# Decision Trees

**Practical Machine Learning (with R)**
UC Berkeley
Fall 2015

# Topics

⮌ Administrativa
  - Github
    - Reorganization → one repository
    - Please put

⮌ Review and Expectations

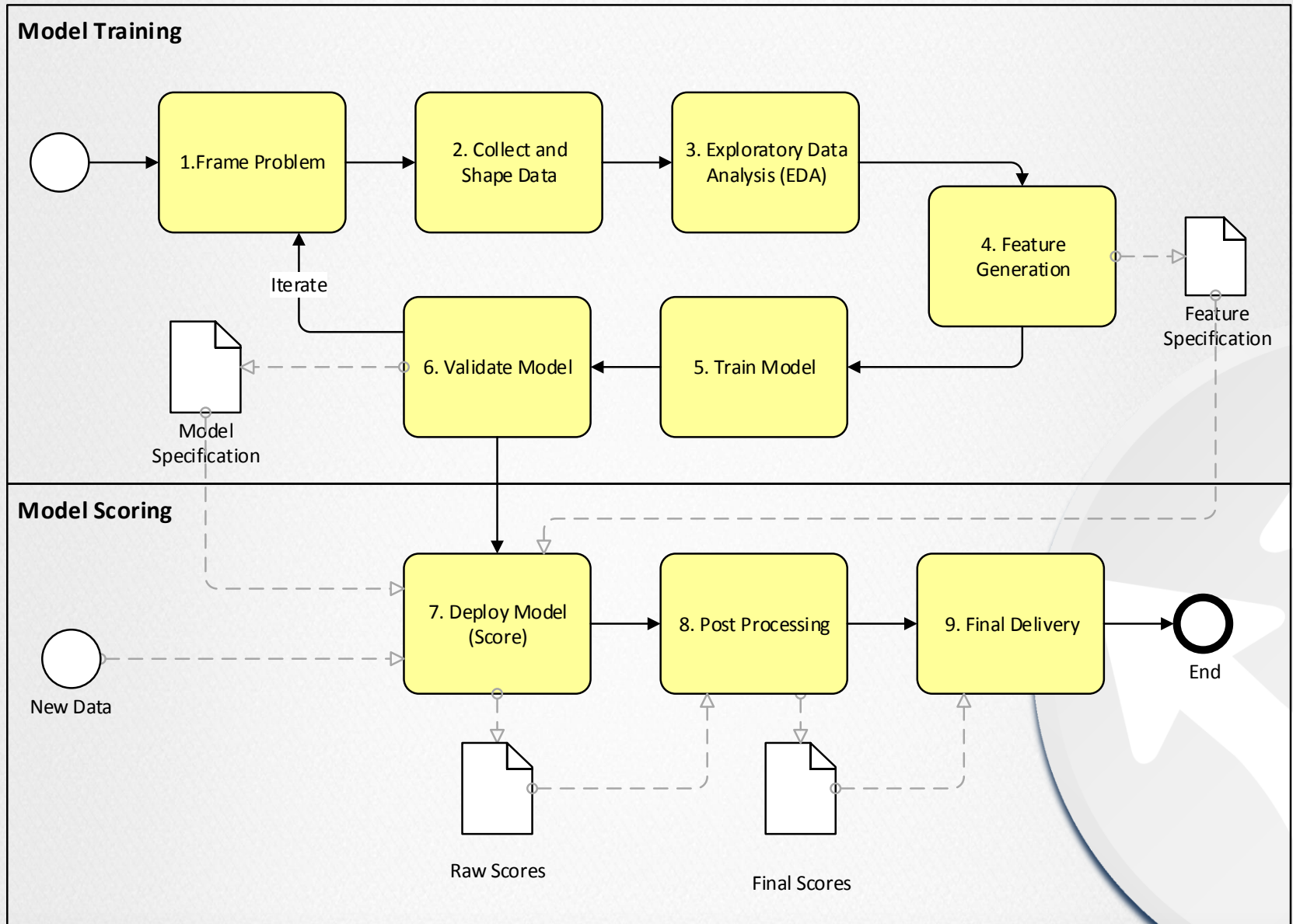⮌ In-Class Assignment

⮌ New Topics

# REVIEW AND EXPECTATIONS

# REVIEW AND EXPECTATION

- Github

- Transformations

- Logistic Regression ~ Linear Regression with the **logit** link function

- `glm( …, family=" ")`

# Comprehensive ML Process



**Model Training**

- 1. Frame Problem
- 2. Collect and Shape Data
- 3. Exploratory Data Analysis (EDA)
- 4. Feature Generation → Feature Specification
- 5. Train Model
- 6. Validate Model → Model Specification
- Iterate

