

Limits and Evaluation I

Dr. Alex Williams

August 28, 2020



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

COSC 425: Introduction to Machine Learning
Fall 2020 (CRN: 44874)

Today's Agenda

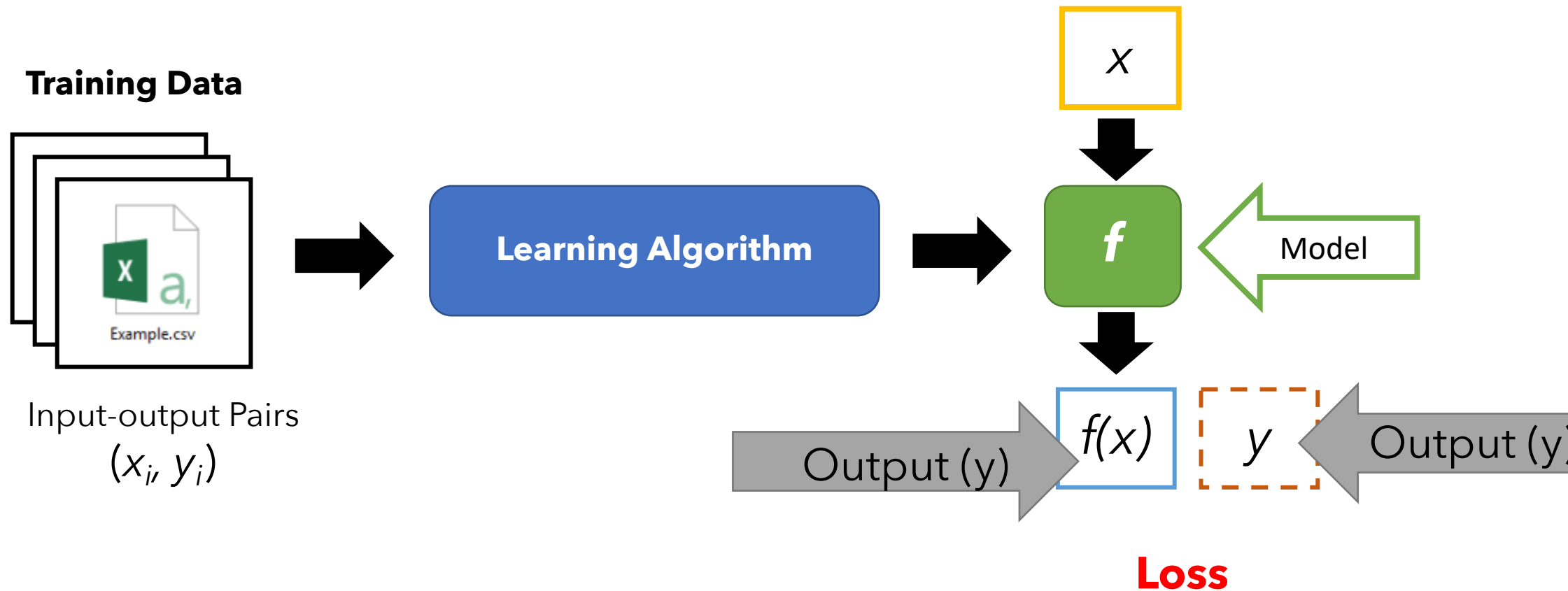


We will address:

1. What's the formal definition of performance?
2. How well can we possibly generalize?

1. What's the formal definition of performance?

Supervised Learning: Performance



Supervised Learning: Performance

To measure how well our learning algorithm is performing, we need to setup a loss function $l(y, \hat{y})$ based on our goal of learning, i.e. between predicted and observed.

squared loss $l(y, \hat{y}) = (y - \hat{y})^2$

absolute loss $l(y, \hat{y}) = |y - \hat{y}|$

zero/one loss $l(y, \hat{y}) = 1_{y \neq \hat{y}}$

Supervised Learning: Assumptions

Supervised learning “models” assume that a probability distribution \mathbf{D} exists over (\mathbf{x}, \mathbf{y}) called the data generating distribution.

