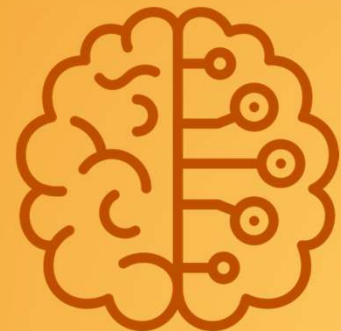

CS 4375 – Introduction to Machine Learning

Bias/Variance Tradeoff & Bagging

Erick Parolin



[Based on the slides from Dr. Nicholas Ruozzi]

Previously...

- **PAC-learnability:** C is said to be *PAC-learnable* by L using H if and only if learner L will with probability at least $1 - \delta$ output a hypothesis $h \in H$, such that:

- $\text{error}_{\mathcal{D}}(h) \leq \epsilon$
- **polynomial time in $1/\epsilon$, $1/\delta$, n and $\text{size}(c)$**

Says: we are willing to tolerate a δ probability of having $> \epsilon$ true error.

- So, that's what we want: $P(\text{error}_{\mathcal{D}}(h) > \epsilon) \leq \delta$
- PAC-learnability is mostly concerned about the number of training examples required
 - Parameters ϵ and δ will determine the number of training examples
 - The number of examples with respect to ϵ and δ must be polynomial
- **Theorem (Haussler, 1988):** For finite H and D ($m \geq 1$ independent random examples), the probability that the version space $VS_{H,D}$ is **not** ϵ -exhausted is less than or equal to $|H|e^{-\epsilon m}$

$$P(\{\exists h \in VS_{H,D} \mid \text{error}_{\mathcal{D}}(h) > \epsilon\}) \leq |H|e^{-\epsilon m} \quad \text{or simply} \quad P(\text{error}_{\mathcal{D}}(h) > \epsilon) \leq |H|e^{-\epsilon m}$$

- **From PAC:** $P(\text{error}_{\mathcal{D}}(h) > \epsilon) \leq |H|e^{-\epsilon m} \leq \delta$

Previously...

- From PAC and Haussler:

$$P(\{\exists h \in \mathcal{H} \mid \text{error}_{\mathcal{D}}(h) > \epsilon\}) \leq |\mathcal{H}|e^{-\epsilon m} \leq \delta$$

- Which gives us the **sample complexity**:

$$m \geq (\ln 1/\delta + \ln |\mathcal{H}|)/\epsilon$$

- This tells us how many training examples suffice to ensure (with probability $(1 - \delta)$) that every hypothesis in \mathcal{H} having zero training error will have a true error of at most ϵ .

- Agnostic Learning** assumes no prior commitment about whether or not $c \in \mathcal{H}$:

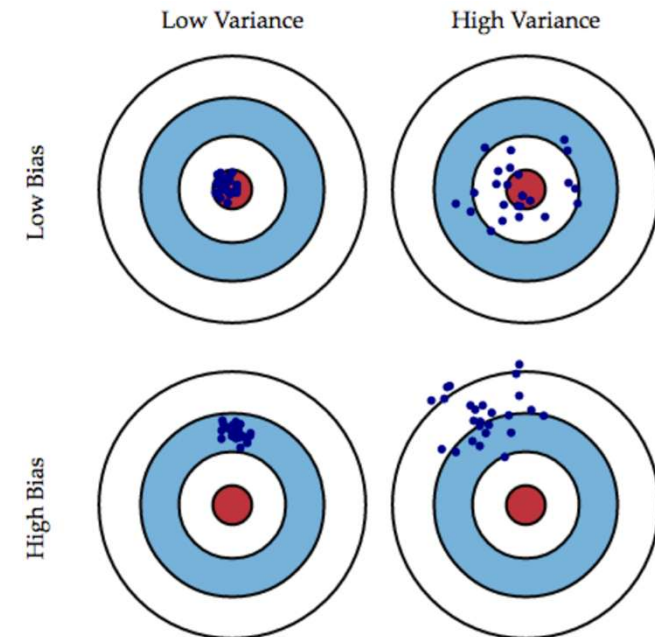
$$P(\text{error}_{\mathcal{D}}(h_{\text{best}}) > \epsilon + \text{error}_{\mathcal{D}}(h_{\text{best}})) \leq |\mathcal{H}|e^{-2m\epsilon^2}$$

- Sample Complexity: $m \geq \frac{1}{2\epsilon^2} (\ln \frac{1}{\delta} + \ln |\mathcal{H}|)$

- For all h , with probability at least $1 - \delta$:
$$\text{error}_{\mathcal{D}}(h_{\text{best}}) \leq \underbrace{\text{error}_{\mathcal{D}}(h_{\text{best}})}_{\text{Bias}} + \underbrace{\sqrt{\frac{\ln \frac{1}{\delta} + \ln |\mathcal{H}|}{2m}}}_{\text{Variance}}$$

Intuition

- **Bias**
 - Measures the accuracy or quality of the learner algorithm
 - High bias means a poor match
- **Variance**
 - Measures the precision or specificity of the match
 - High variance means a weak match
- We would like to minimize each of these
- Unfortunately, we can't do this independently, there is a tradeoff



