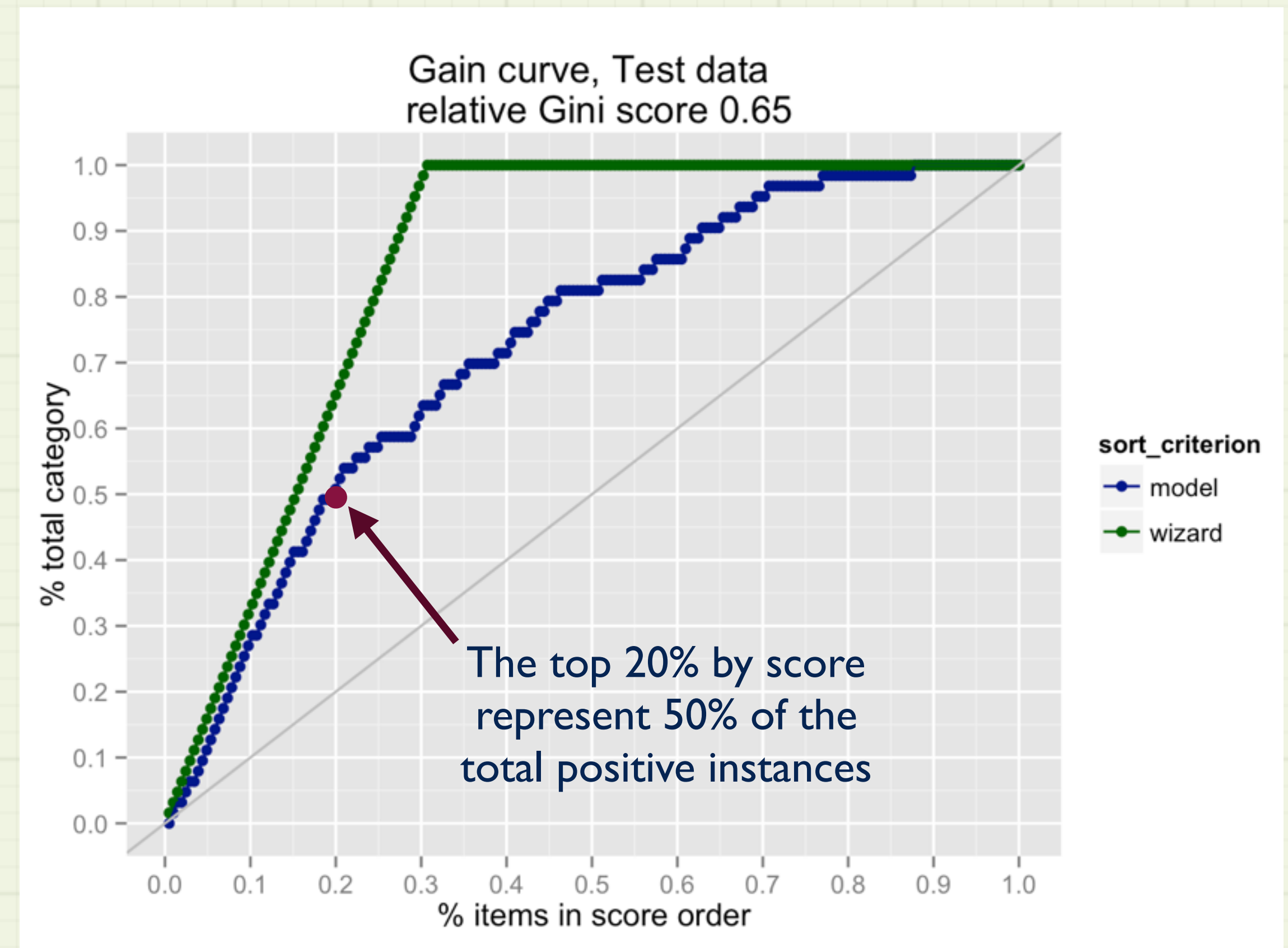
Question	Metric	Example
Is it important that a positive classification is correct?	Precision	If the test comes back positive, is the patient really diabetic?
Is it important we find all positive cases?	Recall Sensitivity	Do we miss any diabetics through this test?
Are false positives expensive?	Precision Specificity	Diagnoses that lead to costly treatment
Are false negatives expensive?	Recall Sensitivity	Diagnosing conditions that are costly if untreated
Is it important to get everything right?	Accuracy	

When Sorting is Important

- Gain Curve
 - Applies for both probability models and general regression models
- Gini Coefficient