Training Data



Note: The data generating distribution does not take any special form. It simply gives high probability to reasonable (\mathbf{x}, \mathbf{y}) pairs, and low probabilities to unreasonable (\mathbf{x}, \mathbf{y}) pairs.

We assume that the training data is a random sample of (\mathbf{x}, \mathbf{y}) pairs drawn from \mathbf{D} .

Input-output Pairs
 (x_i, y_i)

Bayes Optimal Classifier

Suppose that we have a magic function **computeD** that takes two inputs, \mathbf{x} and \mathbf{y} , and returned the probability of (\mathbf{x}, \mathbf{y}) pair under \mathbf{D} .

Then, the Bayes optimal classifier is the classifier that, for any test input x' , returns y' that maximizes **computeD** (x', y') :

$$f^{(BO)}(x') = \arg \max_{y' \in Y} D(x', y')$$

The classifier is optimal because, of all possible deterministic classifiers, it achieves the smallest zero/one error.

Intuition: Classify (i.e. by choosing a label) based on the highest probability.

- Minimizes the likelihood that an error will be made in post-training.

Supervised Learning: Performance Formulation

Based on the training data, we need to induce a function \mathbf{f} that maps new inputs \mathbf{x}' to a corresponding prediction \mathbf{y}' , such that the function does as well as possible on future examples that are also drawn from \mathcal{D} .

$$\epsilon = \mathbb{E}_{(x,y) \sim \mathcal{D}} [l(y, f(x))] = \sum_{(x,y)} \mathcal{D}(x, y) l(y, f(x))$$

All we have are examples from \mathcal{D} . We don't know what \mathcal{D} really is!

Supervised Learning: Performance Formulation

We can estimate future error by computing an error from training data:

$$\hat{\epsilon} = \frac{1}{N} \sum_{n=1}^N l(y_n, f(x_n))$$

The fundamental difficulty in machine learning is that we have access to our training error $\hat{\epsilon}$, but we truly only care about minimizing error on future instances.

In other words:

- We want our function **generalize** beyond the data we have at hand.

Defining Induction

Machine Learning Phrasing

Given a loss function \mathbf{L} and a sample \mathbf{D} from some unknown distribution, you must compute a function \mathbf{f} that has low expected error $\hat{\epsilon}$ over \mathbf{D} with respect to \mathbf{L} .

Alternative Phrasing:

Induction is the process of reaching a general conclusion from specific examples.

2. How well can we possibly generalize?

Generalization: What's Learnable?

Learning might fail because the inductive bias of the learning algorithm is **too "far away" from the concept of being learned**.