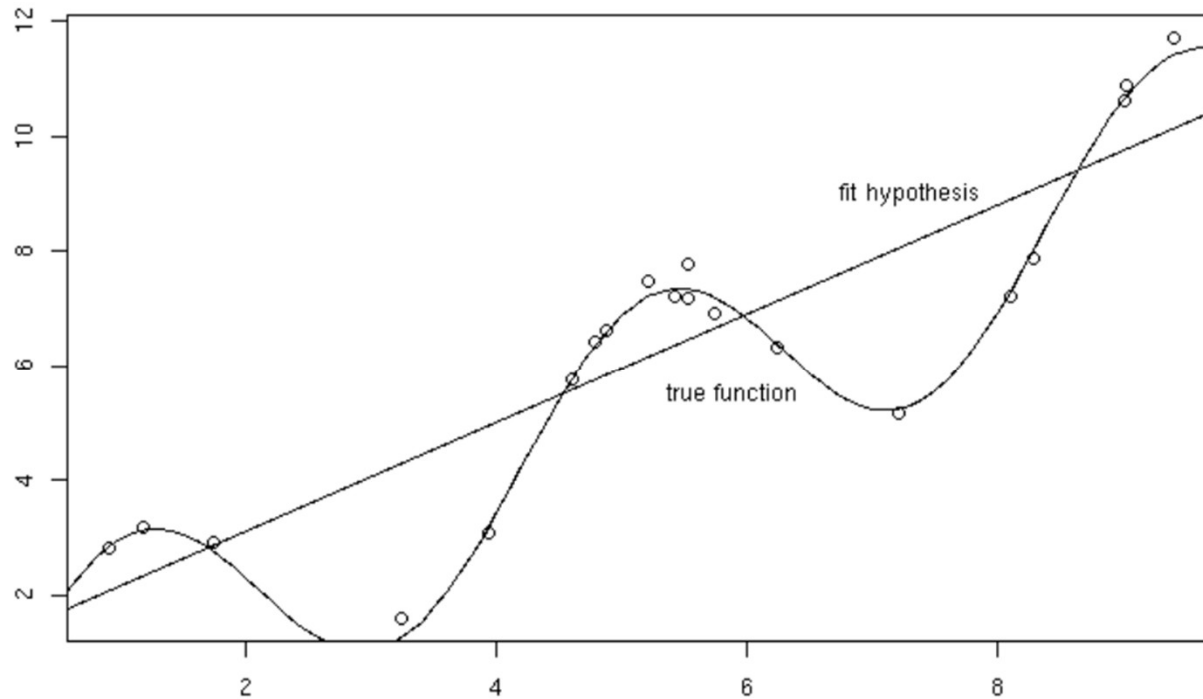
Bias-Variance Analysis in Regression

- **True function is** $y = f(x) + \epsilon$
 - Where noise, ϵ , is normally distributed with zero mean and standard deviation σ
- Given a set of training examples, $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$, we fit a hypothesis $\mathbf{g}(x) = \mathbf{w}^T \mathbf{x} + \mathbf{b}$ to the data to minimize the squared error

$$\sum_i [y^{(i)} - \mathbf{g}(x^{(i)})]^2$$

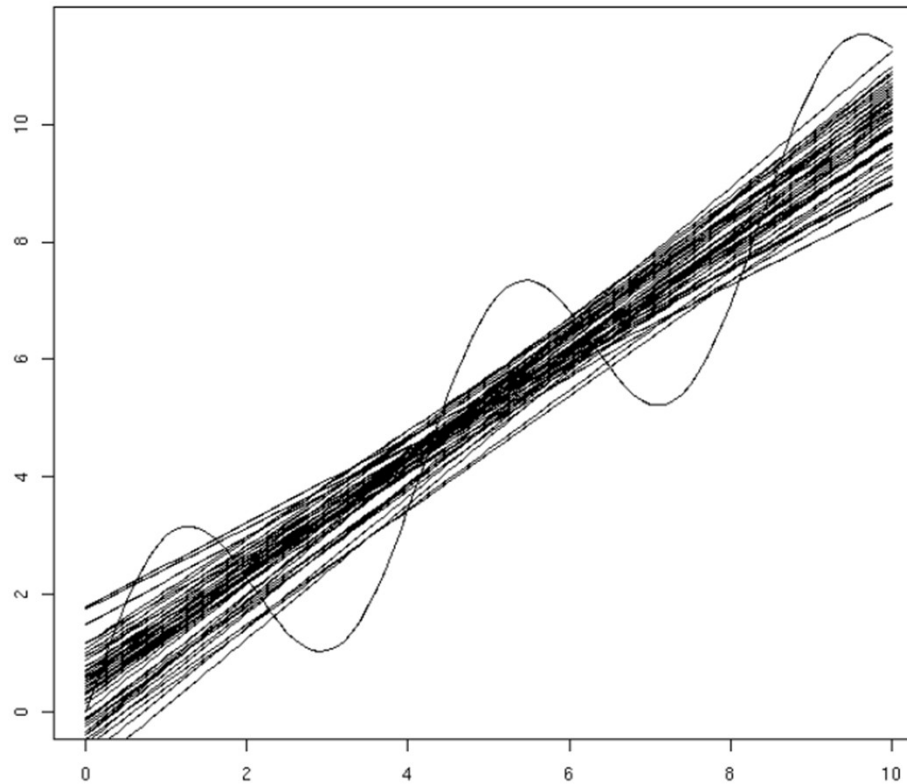
2-D Example

Sample 20 points from $f(x) = x + 2 \sin(1.5x) + N(0,0.2)$



2-D Example

50 fits (20 examples each)