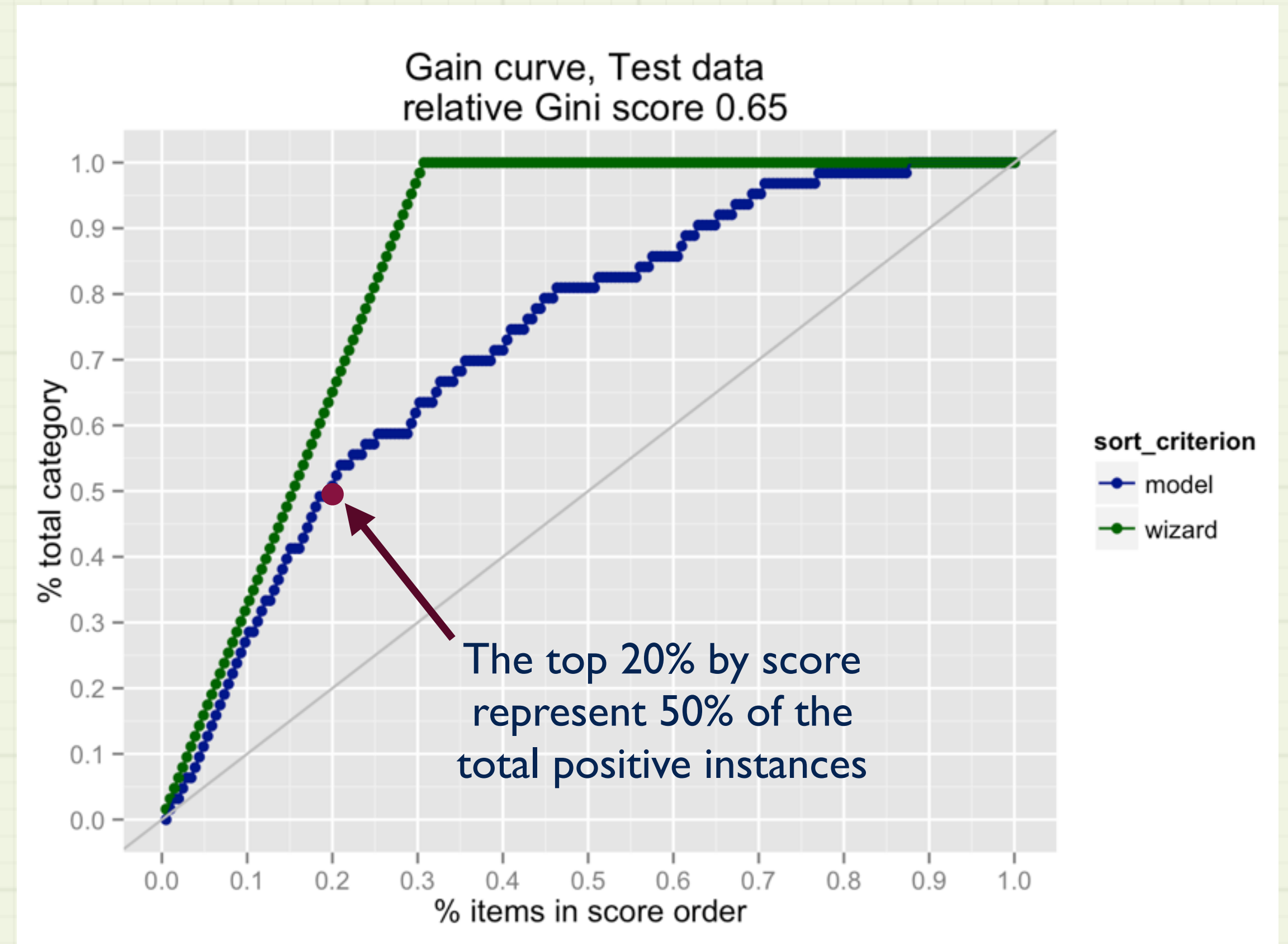
Gain/Lift Curve

- Sort by probability score, descending
- Measures how quickly you identify all positive instances (by sort order)
 - Good when score is used for prioritizing items to act upon
 - E.g. Probability of fraud



Gain Curve (cont.)

- Gini score: $2 \times$ area between the curve and $x=y$ line
- Relative Gini score: ratio of model Gini to ideal Gini score
 - “Wizard” in plot
- **Gain Curve is sensitive to target class prevalence**



labs/Lab02GainCurve

Metrics for Regression

Technical Metrics

- R^2
 - $1 - \sum(y_i - \hat{p}_i)^2 / \sum(y_i - E[y])^2$
 - Fraction of variance “explained” by the model
- Again, useful for comparing candidate models, not always a clear connection to business goals

Domain Metrics

- Root mean squared error
- Mean absolute error
- Mean relative error

Root Mean Squared Error

$$\text{rmse}(y, p) = \sqrt{(\sum (y_i - p_i)^2 / n)}$$

- Estimates the “average” error between outcome and prediction
 - Large differences dominate score
- Related to what linear regression minimizes to fit a model
 - Optimizing RMSE gets expectations and totals correct

Mean Absolute Error

$$\text{mae}(y, p) = (\sum |y_i - p_i| / n)$$

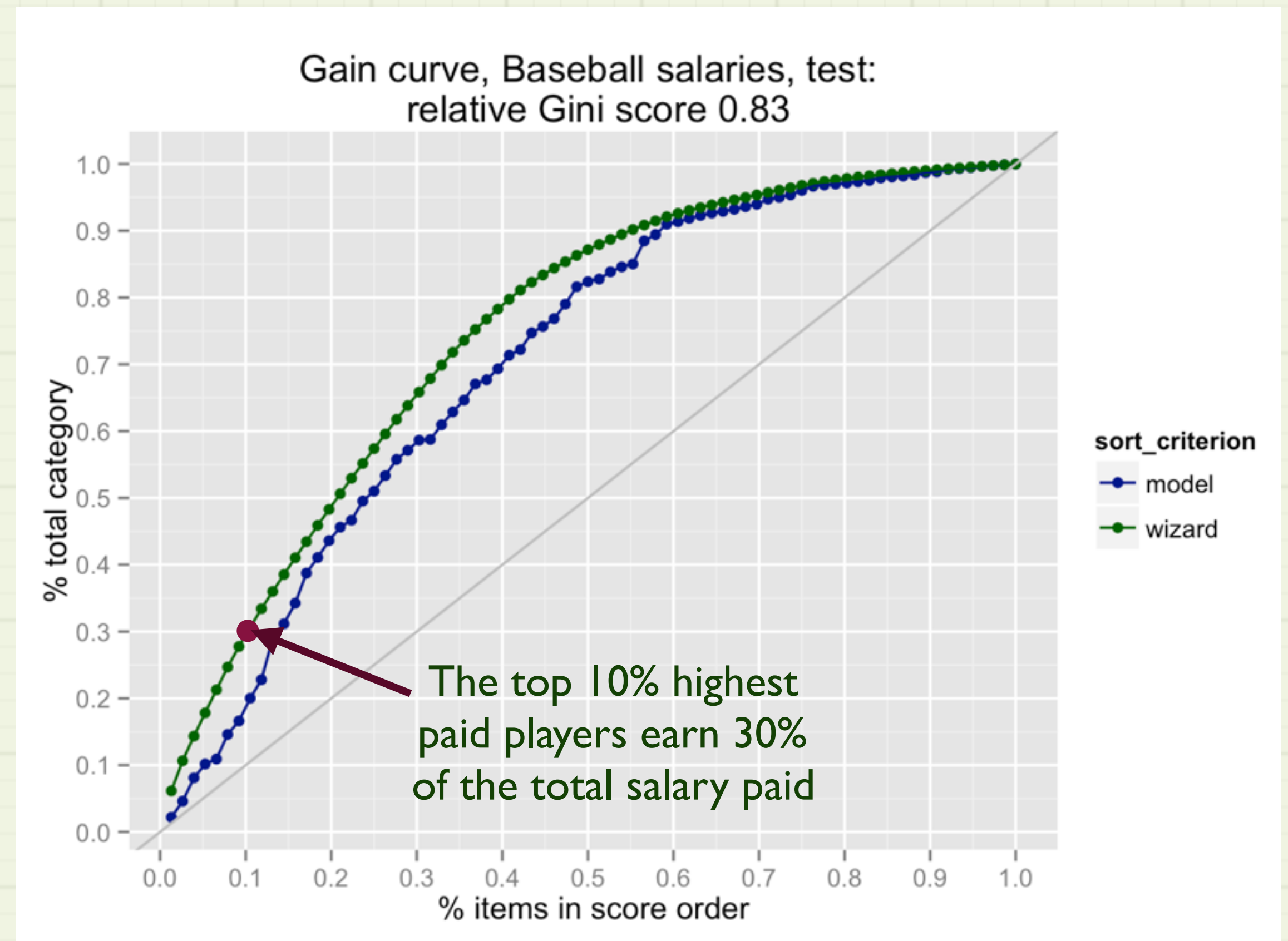
- Estimates the average unsigned error between outcome and prediction
 - Large differences dominate less
 - Arguably what people intuitively consider “average” error
- Optimizing MAE does *not* get expectations and totals correct

