

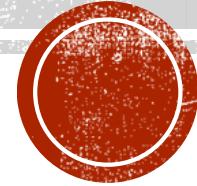
LIMITS OF LEARNING

CS 412 Introduction to Machine Learning

Prof. Zheleva

January 25, 2018

Reading Assignment: CIML: 2, ISL: 2.1-2.2.2



LAST LECTURE: DECISION TREES

- **What is a decision tree?**
 - **Nodes:** test the value of feature x_j
 - **Branches:** correspond to values
 - **Leafs:** provide the class (prediction)
- **How to learn a decision tree from data?**
 - Finding the best tree is NP-hard
 - Scoring functions: Training error, Entropy, Information Gain, Gini Impurity
- **What is the inductive bias?**

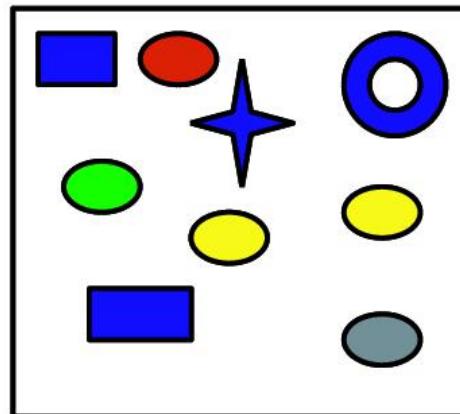


Toy dataset and classification problem

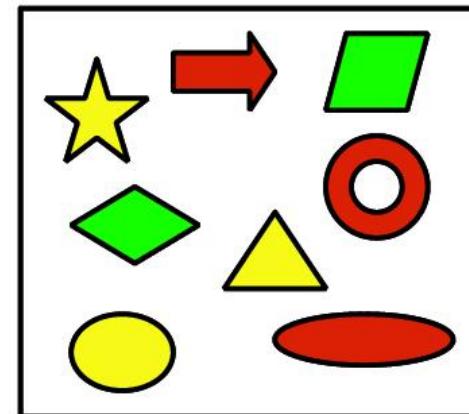
Are these objects from outer space?



yes



no



D features (attributes)