Bias-Variance Analysis

- Given a new data point x' with observed value $y' = f(x') + \epsilon$, want to understand the **expected prediction error**.
- Suppose that training samples are drawn independently from a distribution $p(S)$, want to compute the **expected error** of the estimator

$$E[(y' - g_S(x'))^2]$$

Probability Reminder

- Variance of a random variable, Z

$$\begin{aligned} \text{Var}(Z) &= E[(Z - E[Z])^2] \\ &= E[Z^2 - 2ZE[Z] + E[Z]^2] \\ &= E[Z^2] - E[Z]^2 \end{aligned}$$

Bias-Variance-Noise Decomposition

$$\begin{aligned} E \left[(y' - g_S(x'))^2 \right] &= E[g_S(x')^2 - 2g_S(x')y' + y'^2] \\ &= E[g_S(x')^2] - \boxed{2E[g_S(x')]E[y']} + E[y'^2] \end{aligned}$$

Notes:

- $E[A + B] = E[A] + E[B]$
- $E[AB] = E[A] \cdot E[B]$ only if the two random variables are statistically independent
- The samples S and the noise ϵ are independent, so $g_S(x')$ **and** y' **are independent**

Bias-Variance-Noise Decomposition

$$\begin{aligned} E \left[(y' - g_S(x'))^2 \right] &= E[g_S(x')^2 - 2g_S(x')y' + y'^2] \\ &= E[g_S(x')^2] - 2E[g_S(x')]E[y'] + E[y'^2] \\ &= \text{Var}(g_S(x')) + E[g_S(x')]^2 - 2E[g_S(x')]f(x') \\ &\quad + \text{Var}(y') + f(x')^2 \end{aligned}$$

Notes:

- Remember: $y' = f(x') + \epsilon$ where ϵ is normally distributed with zero mean $\rightarrow E[y'] = f(x')$
- From definition of Variance: $\text{Var}[Z] = E[Z^2] - E[Z]^2$
- $\text{Var}(g_S(x')) = E[g_S(x')^2] - E[g_S(x')]^2 \Rightarrow E[g_S(x')^2] = \text{Var}(g_S(x')) + E[g_S(x')]^2$
- Similarly: $\text{Var}(y') = E[y'^2] - E[y']^2 \Rightarrow E[y'^2] = \text{Var}(y') + f(x')^2$

Bias-Variance-Noise Decomposition

$$\begin{aligned} E \left[(y' - g_S(x'))^2 \right] &= E[g_S(x')^2 - 2g_S(x')y' + y'^2] \\ &= E[g_S(x')^2] - 2E[g_S(x')]E[y'] + E[y'^2] \\ &= \text{Var}(g_S(x')) + E[g_S(x')]^2 - 2E[g_S(x')]f(x') \\ &\quad + \text{Var}(y') + f(x')^2 \\ &= \text{Var}(g_S(x')) + (E[g_S(x')] - f(x'))^2 + \text{Var}(\epsilon) \end{aligned}$$

Notes:

- $E[g_S(x')]^2 - 2E[g_S(x')]f(x') + f(x')^2 = (E[g_S(x')] - f(x'))^2$
- $\text{Var}(y') = \text{Var}(f(x') + \epsilon) = \text{Var}(\epsilon) = \sigma^2$