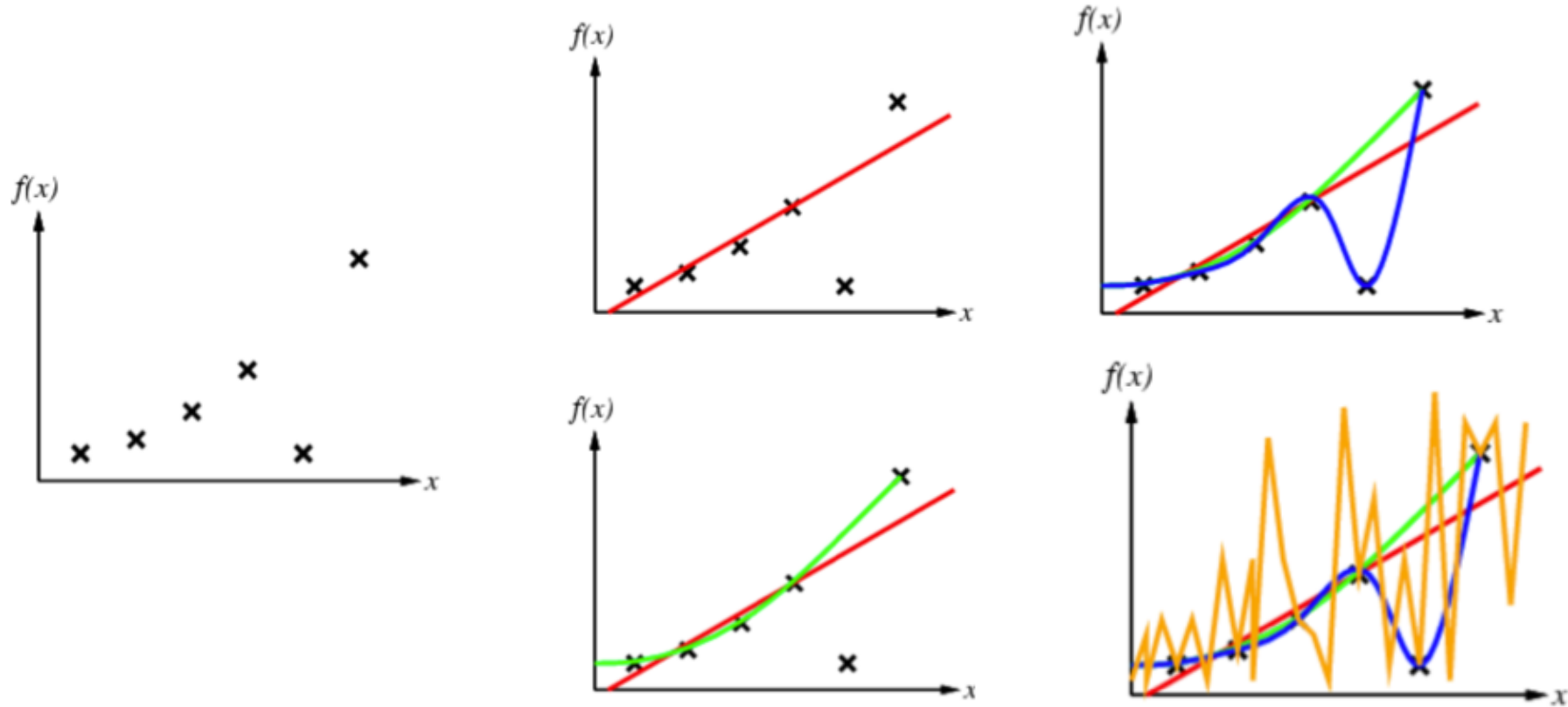
- Your data might reflect "bird" / "no bird".
- The algorithm may be learning something else entirely.

Learning can fail because of **bad data**.

- Input features can be noisy, e.g. sensor failures.
- Output labels can be noisy, e.g. annotator failures.
- Feature space may be insufficient, e.g. information missing.
- Output label can be ambiguous, e.g. fake news / not fake news.

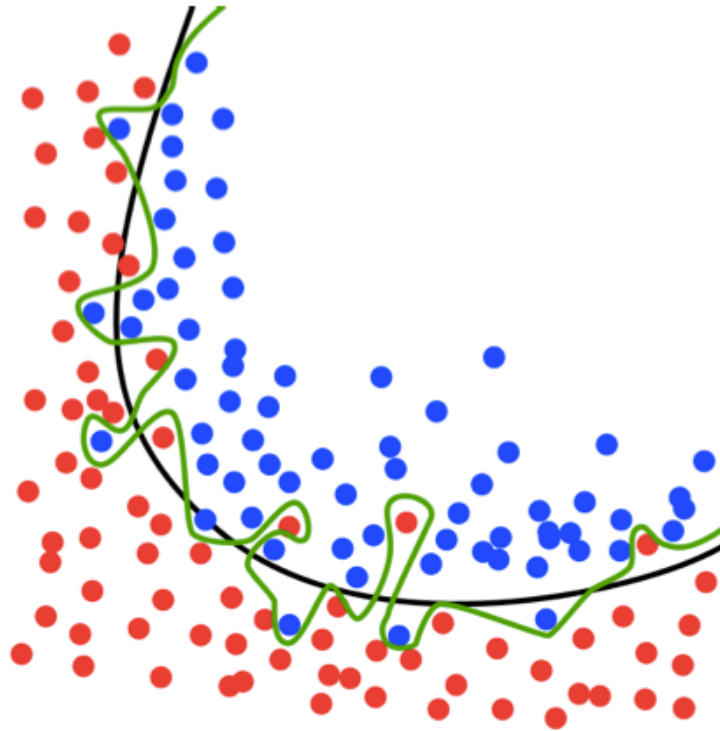
Generalization: Model Selection

Find a function h that fits f at instances \mathbf{x} .