Color	Shape	Size (cm)
Blue	Square	10
Red	Ellipse	2.4
Red	Ellipse	20.7

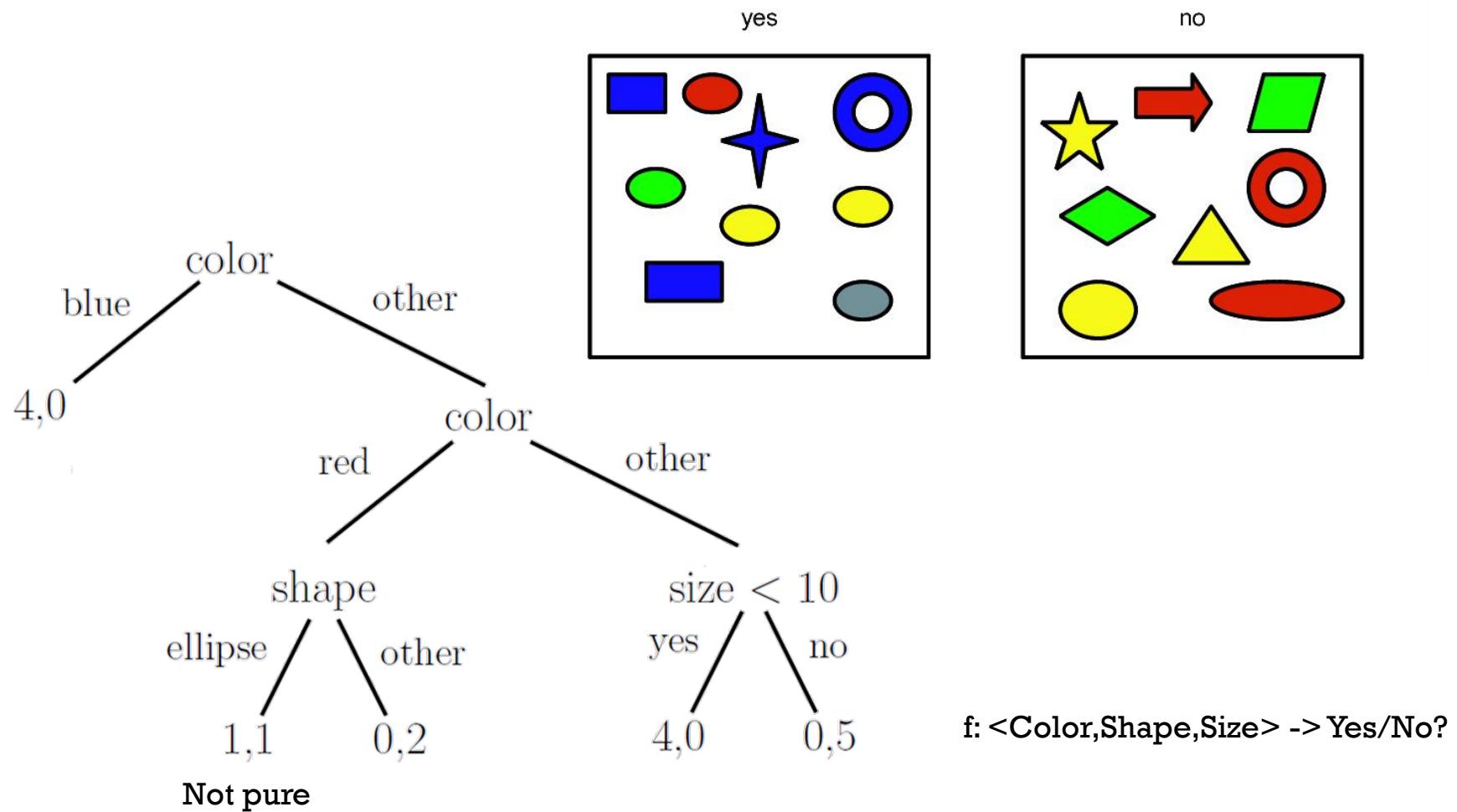
Label

1

1

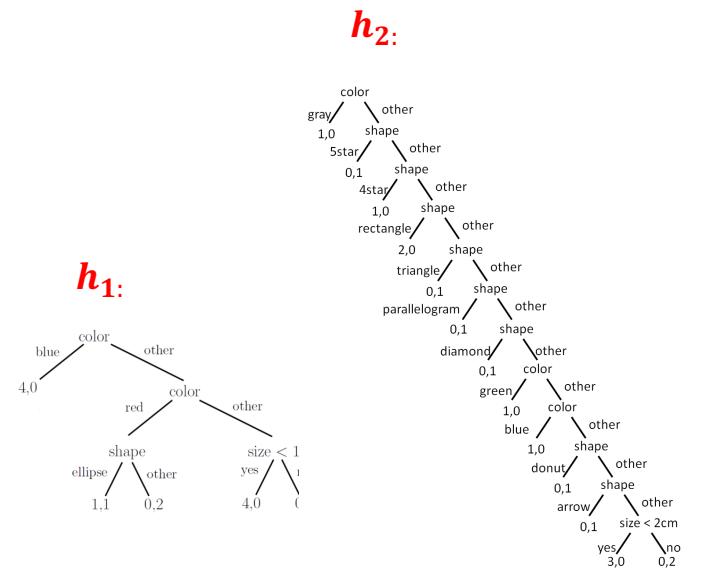
0

N cases



RECAP: SUPERVISED LEARNING SETTING

- Problem setting
 - \mathbf{X} – set of possible instances
 - Each instance $x \in \mathbf{X}$ is a feature vector $x = [x_1, \dots, x_D]$
 - Unknown target function $f: \mathbf{X} \rightarrow \mathbf{Y}$
 - Classification: \mathbf{Y} is discrete valued
 - Set of function hypotheses $H = \{h \mid h: \mathbf{X} \rightarrow \mathbf{Y}\}$
 - Each hypothesis h is a decision tree
- Input
 - Training examples $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ of unknown distribution
- Output
 - Hypothesis $h \in H$ that best approximates target function f



TODAY: LIMITS OF LEARNING

- Decision trees
 - What is the inductive bias?
 - Generalization
- Practical concerns:
 - How to estimate error on unseen examples?
 - Dealing with data: from raw data to well-defined examples
 - (Concerns apply to all models, not just decision trees)



NO FREE LUNCH

“All models are wrong, but some are useful.”

George Box, Statistician

- Modeling assumptions that work well for one problem, may not work well for another
- No “universal learner” exists
 - We need many different techniques for different domains and task characteristics

BAYES OPTIMAL CLASSIFIER

- What if we had full access to the underlying data distribution D over (\mathbf{x}, \mathbf{y}) pairs?
- **Bayes optimal classifier**
 - For any input $\hat{\mathbf{x}}$, we can compute the label \hat{y}

$$f^{(\text{BO})}(\hat{\mathbf{x}}) = \arg \max_{\hat{y} \in \mathcal{Y}} \mathcal{D}(\hat{\mathbf{x}}, \hat{y})$$

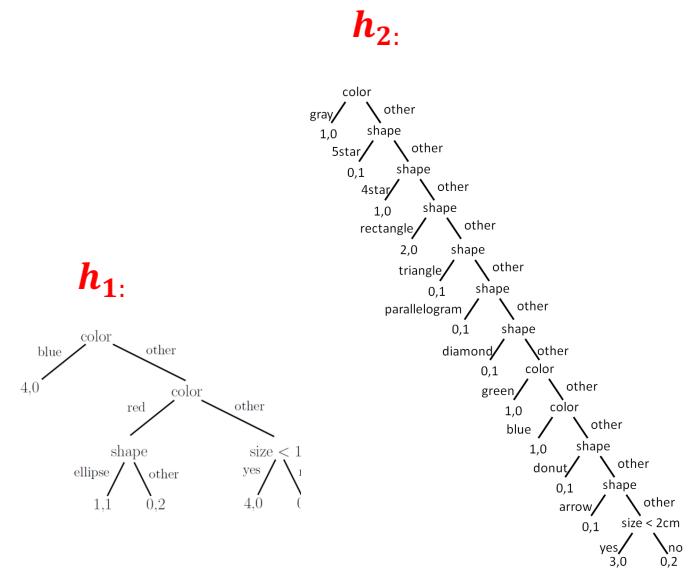
Theorem 1 (Bayes Optimal Classifier). *The Bayes Optimal Classifier $f^{(\text{BO})}$ achieves minimal zero/one error of any deterministic classifier.*

- Bayes error rate: error of Bayes optimal classifier
- Unfortunately, we never have access to D – so what do we do?



INDUCTIVE BIAS IN DECISION TREES

- Our learning algorithm performs heuristic search through space of decision trees
- It stops at smallest acceptable tree
- Why do we prefer small trees?
 - Occam's razor: prefer the simplest hypothesis that fits the data



WHY PREFER SHORT HYPOTHESES?