Bias-Variance-Noise Decomposition

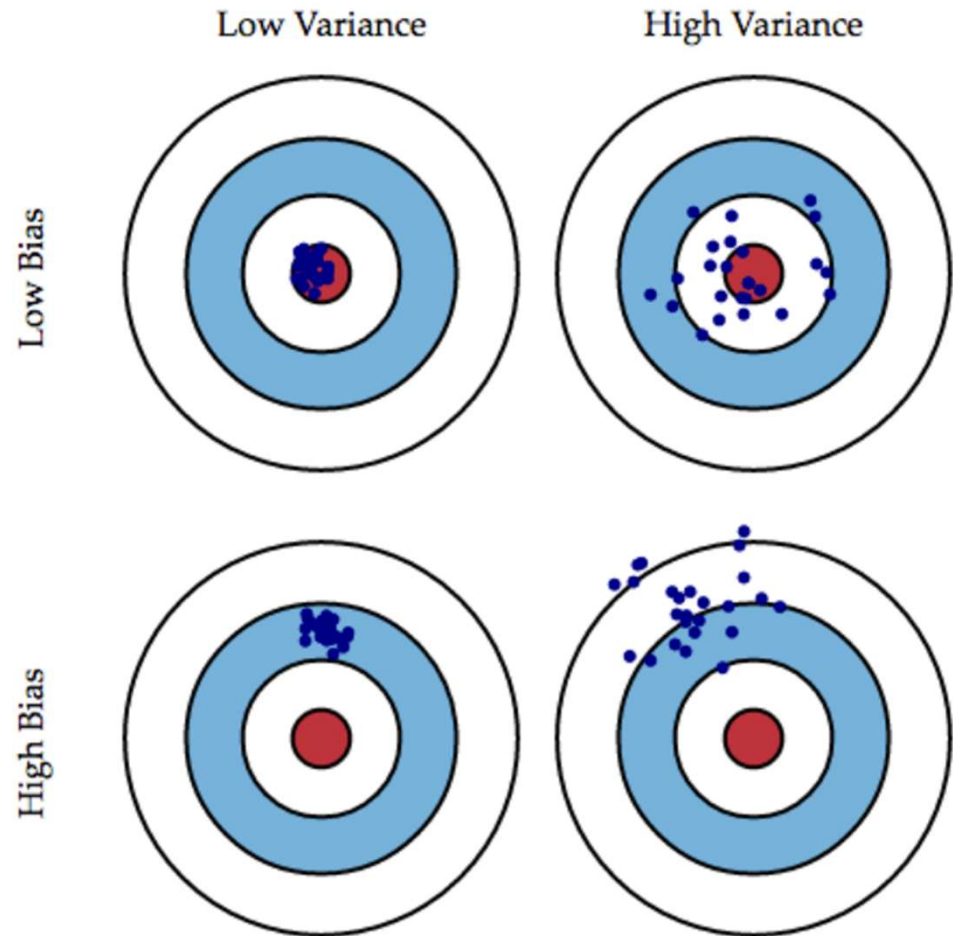
$$\begin{aligned} E \left[(y' - g_S(x'))^2 \right] &= E[g_S(x')^2 - 2g_S(x')y' + y'^2] \\ &= E[g_S(x')^2] - 2E[g_S(x')]E[y'] + E[y'^2] \\ &= \text{Var}(g_S(x')) + E[g_S(x')]^2 - 2E[g_S(x')]f(x') \\ &\quad + \text{Var}(y') + f(x')^2 \\ &= \text{Var}(g_S(x')) + (E[g_S(x')] - f(x'))^2 + \text{Var}(\epsilon) \\ &= \underbrace{\text{Var}(g_S(x'))}_{\text{Variance}} + \underbrace{(E[g_S(x')] - f(x'))^2}_{\text{Bias}} + \underbrace{\sigma^2}_{\text{Noise}} \end{aligned}$$

Bias, Variance, and Noise

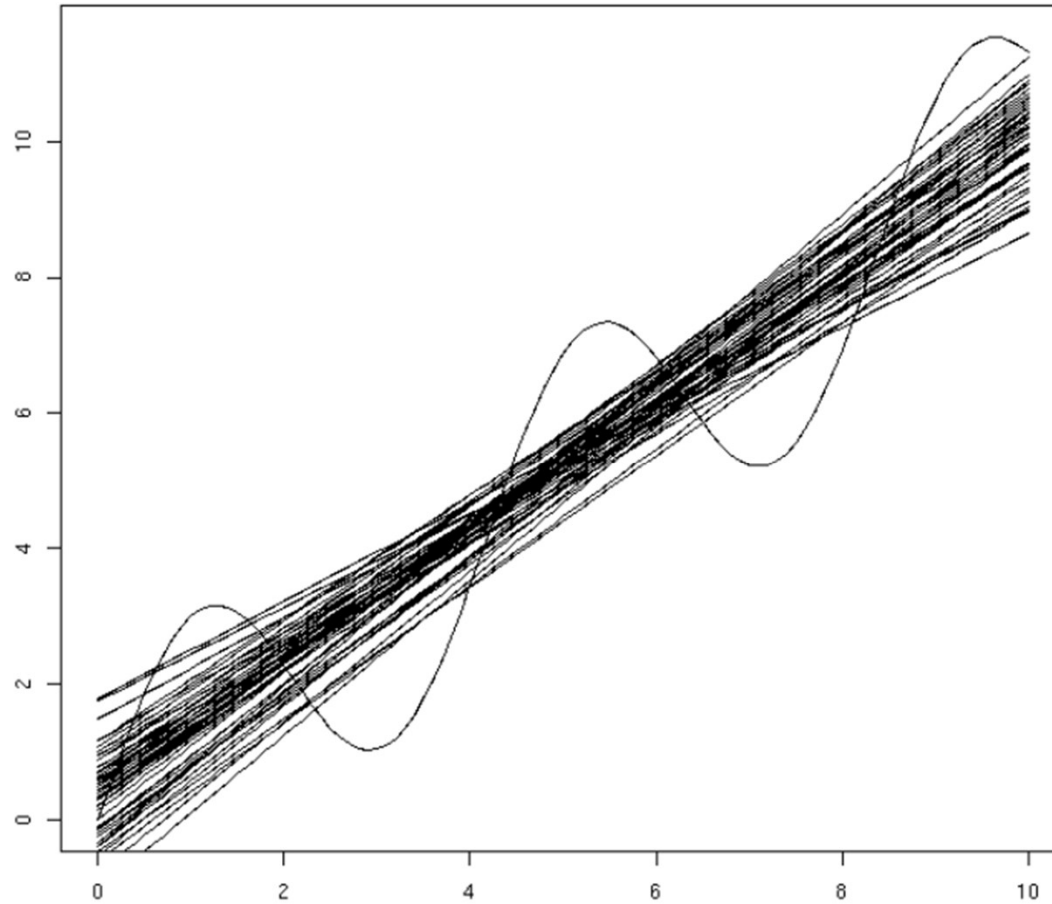
- Variance: $E[(g_{S(x')}) - E[g_{S(x')})]^2]$
 - Describes how much $g_S(x')$ varies from one training set S to another
- Bias: $E[g_S(x')] - f(x')$
 - Describes the average error of $g_S(x')$
- Noise: $E[(y' - f(x'))^2] = E[\epsilon^2] = \sigma^2$
 - Describes how much y' varies from $f(x')$