Minimizing Relative Error

- For example, by predicting $\log(y)$
- Useful when outcome spans several orders of magnitude
 - \$5 error on \$1,000 different from \$5 error on \$10
- Errors on small magnitude outcomes dominate

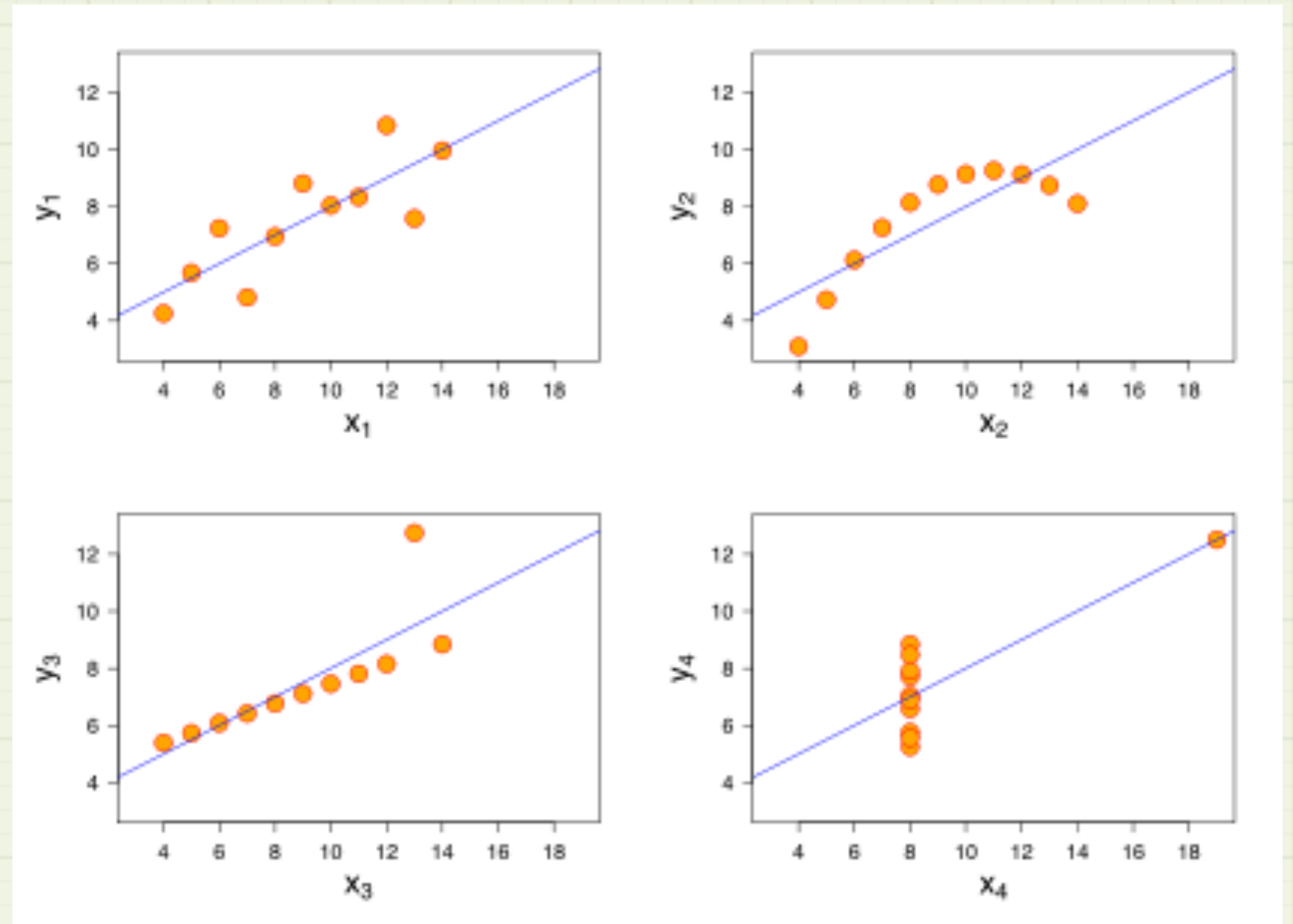
When Order is Important: Gain Curve Revisited

- Sort items by predicted value, descending
- Measure accumulated true value as a fraction of total
- Measures if your predictions are roughly in the right order
 - Do you predict large values as large, and small values as small?
- Random sort: $x=y$ line



Summaries can be deceptive

- Regression: Anscombe's quartet
- Ranking: distraction of indistinguishable pairs



labs/Lab03RankingIssues

Which Metrics Are Appropriate?

Question	Metric	Example
Must we predict the value accurately?	RMSE, mean absolute error	Predicting home sale price
Must we predict the value to a good relative tolerance?	Mean relative error	Predicting home sale price to within 10%
Do we mostly need the order to be right?	Relative Gini Score	Predicting lifetime customer value (Identify most valuable customers)