**Model Scoring**

- New Data
- 7. Deploy Model (Score) → Raw Scores
- 8. Post Processing → Final Scores
- 9. Final Delivery
- End

Goal:

# BUILD UP A TOOL BOX OF SKILLS

# Worked Example: Boston Housing Transformations and Stepwise

# NEW TOPICS

# RMARKDOWN:DEMOSTRATION

# MODEL PERFORMANCE

# Model Performance

⮕ Determine relevant metric, e.g. **RMSE**, **FPR**
⮕ Calculate statistic ("metric")

⮕ On training data
*Training* or *apparent* performance → bias → over-fitting

**Need unbiased estimate for calculating performance**

# Resampling

⟳ Best Solution: Data Splitting
Split data into training and test data

- Easy to interpret defend
- Requires data not be consumed by model
- Computationally easy
- Is generally not (by itself) the most accurate → no confidence

⟳ Resampling Strategies

- Repeated Splitting
- K-Fold Cross Validation
- Bootstrap

# REPEATED SPLITTING

AKA Monte Carlo Splitting

⮱ Split data 75%-25%
- Fit Model
- Calculate Metric
- Repeat with Different Split(30+ times)

⮱ Calculate Metric

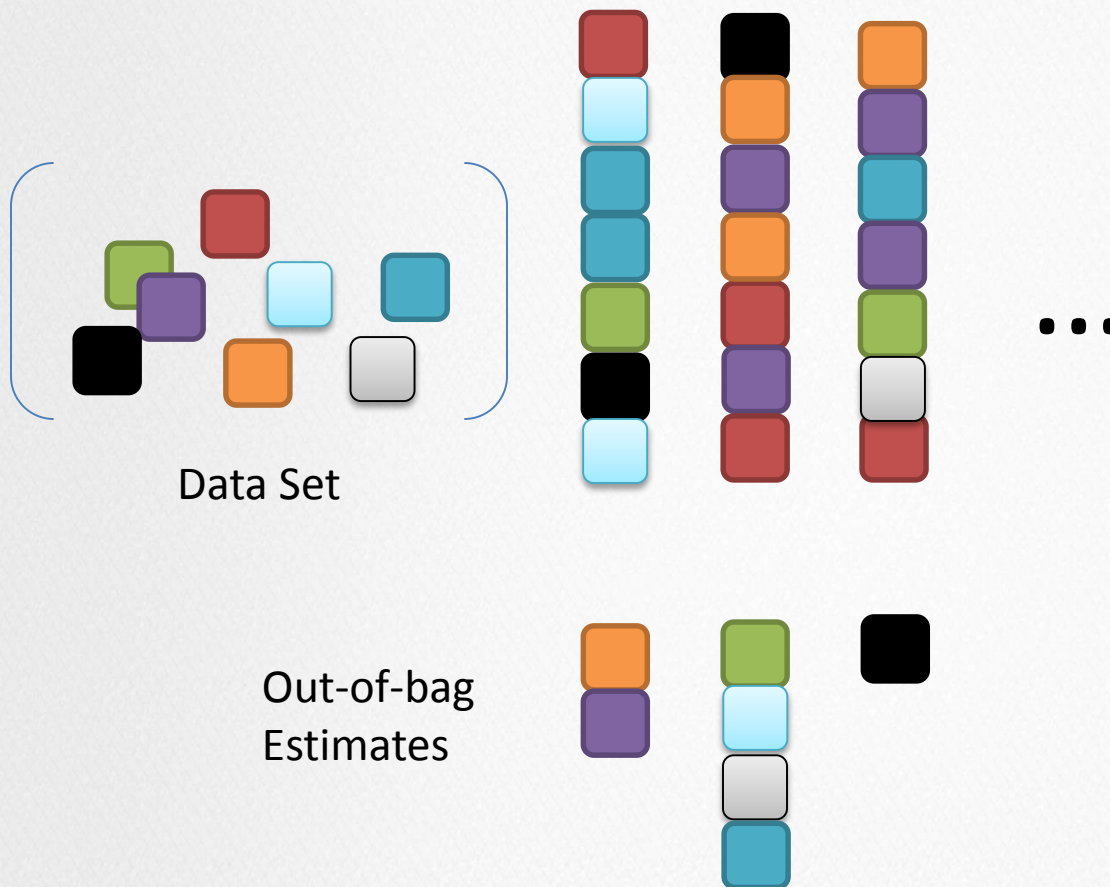$$Metric = AVG_i(metric)$$

# 10-Fold Cross Validation



LOOCV : K➔n

- Split the data set into 10 equal sized samples.
- Leave one sample out (fold)
  - Fit the model
  - calculate the metric on the fold
  - Repeat choosing another sample until

- Calculate Metric

$$Metric = AVG_i(metric)$$

- 5 or 10-fold common

# Bootstrap

⊃ "Sampling with Replacement"

Data Set

Out-of-bag
Estimates

# Which Is Best?



→ There isn't one.