Bias, Variance, and Noise

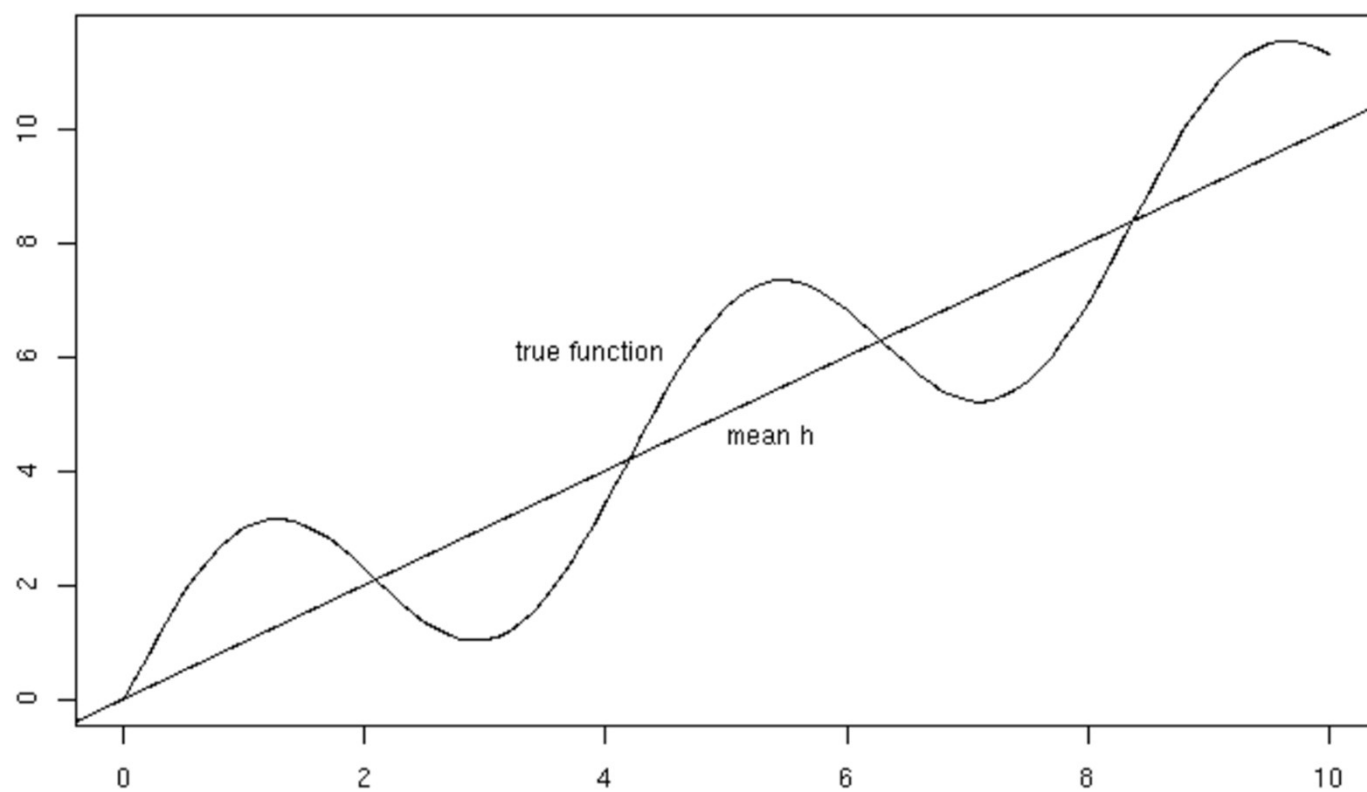
$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Noise}$$