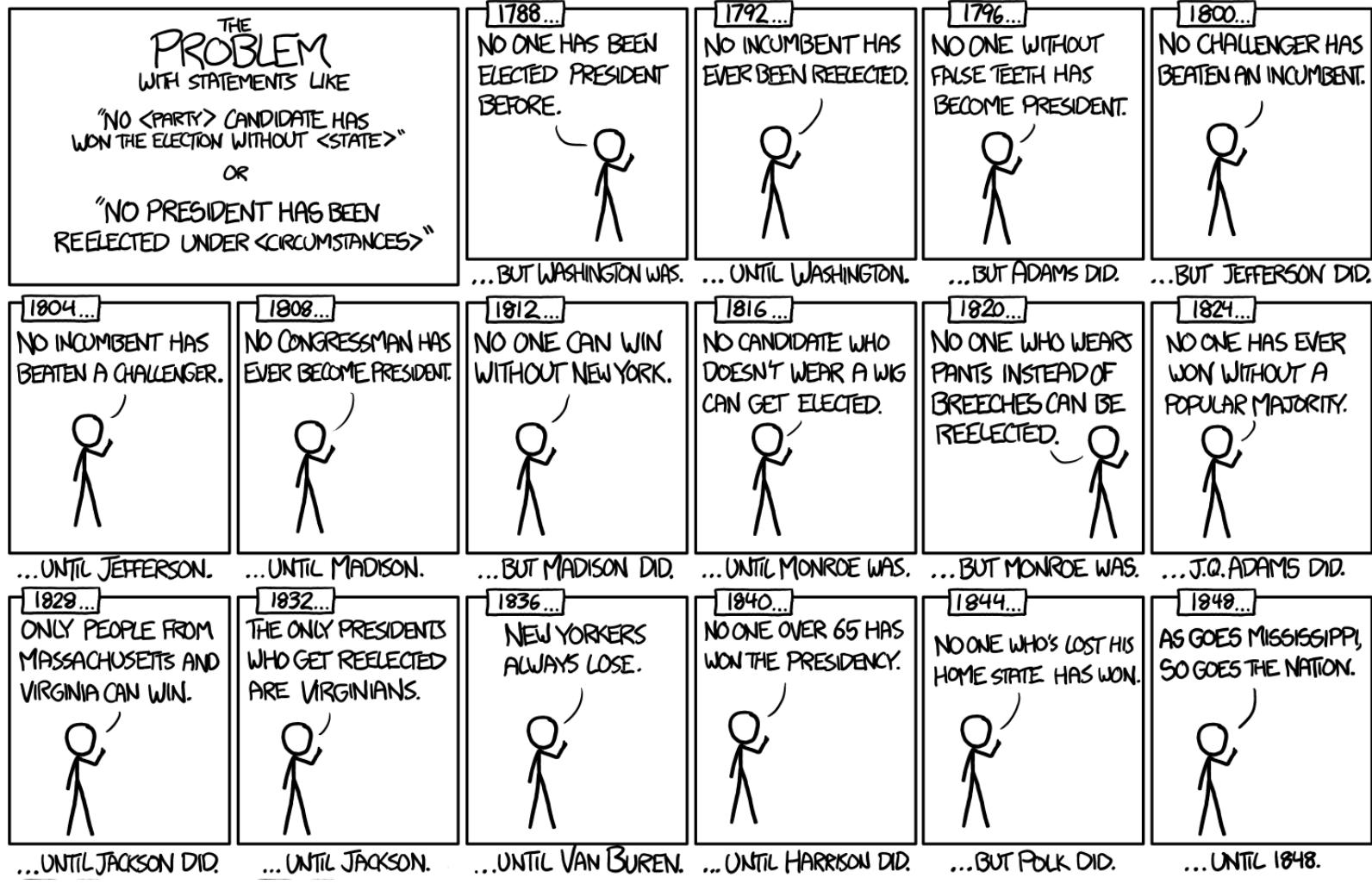
- **Pros**

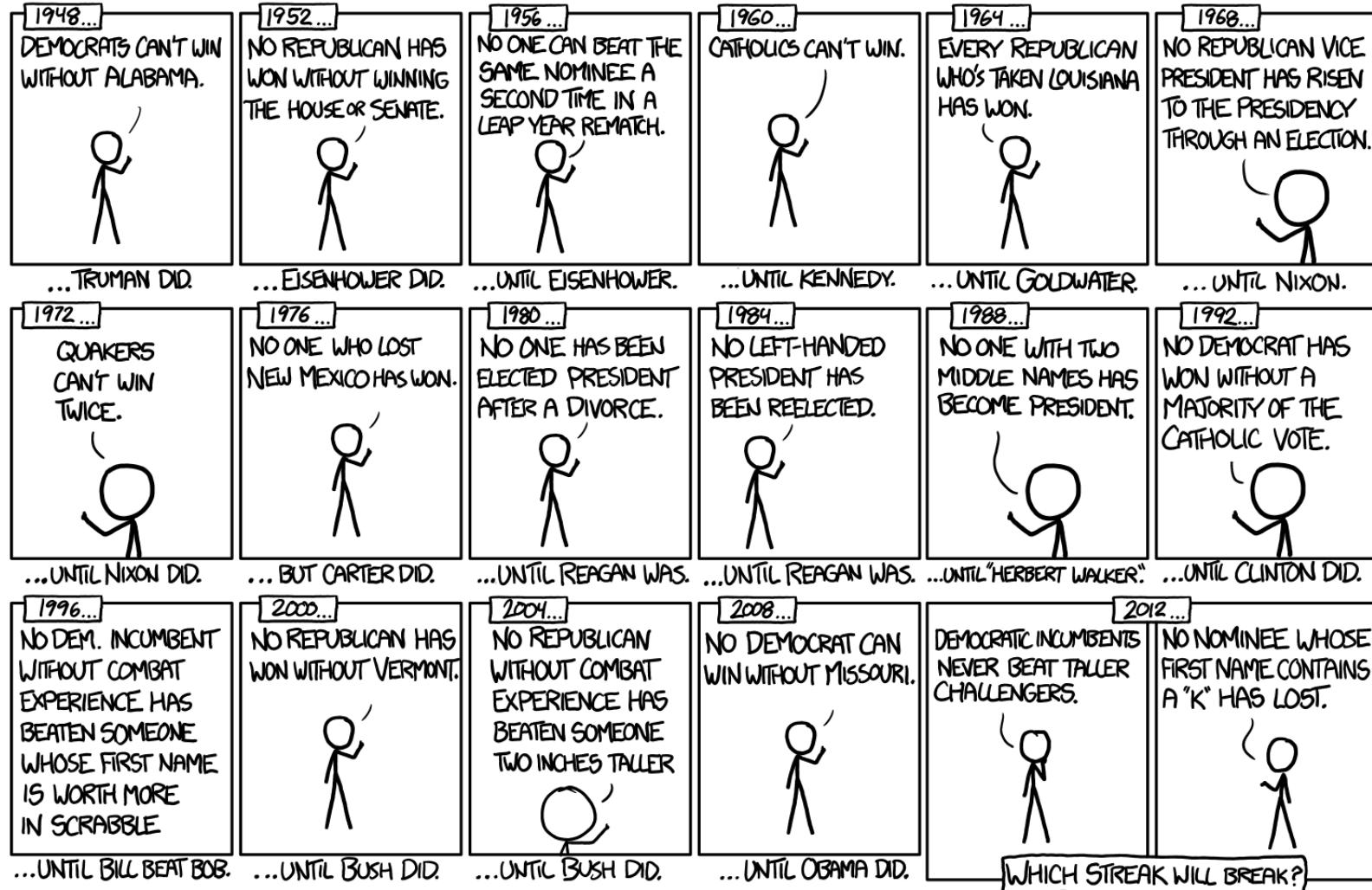
- Fewer short hypotheses than long ones
 - A short hypothesis that fits the data is less likely to be a statistical coincidence

- **Cons**

- What's so special about short hypotheses?