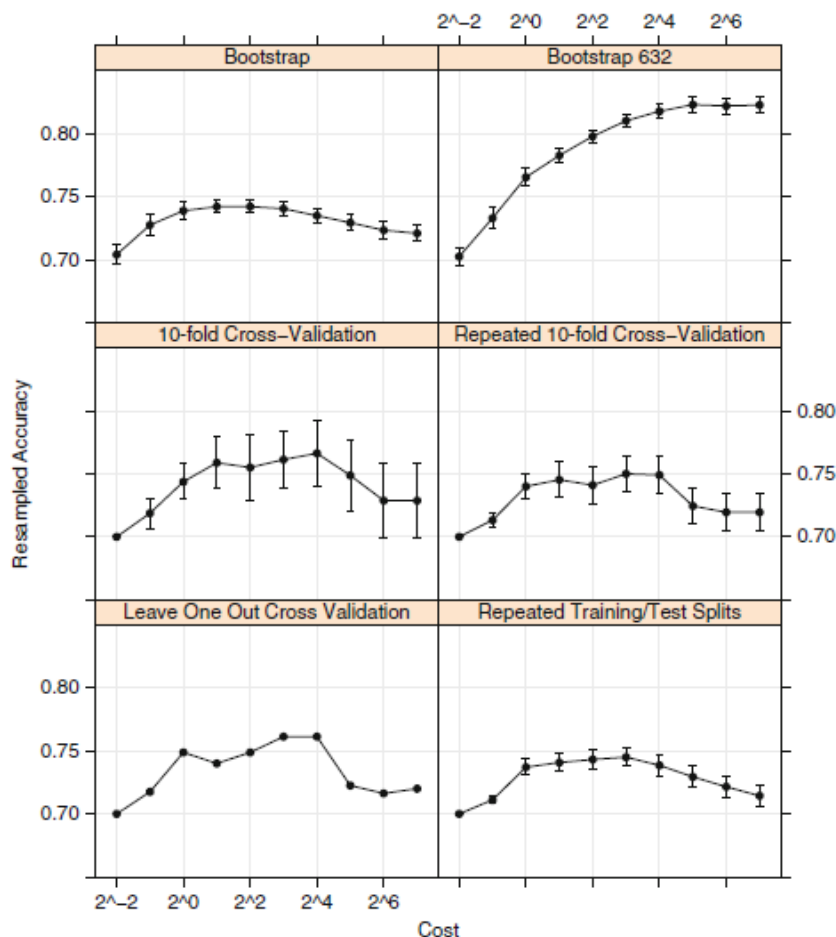K-fold cross validation
  Higher Variance
  Lower Bias