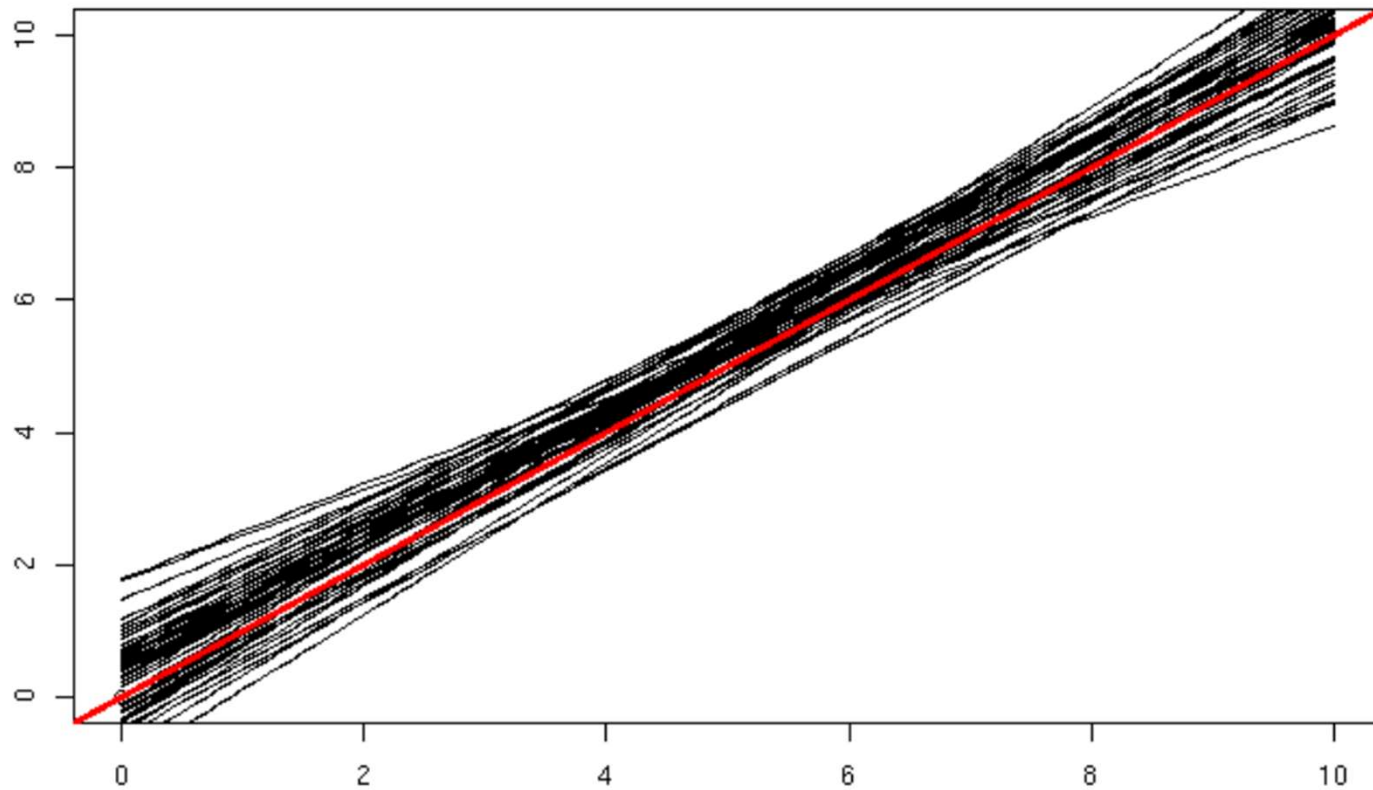
50 fits (20 examples each)