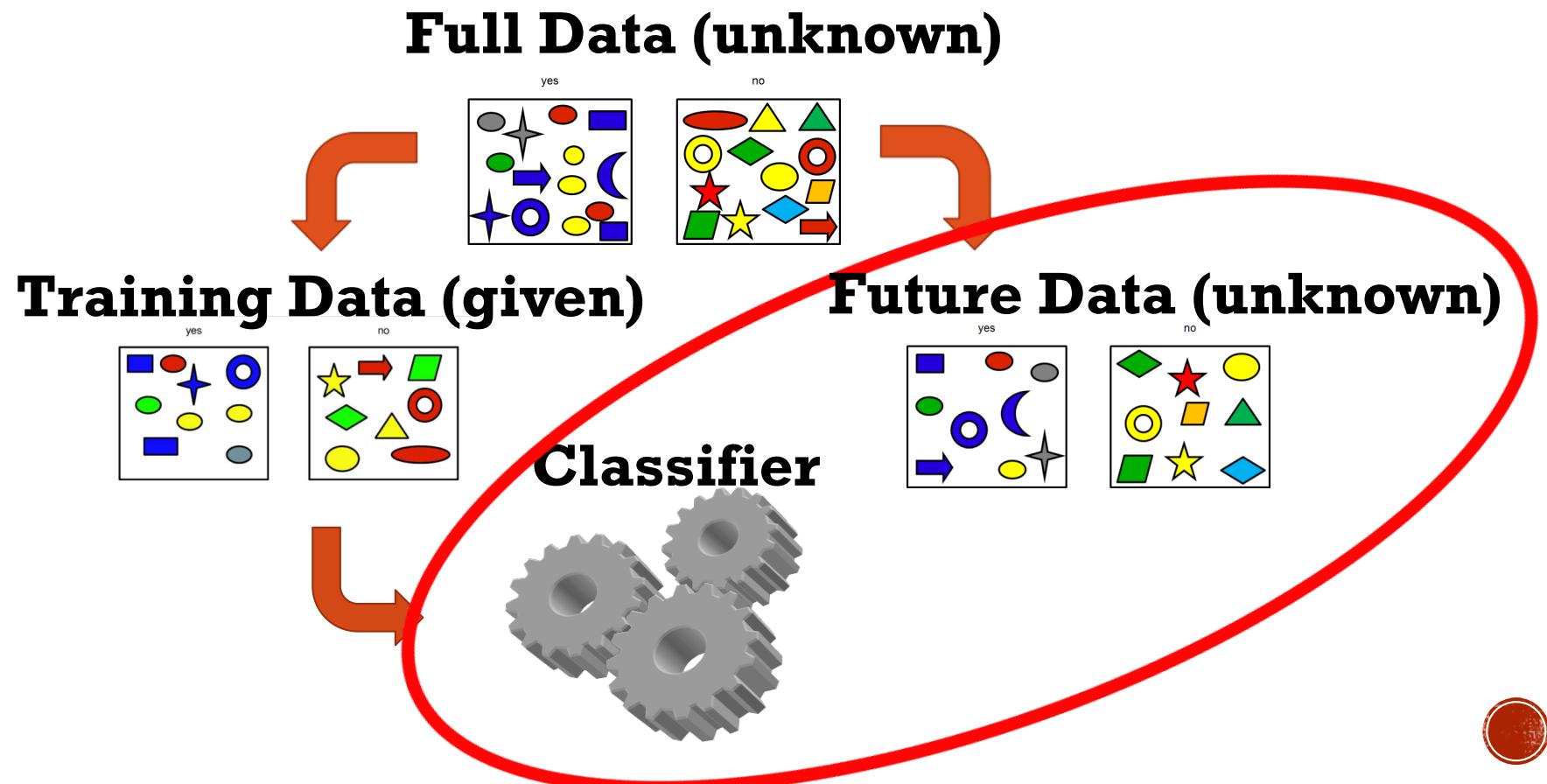


EVALUATING THE LEARNED HYPOTHESIS

- Assume
 - We've learned a tree h using the top-down induction algorithm
 - It fits the data perfectly
- Are we done? Can we guarantee we have found a good hypothesis?
- Ideally, we want a hypothesis that can explain the past but is also simple enough to generalize to the future



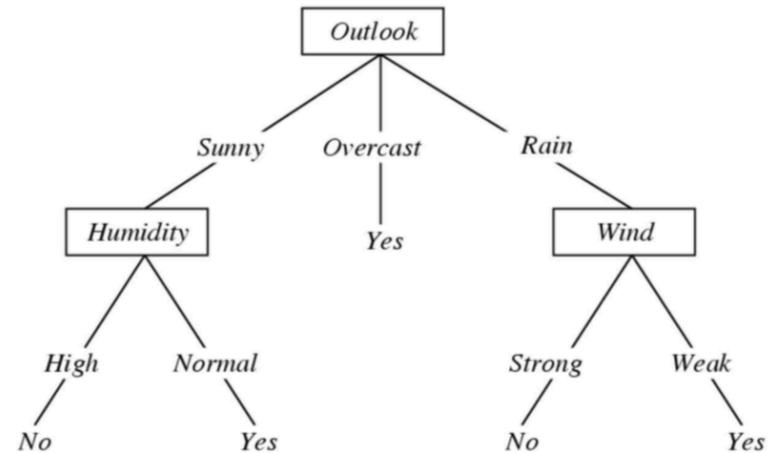
SUPERVISED LEARNING



SHOULD WE PLAY TENNIS?

Day Outlook Temperature Humidity Wind PlayTennis?

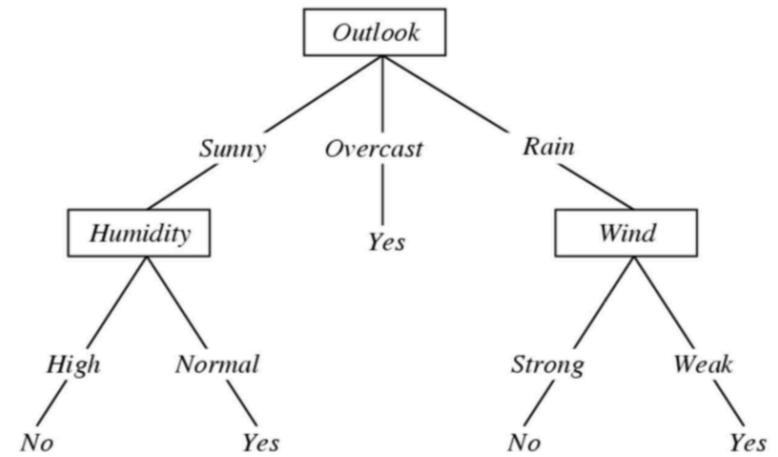
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



LET'S ADD A NOISY EXAMPLE

Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No
D15	Sunny	Hot	Normal	Strong	No

How does this affect the decision tree?