Speaker Change

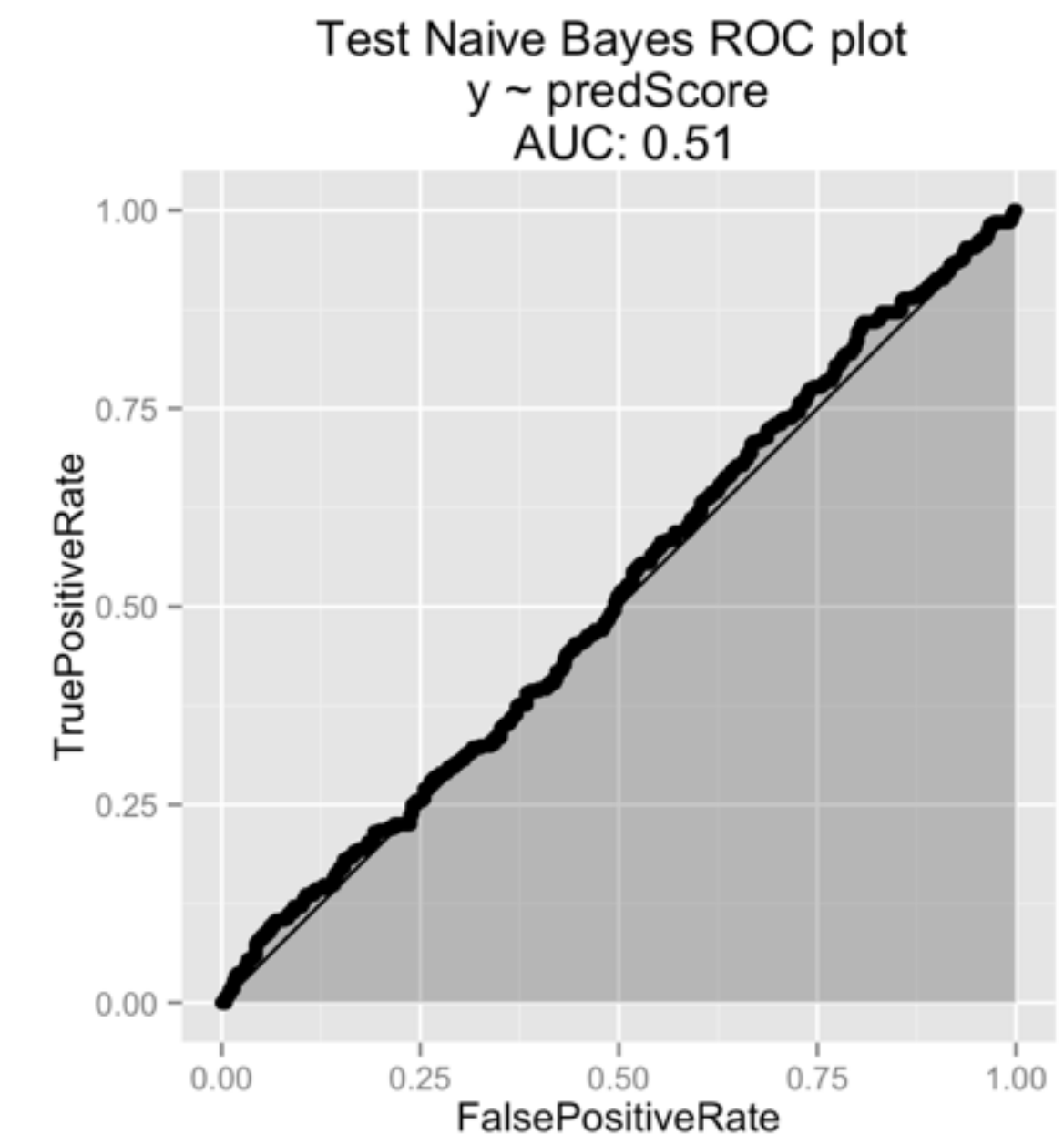
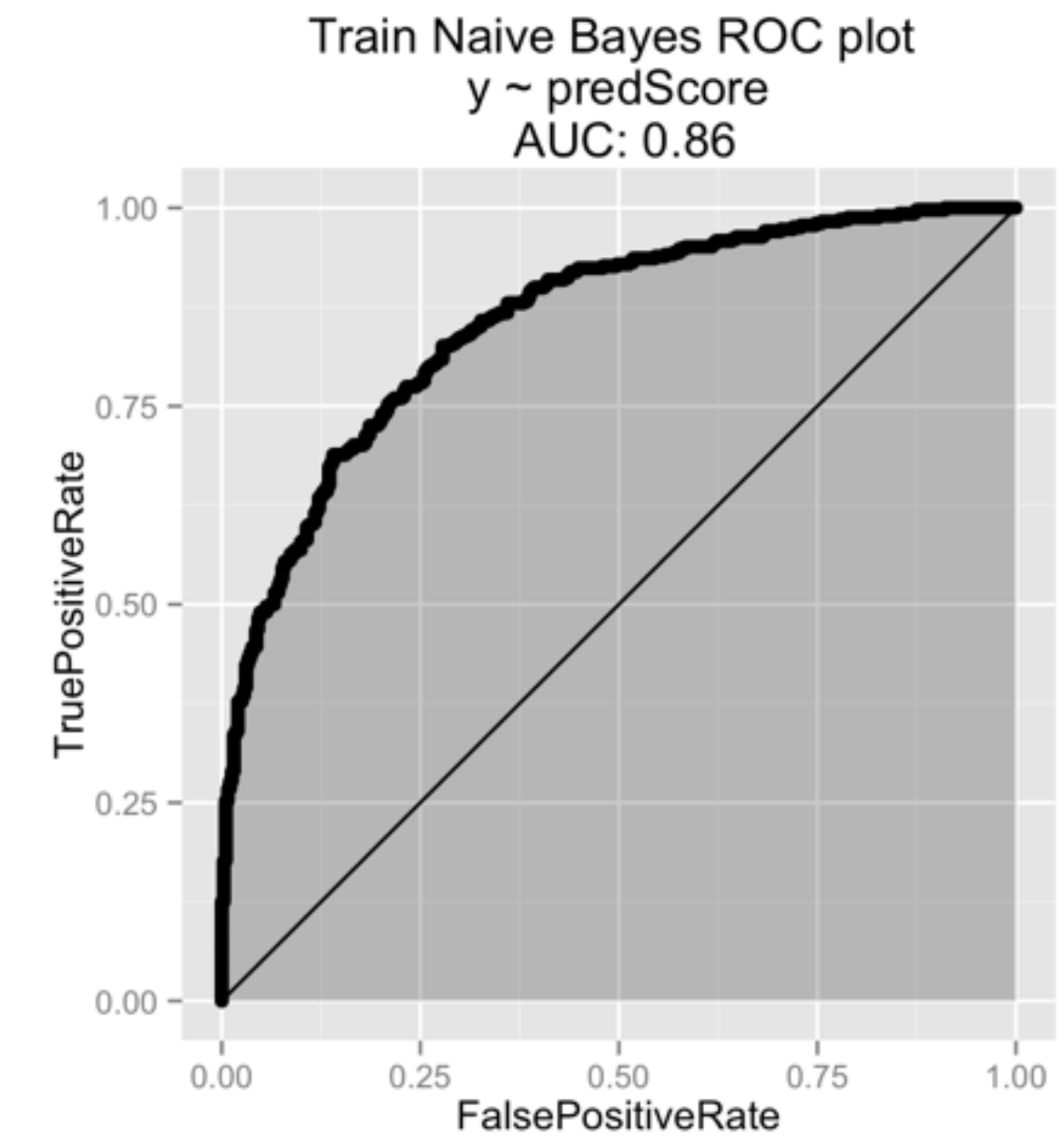
Is it *really* a Good Model?
Estimating *out of sample*
performance