Generalization: Model Selection

Degrees of freedom (i.e. Features) \rightarrow Training performance improves.



Generalization: Model Selection

Many learning algorithms are based on the notion of “modelling”.

- i.e. We have some formal “model” of our data.
- A decision tree is a model of our data.

Each method tells us what we can learn and what assumptions we make.

Models have **parameters** that are a lot like knobs / settings.

- The learning algorithm’s job is to find out what parameters work best in a given model, e.g. which questions to ask.

Models also have **hyper-parameters** that are controlled manually.

- e.g., max depth in a decision tree.

Generalization: Trading Off Bias & Variance

Suppose we're trying to learn a regression function $f(\mathbf{x})$ to approximate the true function \mathbf{y} .

We can prove that the expectation generalization error (based on squared loss) is a combination of noise, bias, and variance.

$$E[(y - f(x))^2] = \text{Noise}^2 + \text{Bias}^2 + \text{Variance}$$

https://en.wikipedia.org/wiki/Bias-variance_tradeoff

Generalization: Trading Off Bias & Variance

Noise: Irreducible Error

→ Our model will capture some level of irreducible error.

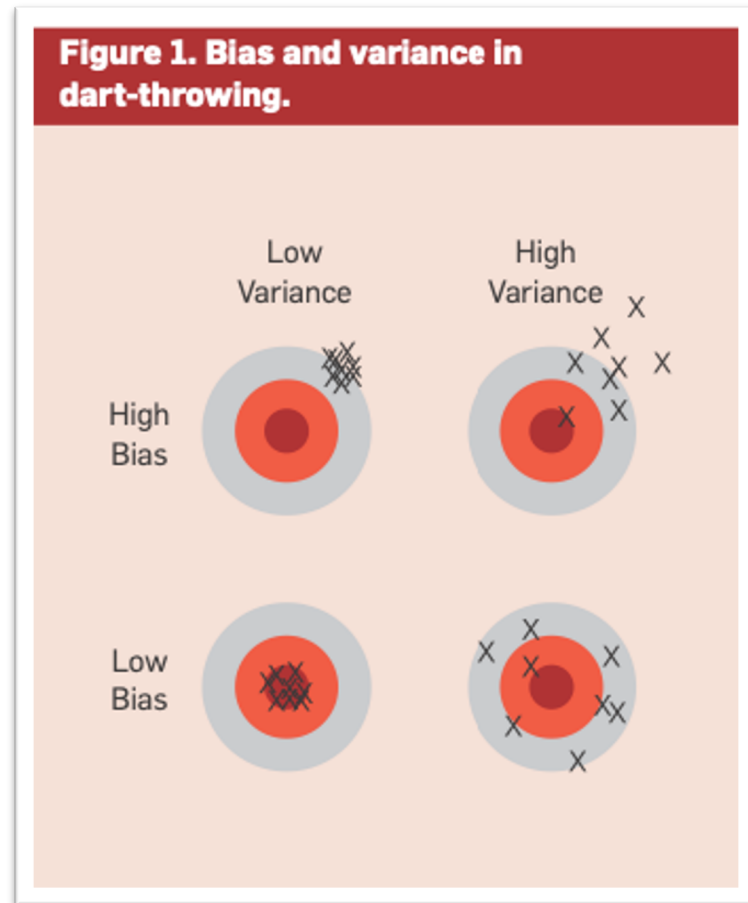
Bias: Tendency to consistency learn the wrong thing.

→ Our model can make erroneous assumptions about our data which may cause our model to be overly simplistic (Underfit!)

Variance: Tendency to consistently learn random things.