RECALL: FORMALIZING INDUCTION

- $f(\mathbf{x})$ should make good predictions
 - As measured by loss function
 - On future examples (typically) drawn from D
- Given
 - A loss function
 - A sample from some unknown data distribution $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
- Compute a function f that has a low expected error over D with respect to the loss

$$\varepsilon \triangleq \mathbb{E}_{(x,y) \sim D} \{l(y, f(x))\} = \sum_{(x,y)} D(x,y) l(y, f(x))$$



TRAINING ERROR IS NOT SUFFICIENT

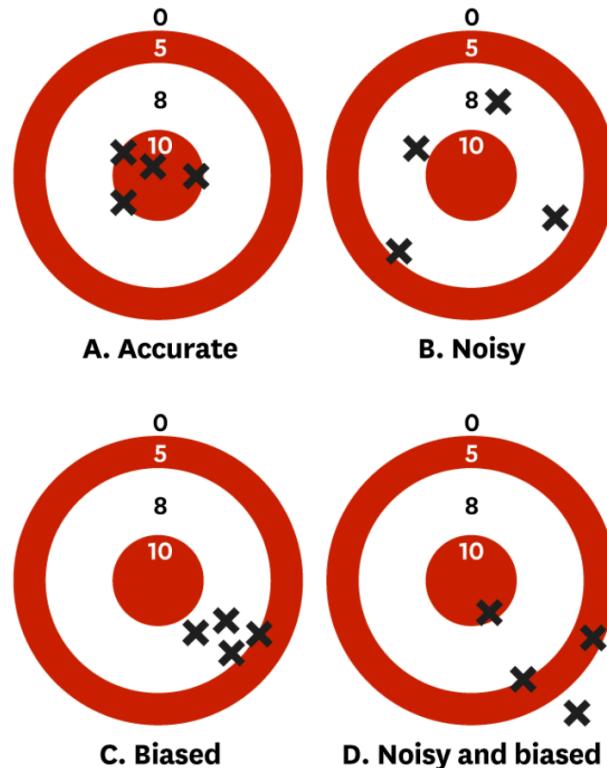
- We care about **generalization** to new examples
- A tree can classify training data perfectly, yet classify new examples incorrectly
 - Because **training examples are only a sample** of data distribution
 - a feature might correlate with class by coincidence
 - Because **training examples could be noisy**
 - e.g., accident in labeling



NOISE VS. BIAS IN TRAINING DATA

- Goal: estimate target in dart-throwing
- Training data: four sets of four examples
- Noise = variance

How Noise and Bias Affect Accuracy



SOURCE DANIEL KAHNEMAN,
ANDREW M. ROSENFIELD,
LINNEA GANDHI, AND TOM BLASER
FROM "NOISE," OCTOBER 2016

© HBR.ORG

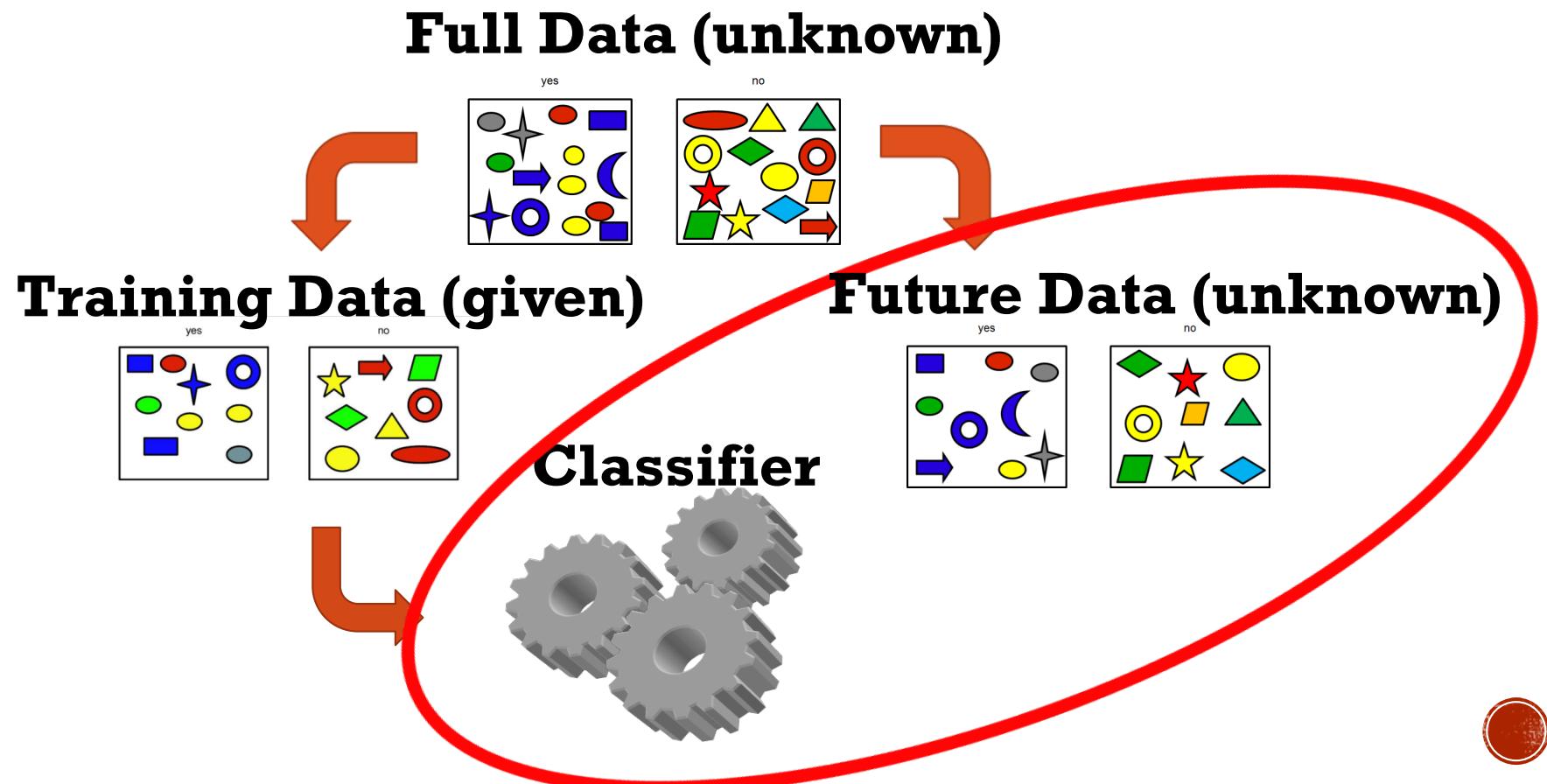
Example from <https://hbr.org/2016/10/noise>

NEXT: ERROR ON UNSEEN EXAMPLES

- Decision trees
 - What is the inductive bias?
 - Generalization
- **Practical concerns**
 - How to estimate error on unseen examples?
 - Dealing with data: from raw data to well-defined examples