Principles

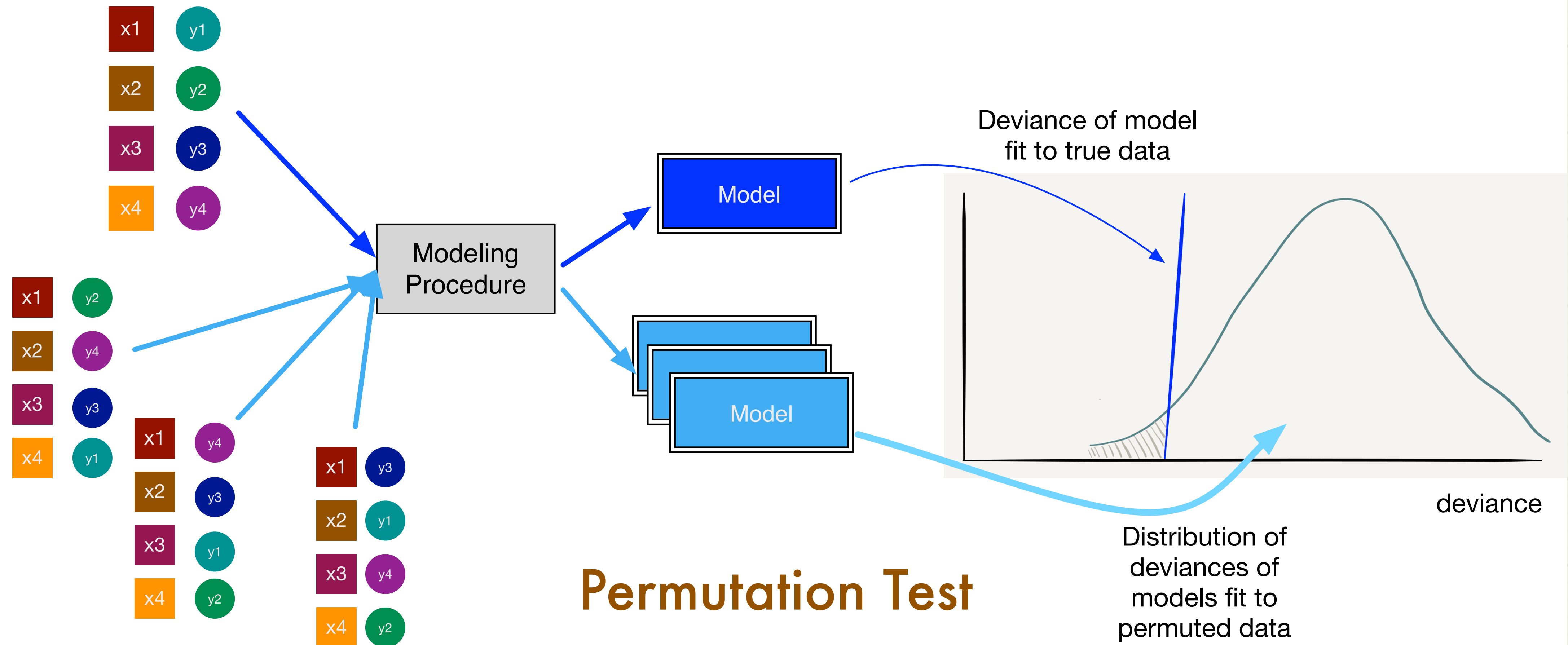
- Evaluation is comparison: distinguish good models from bad (or less good).
- Throughout: Estimate out-of-sample performance by repartitioning training data.
- Treat variable selection as choosing models by building throw-away single variable models.
- Work in probability units (not in effect sizes) wherever practical.

Are you predicting *anything* at all?

- Or: would a model fit to noise score this well?
- Or: Is the output *really* related to the inputs?



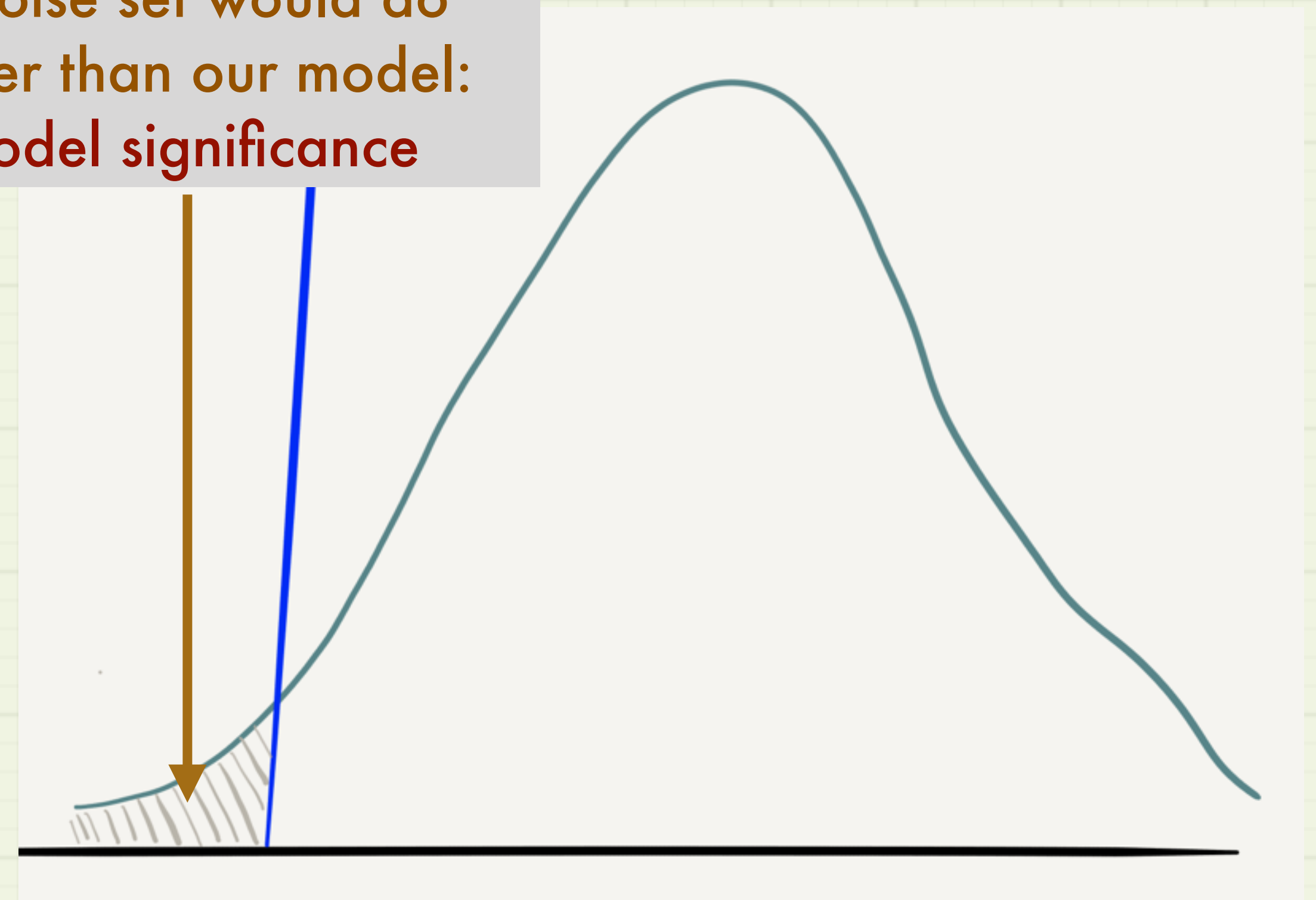
Thought Experiment: Is the input related to the output?