Bootstrap
  Lower Variance
  Higher Bias

**CALCULATING PERFORMANCE IS *NOT* THE SAME AS FITTING THE MODEL**

# EXERCISE: LECTURES/04-DECISION-TREES/RESAMPLING.RMD

KNOW YOUR DATA

# Model Formula (higher Order Terms)

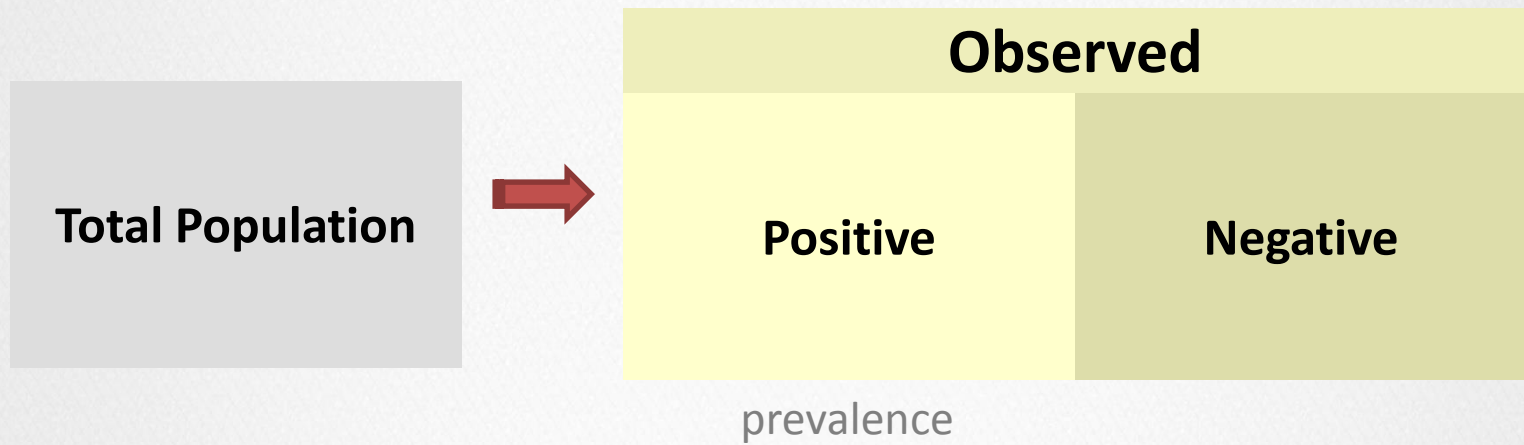➲ Model Formula …

# CLASSIFICATION PERFORMANCE

# METRICS FOR BI-NOMIAL CLASSIFICATION

**Total Population**

**Total Population** → **Observed**

| Positive | Negative |
|----------|----------|

prevalence

Total Population

Observed

Positive | Negative

Predicted

Positive

Negative

Accuracy

Error Rate
or Misclassification Rate

|  | | Observed | |
|---|---|---|---|
| **Total Population** | | **Positive** | **Negative** |
| **Predicted** | **Positive** | | |
| | **Negative** | | |

- https://en.wikipedia.org/wiki/Sensitivity_and_specificity

# Alternatives: Norm by Observed

| Total Population | Observed | |
|---|---|---|
| | **Positive** | **Negative** |
| **Predicted** / **Positive** | True Positive Rate (TPR), **Sensitivity**, Recall<br><br>$\dfrac{\text{True Positives}}{\textbf{Observed Positives}}$ | False Positive Rate (FPR), Fall-Out<br><br>$\dfrac{\text{False Positives}}{\textbf{Observed Negatives}}$ |
| **Negative** | False Neg. Rate (FNR), Miss rate<br><br>$\dfrac{\text{False Negatives}}{\textbf{Observed Positives}}$ | True Neg. Rate (TNR), **Specificity** (SPC)<br><br>$\dfrac{\text{True Negatives}}{\textbf{Observed Negatives}}$ |

# Alternatives: Norm by Predicted

| Total Population | → | **Observed** | |
|---|---|---|---|
| | | **Positive** | **Negative** |
| **Predicted** | **Positive** | Pos. Predictive Value (PPV), **Precision** <br><br> $\dfrac{\text{True Positives}}{\textbf{Predicted Positives}}$ | False Discovery Rate (FDR) <br><br> $\dfrac{\text{False Positives}}{\textbf{Predicted Positives}}$ |
| | **Negative** | False Omission Rate(FOR) <br><br> $\dfrac{\text{False Negatives}}{\textbf{Predicted Negatives}}$ | Negative Predictive Value (NPV) <br><br> $\dfrac{\text{True Negatives}}{\textbf{Predicted Negatives}}$ |

- https://en.wikipedia.org/wiki/Sensitivity_and_specificity

# MORE FUN …

[https://en.wikipedia.org/wiki/Sensitivity_and_specificity](https://en.wikipedia.org/wiki/Sensitivity_and_specificity)

# EXERCISE BINOMIAL METRICS: SKIN-NON SKIN

# Even More Complication

- Not all errors need count "*equivocal zone*" or "*intermediate zone*"