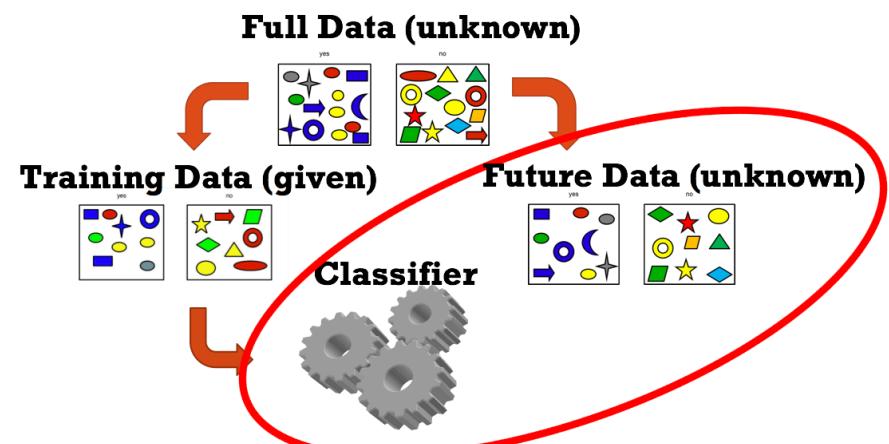


SUPERVISED LEARNING



OVERFITTING

- Consider a hypothesis h and its:
 - Error rate over training data $\text{error}_{\text{train}}(h)$
 - True error rate over full data $\text{error}_{\text{true}}(h)$
- We say that h overfits the training data if
 - $\text{error}_{\text{true}}(h) > \text{error}_{\text{train}}(h)$
- Amount of overfitting
 - $\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)$

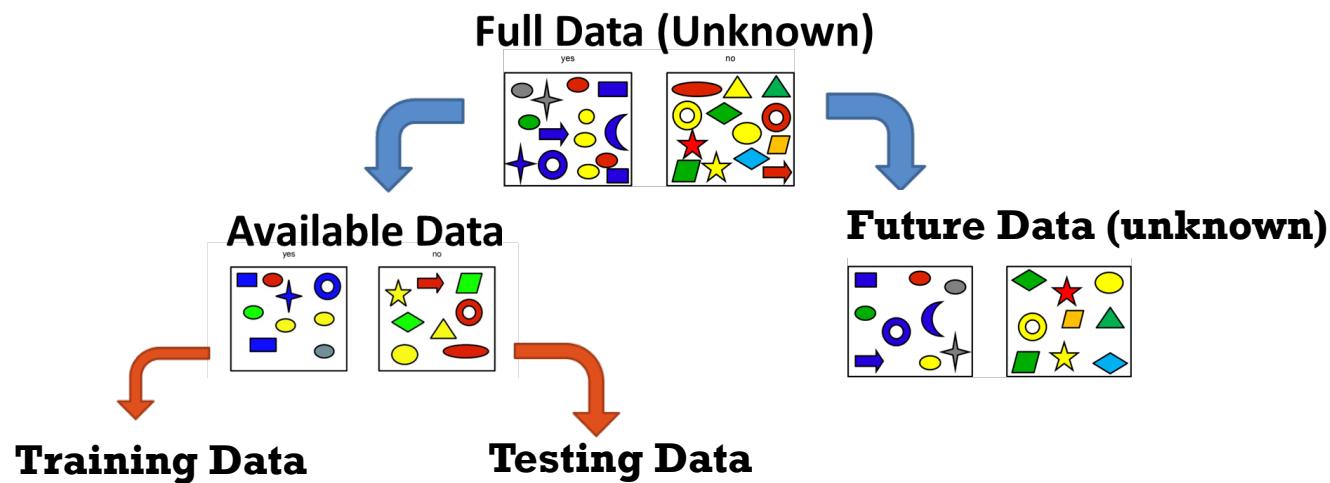


HOW DO WE EVALUATE PERFORMANCE?

- **Problem**
 - we don't know $error_{true}(h)$!
- **How do we evaluate then?**
 - we set aside a testing set from the available data
 - some examples that will be used for evaluation
 - we don't look at them during training! (pretend we have not seen them)
 - after learning a model, we calculate $error_{test}(h)$