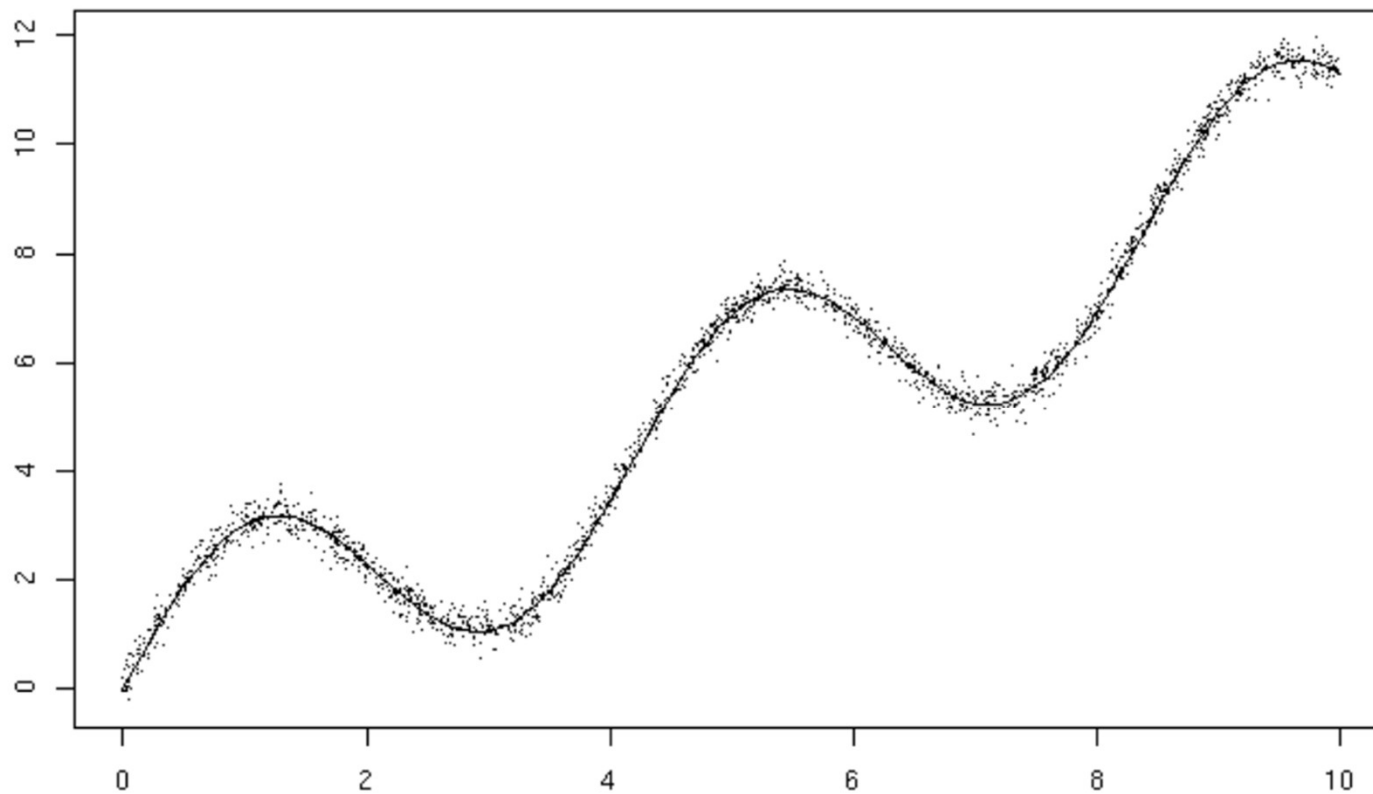
Bias



Variance



Noise



Bias

- **Low bias**
 - Linear regression applied to linear data
 - 2nd degree polynomial applied to quadratic data
- **High bias**
 - Constant function
 - Linear regression applied to highly non-linear data

Variance

- **Low variance**
 - Constant function
 - Simple model independent of training data
- **High variance**
 - Complex model, e.g., neural nets
 - High degree polynomial

Bias/Variance Tradeoff

- **To get low bias**
 - You need to **approximate** as much as possible to the **real function**.
 - Usually requires a **complex model** that has larger number of parameters.
 - Higher complexity
 - Parameters are learned based on **training data**
 - So, for **each sample**, you are likely to get a **different model**.
 - Different models for different samples \Rightarrow **high variance**
 - **So, you have a low bias with high variance model.**

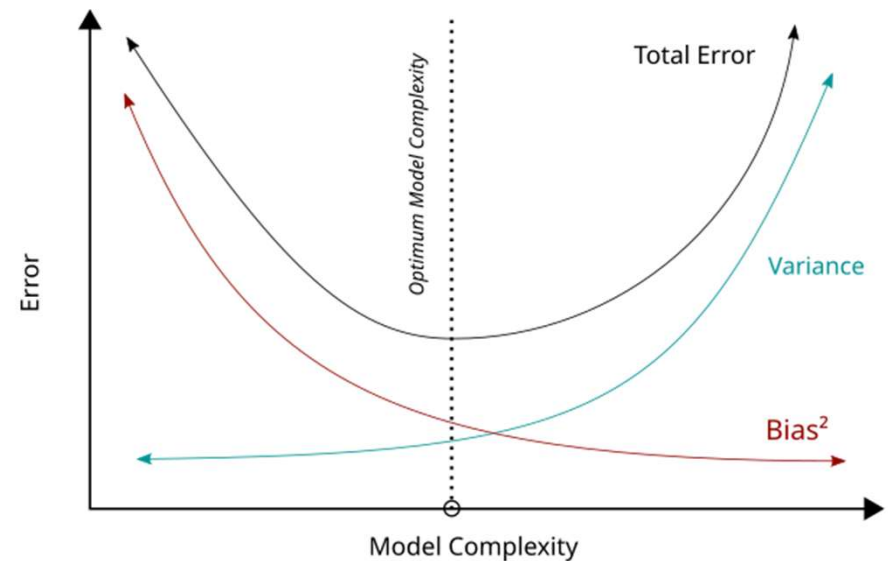
Bias/Variance Tradeoff

- **To get low variance**

- You need a model more “**stable**”, whose output doesn’t change much on different samples.
- A very simple model (e.g., a constant) **doesn’t change its outputs for different samples.**
 - Very low complexity
- But the output will be very **different from the real function** that you are trying to approximate ⇒ **high bias**
- **So, you have a low variance with high bias model.**

Bias/Variance Tradeoff

- $(\text{bias}^2 + \text{variance})$ is what counts for prediction
- Often:
 - Low bias \Rightarrow high variance
 - Low variance \Rightarrow high bias
- Tradeoff
 - bias² vs. variance
- Is this a firm rule?