Model Significance

- In (data science) practice: hold-out set
- Useful for finding variables with signal
 - Fit a one-variable model
 - linear regression: F test
 - logistic regression: χ^2 test

Probability that model on a noise set would do better than our model:
model significance



Model Significance in R

Logistic Regression

```
# get the significance of glm model
get_glm_significance = function(model) {
  delta_deviance = model$null.deviance - model$deviance
  df = model$df.null - model$df.residual
  pchisq(delta_deviance, df, lower.tail=FALSE)
}
```

Linear Regression (from `summary(model)`)

```
# get the significance of lm model
get_lm_significance = function(model) {
  fs = summary(model)$fstatistic
  pf(fs["value"], fs["numdf"], fs["dendf"], lower.tail=FALSE)
}
```

labs/ Lab04PermTestVarSel

“Ok. I know the model is predicting something.
But how well is it doing?”

Can I trust my model evaluation?

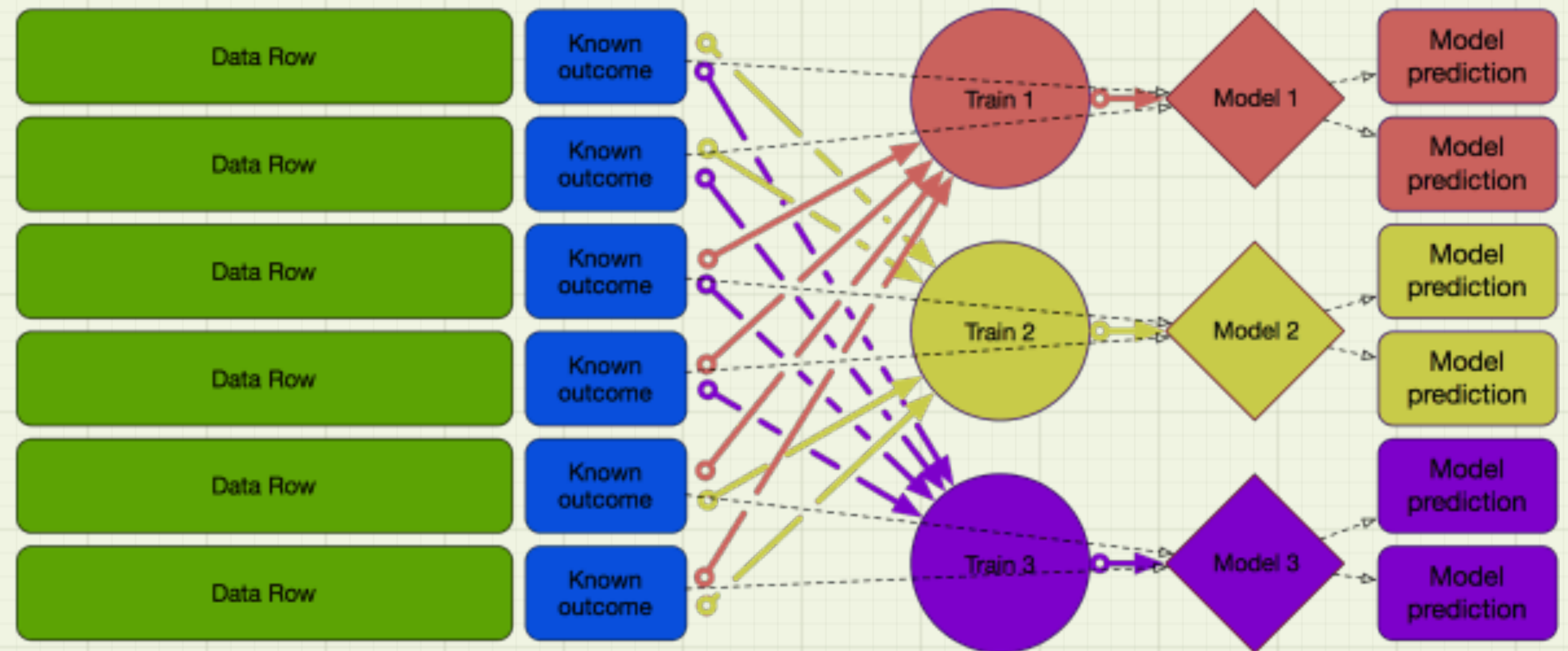
In-sample vs. Out-of-sample estimates

- Many common model metrics are in-sample
 - R^2 , model likelihood
 - RMSE or deviance on training set
 - **Performance estimates are upwardly biased**
- In-sample metrics that try to compensate for bias
 - Adjusted R^2 , AIC
 - *Only* consider bias due to model's degrees of freedom (# of parameters). Miss other multiple comparison issues that we don't tell them about!