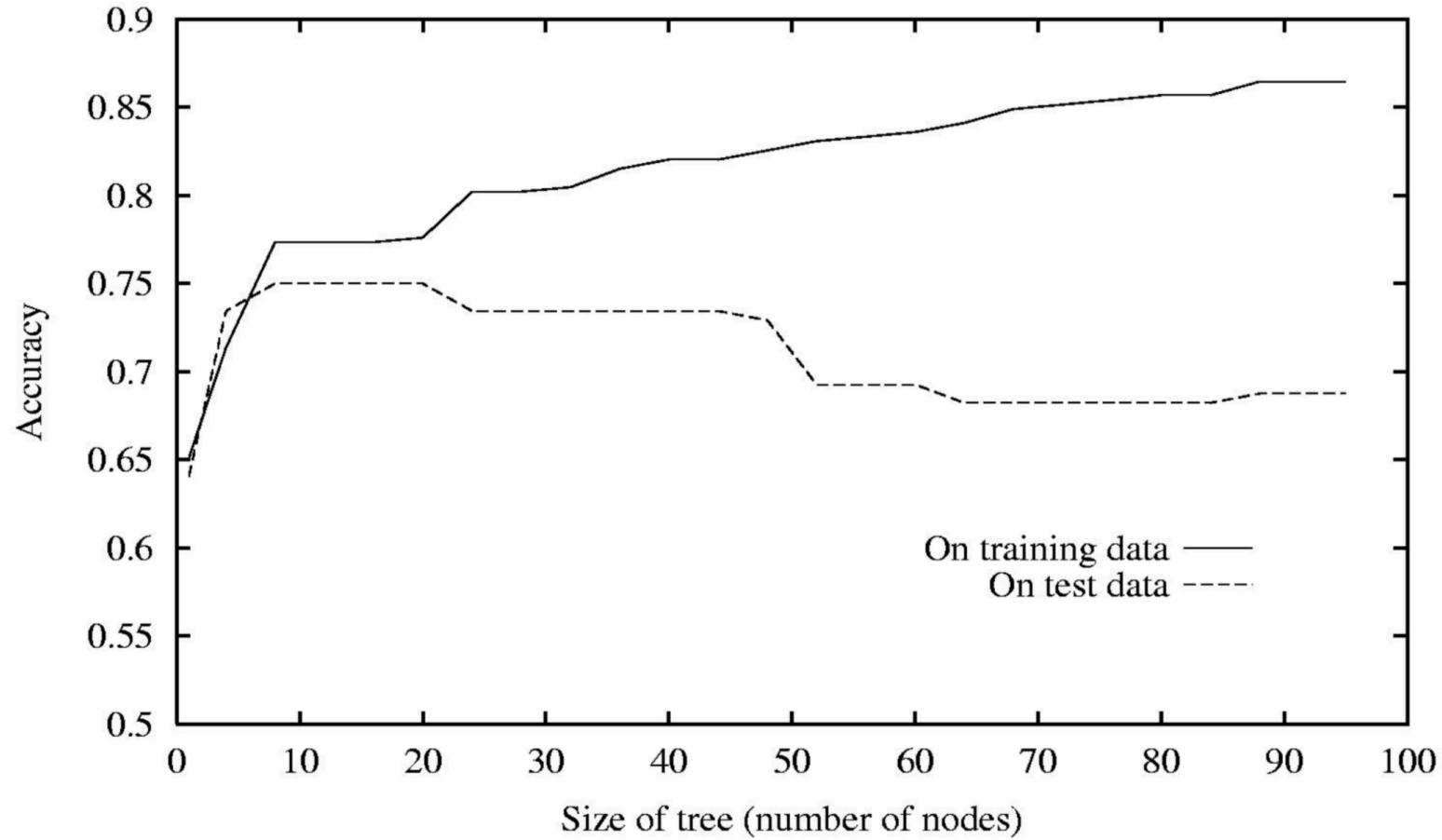


SUPERVISED LEARNING EVALUATION

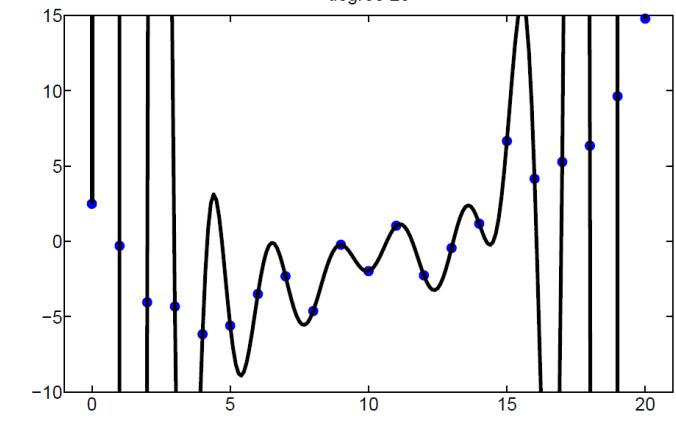
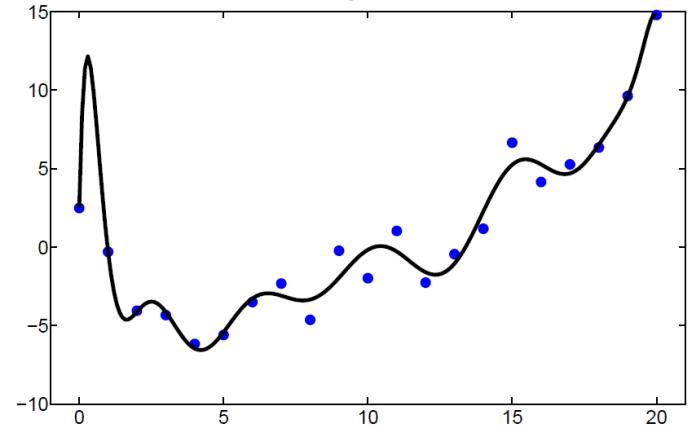
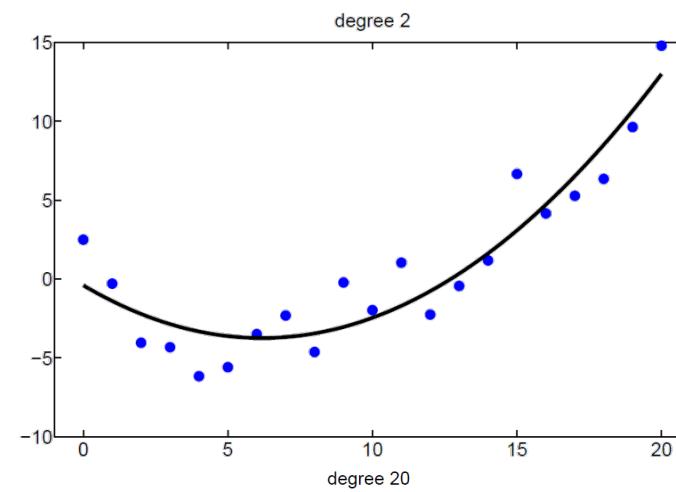
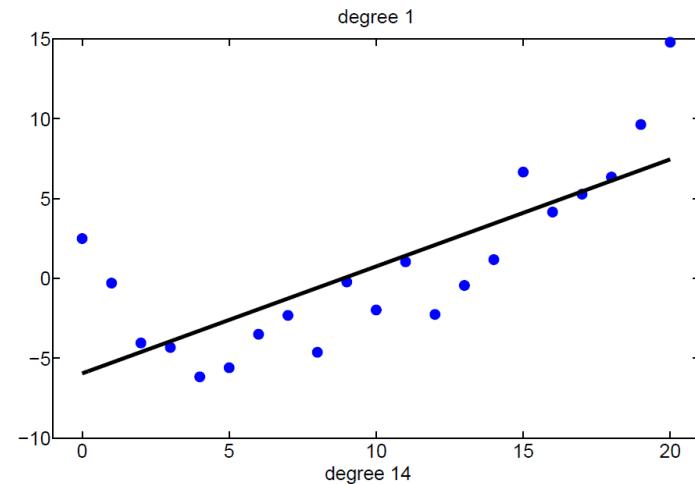


We use testing data for evaluation
We don't look at testing data during training!

MEASURING OVERFITTING EFFECT



OVERFITTING IN REGRESSION



OVERFITTING

- We say that h overfits the training data if
 - $error_{true}(h) > error_{train}(h)$
- Another way of putting it: A hypothesis h is said to overfit the training data, if there is another hypothesis h' , such that
 - h has a smaller error than h' on the training data
 - but h has larger error on the test data than h' .