→ Our model can be erroneously sensitive to small changes in our data which may cause our model to map to its patterns too tightly. (Overfit!)

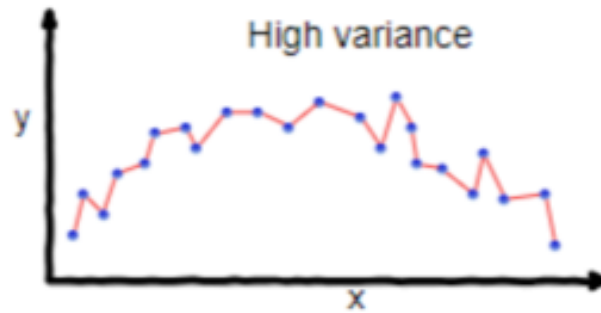
Generalization: Trading Off Bias & Variance



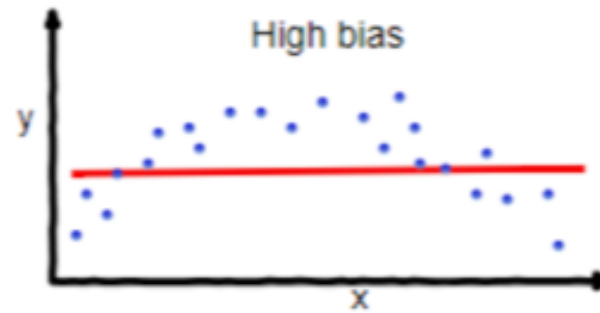
A Few Useful Things to Know about Machine Learning.

Domingos, 2012.

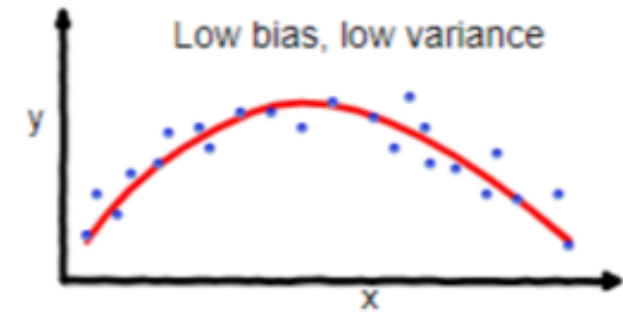
Generalization: Trading Off Bias & Variance