Estimating Out-of-sample performance with training data

Cross-validation

- No point is evaluated on a model it helped to train



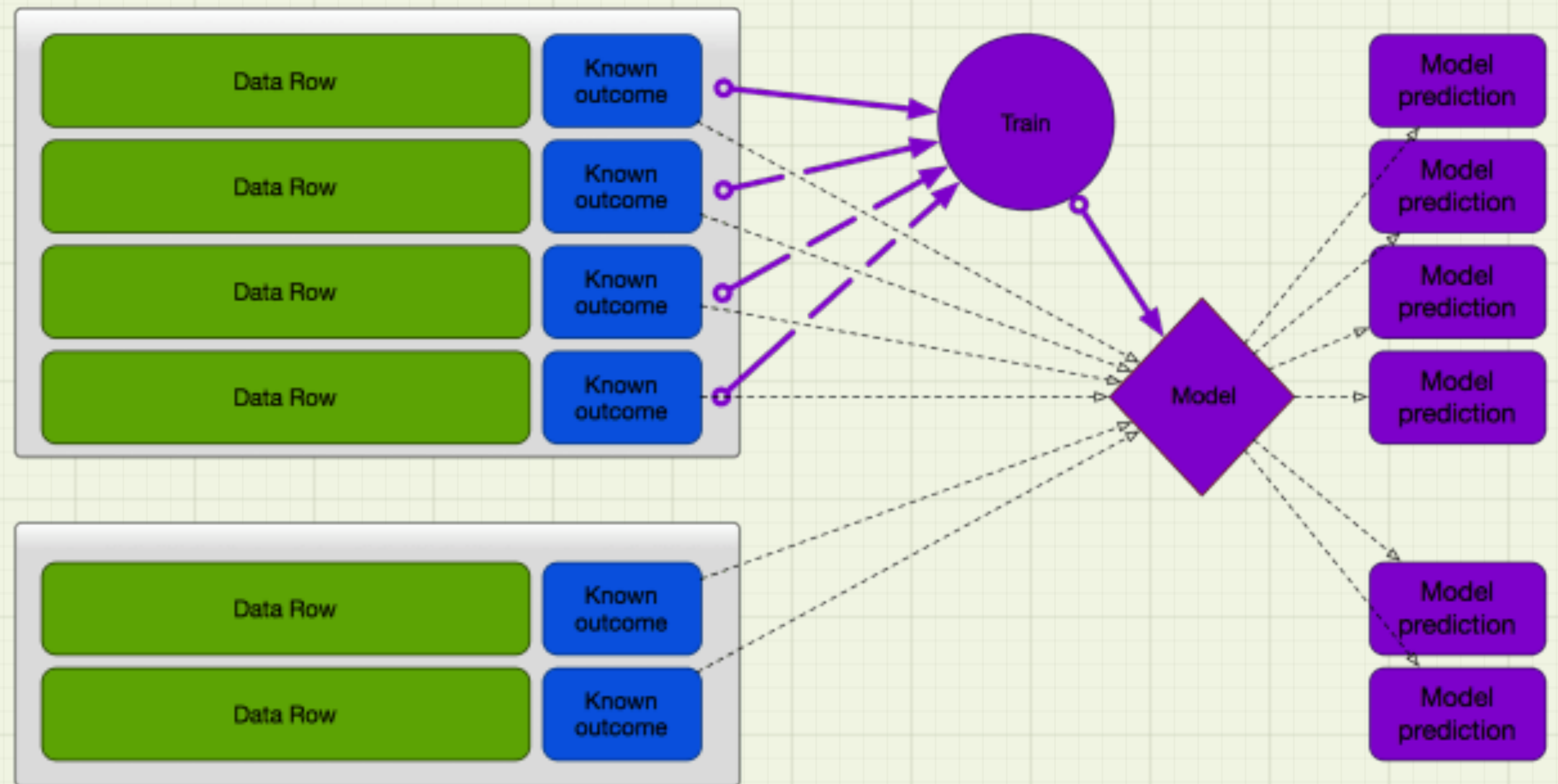
Cross-validation

- Unbiased estimate of out-of-sample error
- Statistically efficient
- Typically gives a point estimate of performance
- Computationally inefficient
 - With some exceptions: PRESS for linear models
- **Evaluates *modeling procedure* — NOT the production model.**

True estimate of out-of-sample performance: holdout data

Test-train split

- Subset of data only used for model evaluation



Test-train split

- Unbiased estimate of out-of-sample error
 - With no peeking!
- Statistically inefficient
 - Shouldn't split a small data set
- Point estimate only (in standard practice)
- Computationally efficient
- **Evaluates a specific model**

Cross-val vs. Holdout

	Statistically Efficient	Computationally Efficient	Evaluates Model	Evaluates Procedure
Cross-val				
Holdout				

- Mnemonic: individual researchers are Bayesians (want to know *their* model works), bosses are frequentists (want to know *the modeling process* works).

Data Science : Data-rich

- Generally, we will prefer test-train 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

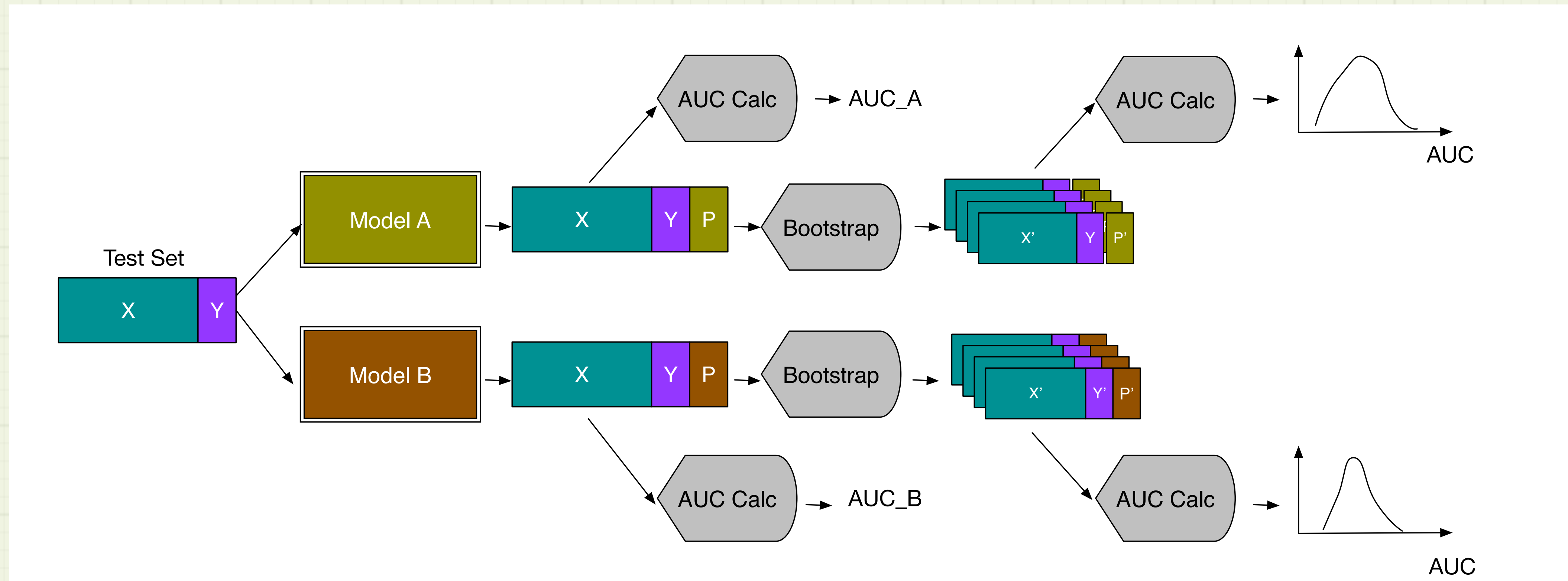
- Rare target or rare features of interest
- Setting modeling parameters
 - Especially if data set too small for test-train-calibration split
- Variable selection
- Small data sets