UNDERFITTING VS. OVERFITTING

- **Underfitting**

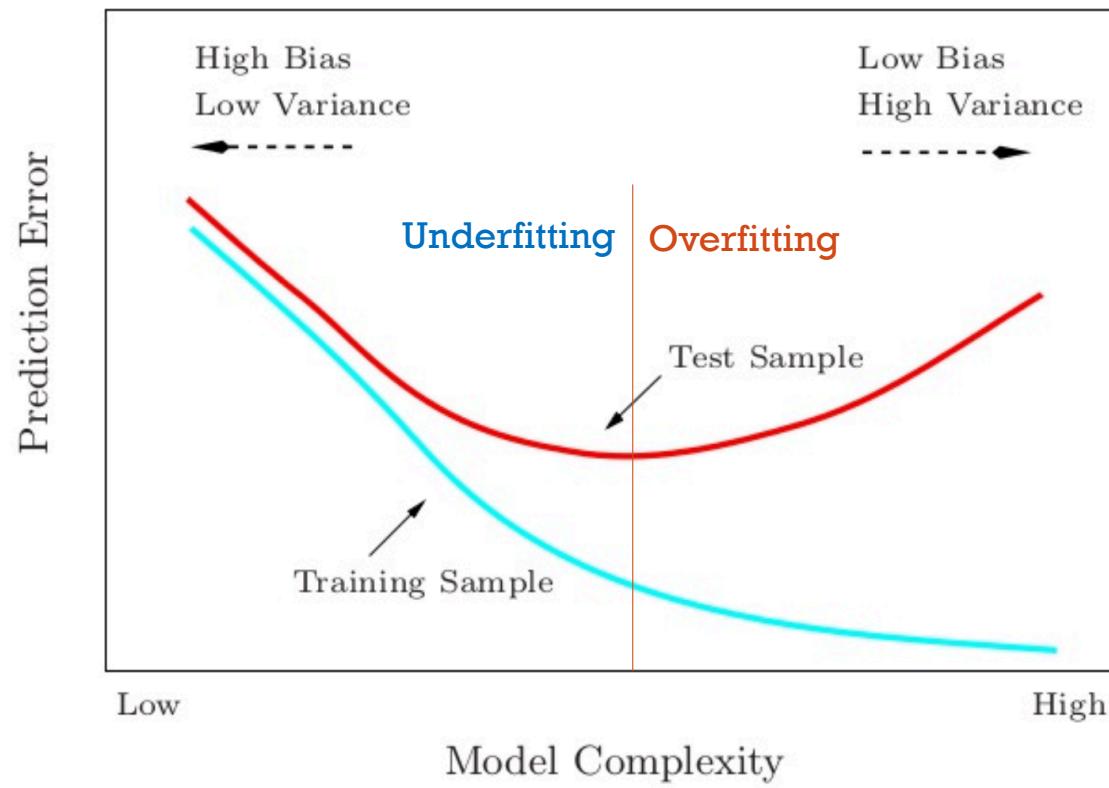
- Learning algorithm had the opportunity to learn more from training data, but didn't

- **Overfitting**

- Learning algorithm paid too much attention to idiosyncracies of the training data; the resulting tree doesn't generalize



UNDERFITTING VS. OVERFITTING



PRACTICAL IMPACT

- What we want:
 - A decision tree that neither underfits nor overfits
 - Because it is expected to do best in the future
- How can we encourage that behavior?
 - Set a maximum tree depth D
- How do we learn this maximum tree depth D?
 - Hint: We can't use testing data during training!

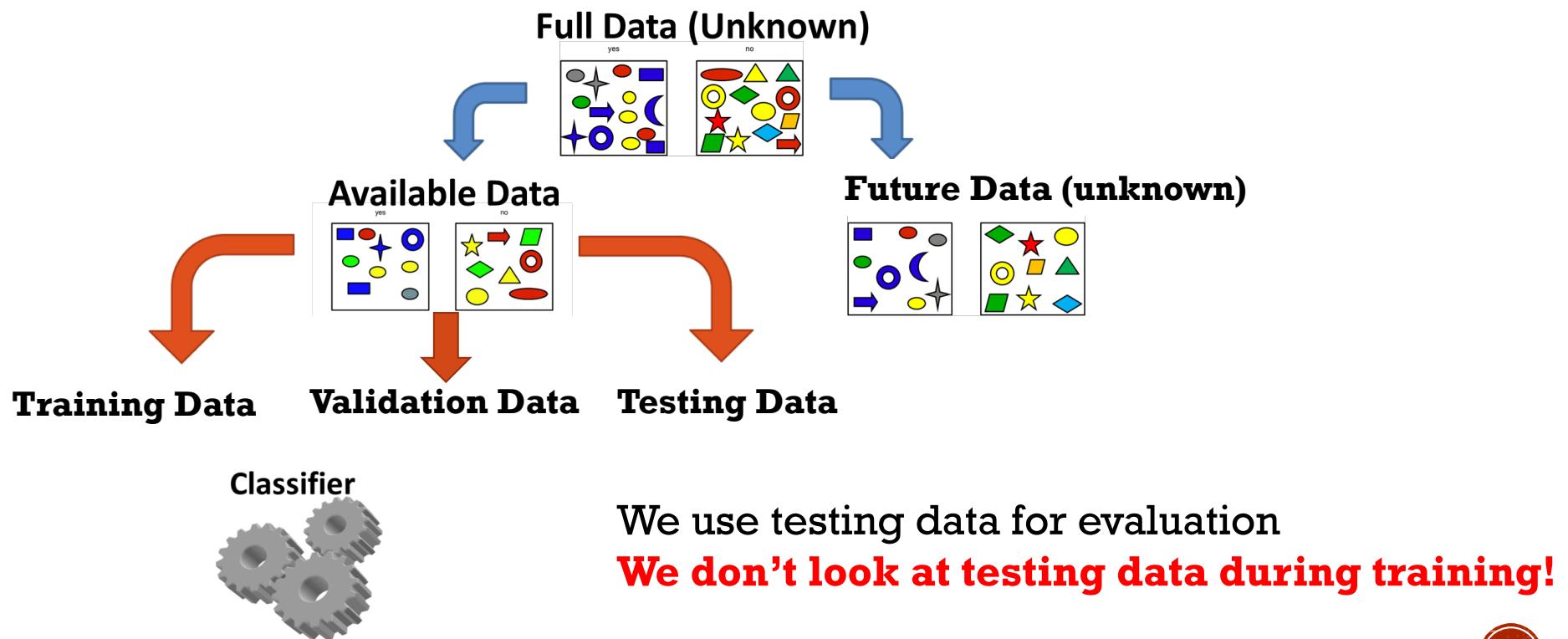


HOW TO LEARN FROM DATA

- In practice, we always split examples into 3 distinct sets
- **Training set**
 - Used to learn the **parameters** of the ML model
 - e.g., what are the nodes and branches of the decision tree
- **Validation set**
 - aka tuning set, aka development set, aka held-out data
 - Used to learn **hyperparameters**
 - Parameter that controls other parameters of the model
 - e.g., max depth of decision tree
- **Test set**
 - Used to evaluate how well we're doing on new unseen examples



SUPERVISED LEARNING



SUMMARY

- Decision trees
 - What they are, and how to learn one
- Fundamental machine learning concepts
 - What inductive bias is and what is its role in learning
 - What underfitting and overfitting means, and how it's related to the bias-variance tradeoff
 - Why you should never touch your test data

ANNOUNCEMENTS

- Quizzes have been graded
 - If you signed up for Gradescope, you would be able to see your result there
- If you are still trying to register and you have taken the quiz, fill out the form on Piazza
 - I will send individual emails to students the class can accommodate
- Any questions about the homework?



ACKNOWLEDGEMENTS

- These slides use materials by Brian Ziebart, Marine Carpuat, Tom Mitchell, and Jake Hofman