overfitting



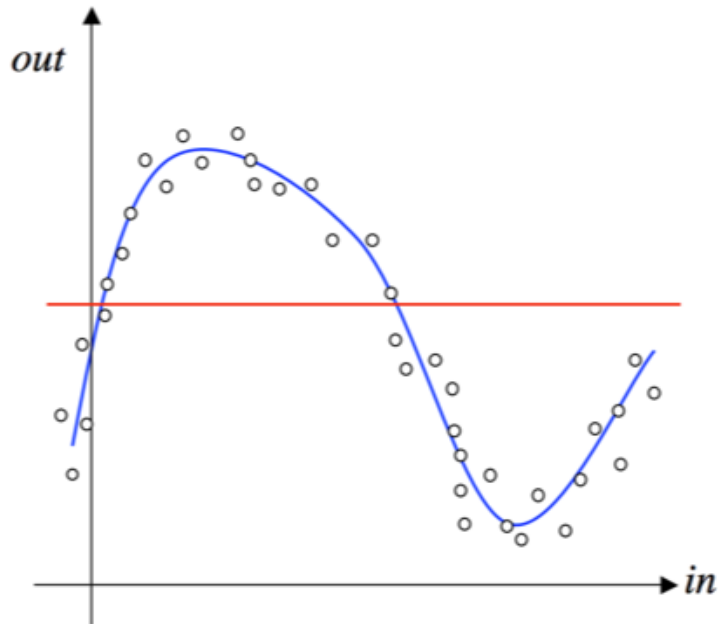
underfitting



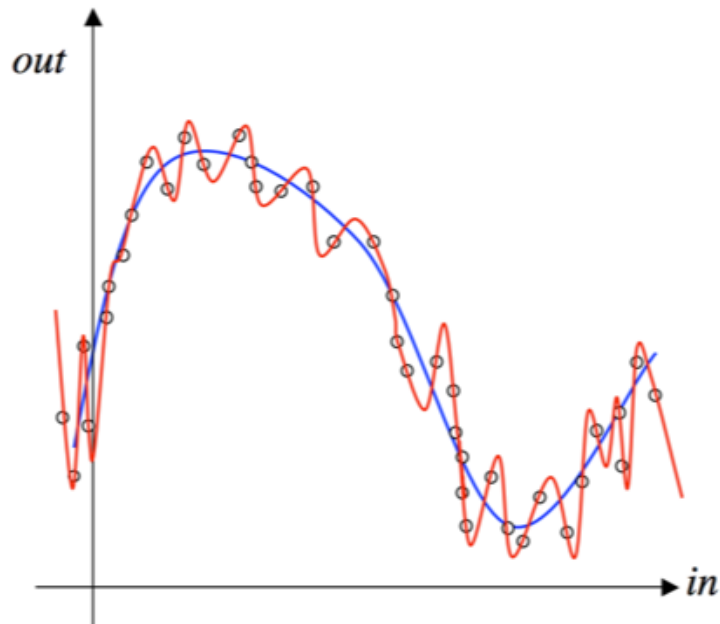
Good balance

<https://towardsdatascience.com/understanding-the-bias-variance-tradeoff-165e6942b229>

Generalization: Trading Off Bias & Variance



Ignore the data \Rightarrow
Big approximation error (high bias)
No variation between data sets (no variance)



Fit every data point \Rightarrow
No approximation error (zero bias)
Variation between data sets (high variance)

<https://www.cs.bham.ac.uk/~jxb/INC/l9.pdf>

Decision Trees: Bias and Variance

Consider an **empty** vs **full** decision tree.

Empty Tree → Same prediction regardless of its input.
(It doesn't ask questions!)

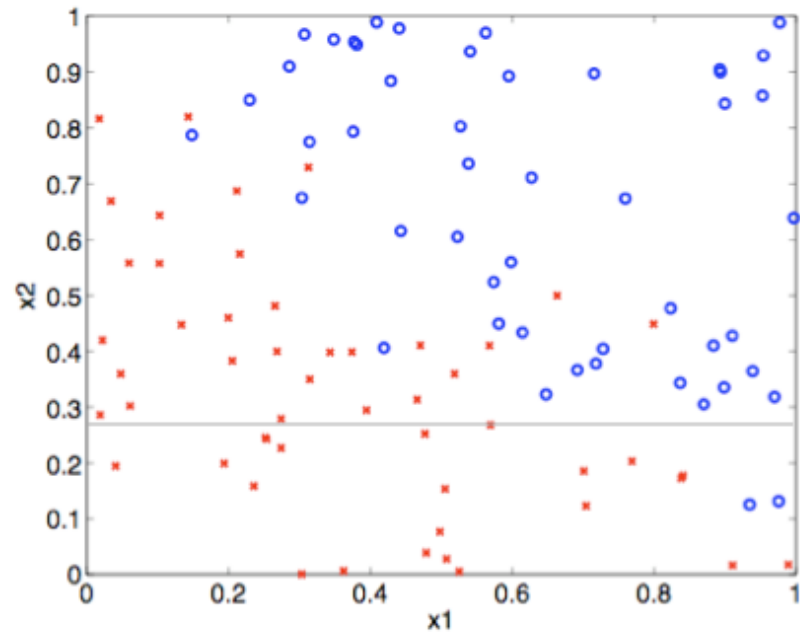
Full Tree → Will always make correct predictions.
(Assuming every leaf has 0 or 1 examples assigned)

How would these trees perform on future data?