- *Prevalent when the model has three choices, e.g. A or B or Nothing.*

# MUTLINOMIAL CLASSIFICATION

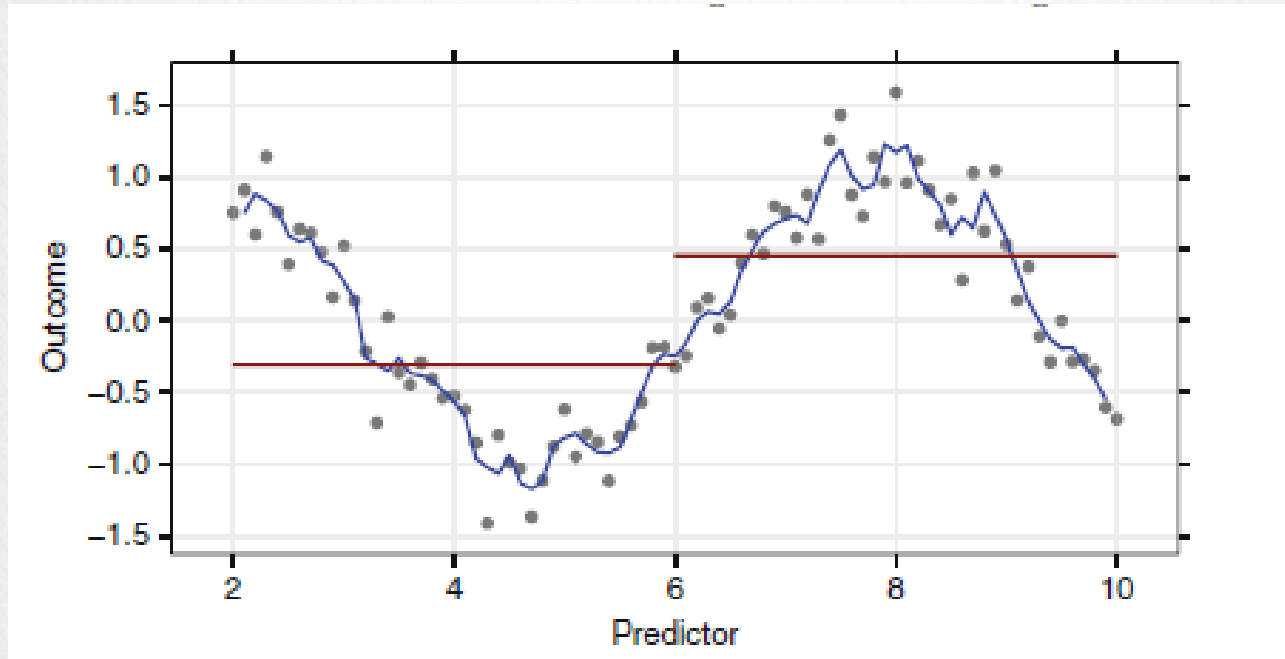# TERMS

- Kappa Statistic,
- S-Statistics, F-Statistic

# Multiclass Classification with Logistic Regression

# Bias Variance Trade-Off

$$E[MSE] = \sigma^2 + (model\ bias)^2 + model\ variance$$

# Process

# EXAMPLE OF ML ALGORITHM(S)

- Spam Filter
- handwriting recognition (svm)
- Traffic engineering (lights)
- Weather prediction
- Sentiment analysis (social media)
- Netflix Recommender
- Fraud detection (Visa)
- Imaging processing
- (network) Intrution detection
- Self-driving cars

# LOGISTIC REGRESSION

# DECISION TREES

# APPENDIX

# Comparison of Models (Chart)

# TRANSFORMATIONS

- Centering and Scaling: `scale`*
- Resolve skewness: `log, sqrt, inv`
- Resolve outliers: spatial sign, `PCA`

Some algorithms require scaling

Some are insensitive

Time consuming

Somewhat of an art

- Genetic algorithms (GA)

| | | True condition | | | |
|---|---|---|---|---|---|
| **Total population** | | Condition positive | Condition negative | Prevalence = Σ Condition positive/Σ Total population | |
| **Predicted condition** | Predicted condition positive | **True positive** | **False positive** (Type I error) | Positive predictive value (PPV), Precision = Σ True positive/Σ Test outcome positive | False discovery rate (FDR) = Σ False positive/Σ Test outcome positive |
| | Predicted condition negative | **False negative** (Type II error) | **True negative** | False omission rate (FOR) = Σ False negative/Σ Test outcome negative | Negative predictive value (NPV) = Σ True negative/Σ Test outcome negative |
| | Accuracy (ACC) = Σ True positive + Σ True negative/Σ Total population | True positive rate (TPR), Sensitivity, Recall = Σ True positive/Σ Condition positive | False positive rate (FPR), Fall-out = Σ False positive/Σ Condition negative | Positive likelihood ratio (LR+) = TPR/FPR | Diagnostic odds ratio (DOR) = LR+/LR− |
| | | False negative rate (FNR), Miss rate = Σ False negative/Σ Condition positive | True negative rate (TNR), Specificity (SPC) = Σ True negative/Σ Condition negative | Negative likelihood ratio (LR−) = FNR/TNR | |