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 set is small
- Recent differential-privacy related results to alleviate this
 - <http://www.win-vector.com/blog/2015/10/a-simpler-explanation-of-differential-privacy/>

Is Model A decisively better
than Model B?
Distribution Estimates

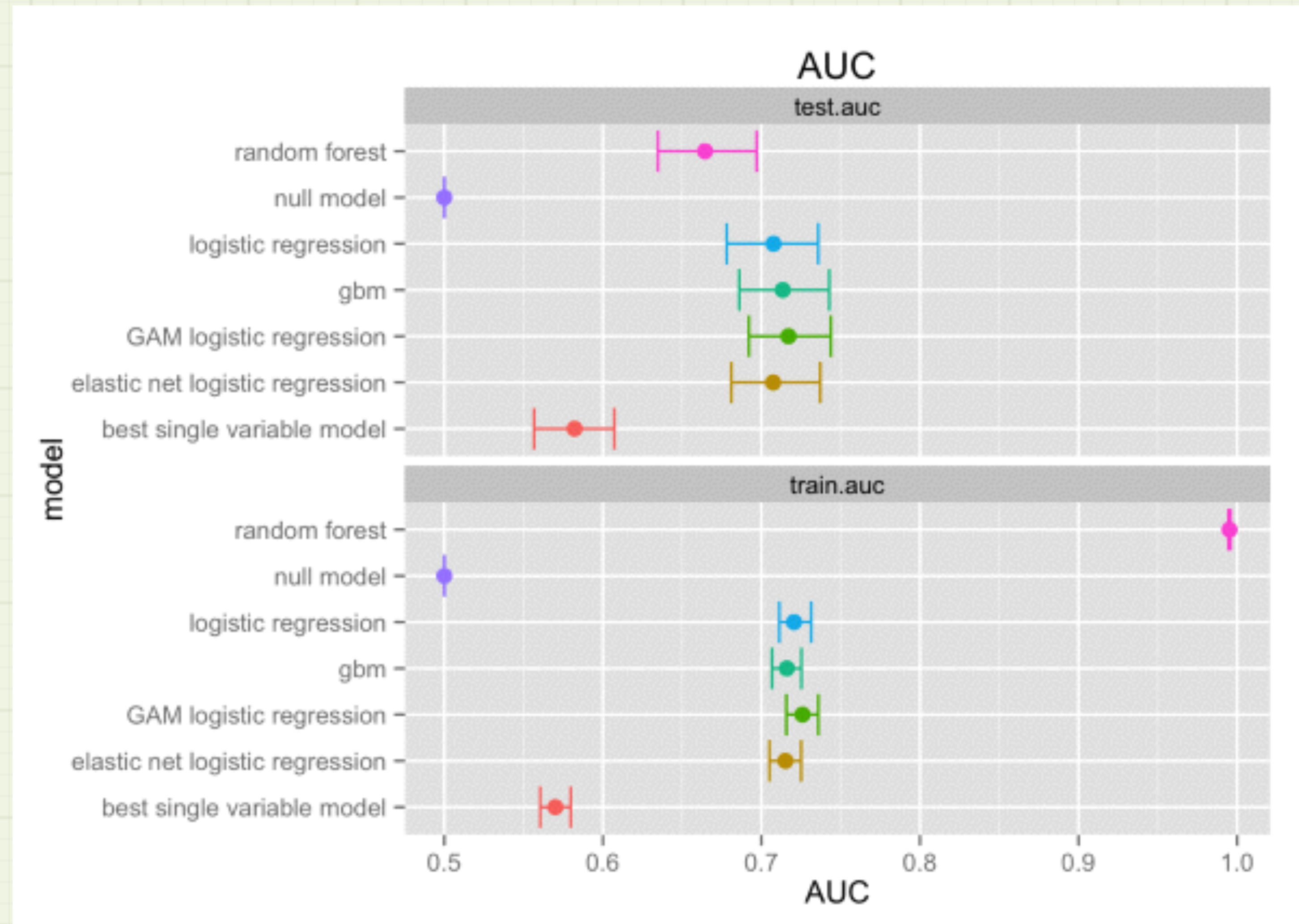
Bootstrapped Scoring



- One test-train split
- Build **one** model
- Bootstrap score the test (and training) sets

Bootstrapped Scoring

- Confidence Interval or Std. Dev. as distance unit
- Detect high variance
 - Situations the model gets really wrong (or really right)
- Does NOT measure stability of training procedure



labs/Lab05BootstrapTest

Conclusions / Take aways

- Model evaluation needs to be integrated into the entire data science process, including project proposal and result delivery.
- Metrics split between “technical” (for the data scientist, more appropriate early in a project) and “business” (adapted to intended use).
- R gives a flexible highly visual framework that implements classical statistical tests, important empirical tests, and any necessary ad-hoc procedures.

Some More Reading

- “How do you know if your model is going to work?”
 - <http://www.win-vector.com/blog/2015/09/isyourmodelgoingtowork/>
- “How Do You Know if Your Data Has Signal?”
 - <http://www.win-vector.com/blog/2015/08/how-do-you-know-if-your-data-has-signal/>
- “Statistics to English Translation”
 - **Part 1: Accuracy Measures**
 - **Part 2a: 'Significant' Doesn't Always Mean 'Important'**
 - **Part 2b: Calculating Significance**
- More testing in R:
 - **Finding the K in K-means by Parametric Bootstrap**
- Win-Vector LLC ODSC 2015 “Preparing Data workshop”
 - <https://github.com/WinVector/PreparingDataWorkshop/>
- White paper and video: Data Preparation in R (require registration)
 - **White Paper**
 - **Video (Pre-recorded webinar)**

Thank You

- Please stay in touch:
 - [@WinVectorLLC](#)
 - <http://www.win-vector.com/>
- Materials: <https://github.com/WinVector/ValidatingModelsInR>