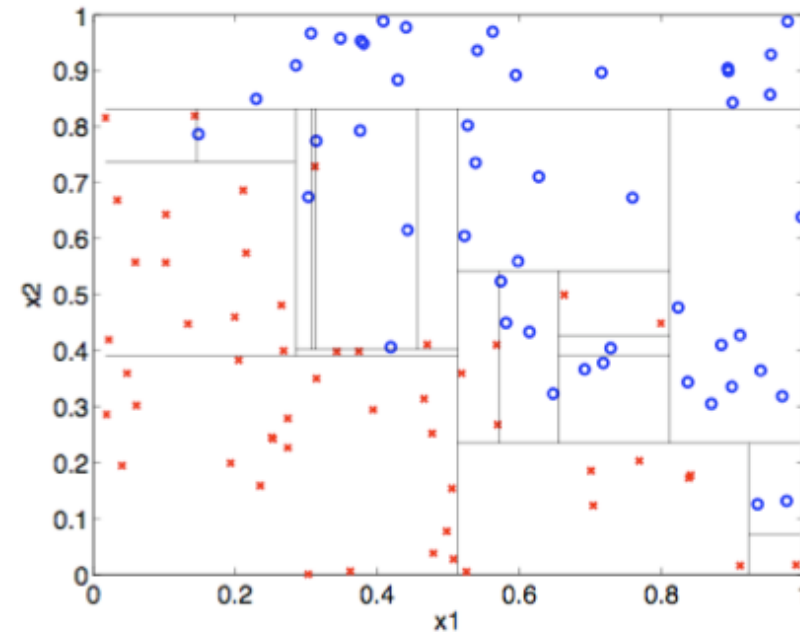
Empty Tree → We expect it to make the same error on all future examples, (i.e. high bias, low variance)

Full Tree → For examples identical to training, it will do the right thing. Otherwise, significantly large error (i.e. low bias, high variance)

Decision Trees: Bias and Variance



High bias, Low variance



Low bias, high variance

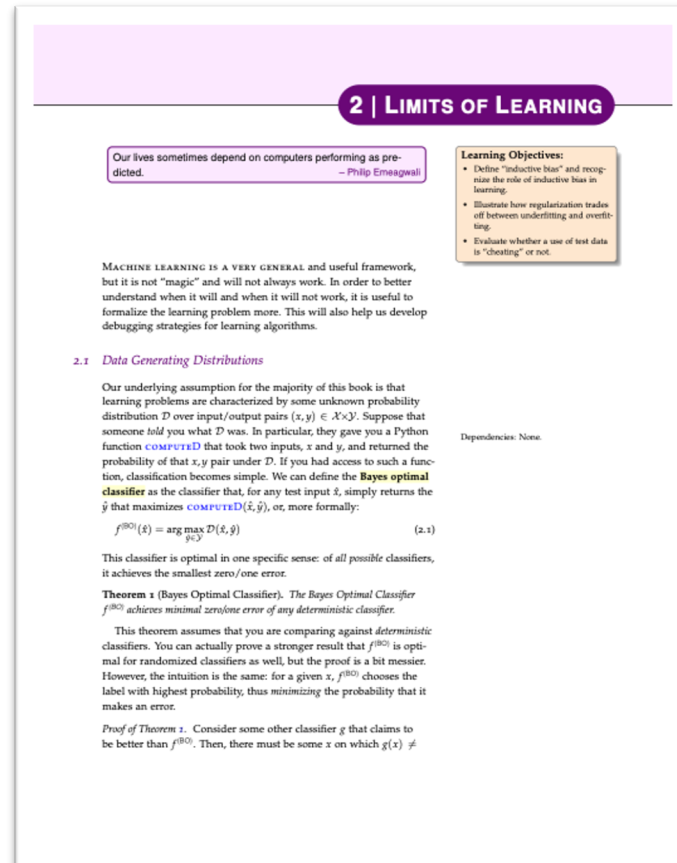
Today's Agenda



We have addressed:

1. What's the formal definition of performance?
2. How well can we possibly generalize?

Reading



Daumé, Sec 2.1 – 2.4.



Optional: Domingos, 2012.

Next Time

We will address:



1. What should we do to improve generalization?
2. What are common performance metrics in machine learning?