Reduce Variance Without Increasing Bias

- **Averaging** reduces variance:

$$\text{Var}\left(\frac{1}{N}\sum_i z_i\right) = \frac{1}{N}\text{Var}(Z)$$

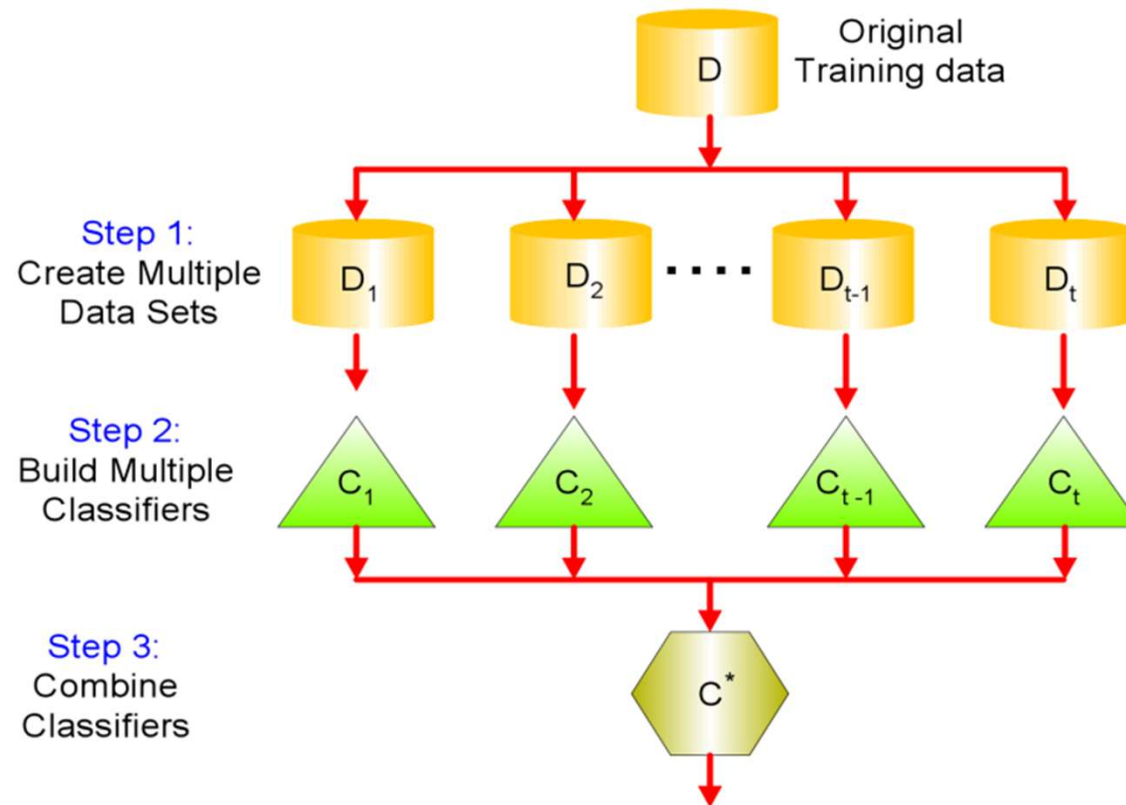
where $Z = z_1, \dots, z_N$ is i.i.d. random variables

- Problem:
 - Only one training set
 - Where do multiple models come from?

Bagging: Bootstrap Aggregation

- Take repeated bootstrap samples from training set D (Breiman, 1994)
- **Bootstrap sampling:** Given set D containing N training examples, create D' by drawing N examples at random with replacement from D
- **Bagging:**
 - Create k bootstrap samples D_1, \dots, D_k
 - Train distinct classifier on each D_i
 - Classify new instance by majority vote / average

Bagging: Bootstrap Aggregation



Bootstrap Sampling

- In each bootstrap sample, each data point has probability $\left(1 - \frac{1}{N}\right)^N$ of **not** being selected
- So, probability of being selected is $1 - \left(1 - \frac{1}{N}\right)^N$
- Expected number of distinct data points in each sample is then

$$N \cdot \left(1 - \left(1 - \frac{1}{N}\right)^N\right) \approx N \cdot (1 - \exp(-1)) = 0.632N$$

Remember: e^x can be written as $\exp(x) = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n$

Analysis

- **Bagging:** Output an aggregated model with less variance.
- **Ideally:**

$$\text{Var}\left(\text{Bagging}(L(x, D))\right) = \frac{\text{Var}(L(x, D))}{N}$$

- **In practice:**
 - Models are correlated, so reduction is smaller than $1/N$

Bagging Empirical Results

Experiments design:

- Many datasets were split into Train and Test.
- Two learners/classifiers were tested:
 - Decision Tree
 - Bagging (using Decision Tree)
- Observed error (misclassification rate):
 - \bar{e}_S is the error for Decision Tree learner.
 - \bar{e}_B is the error for Bagging.
- Experiments were repeated 100 times, and the reported \bar{e}_S and \bar{e}_B are the average over the 100 iterations.

| Data Set | \bar{e}_S | \bar{e}_B | Decrease |
|---------------|-------------|-------------|----------|
| waveform | 29.1 | 19.3 | 34% |
| heart | 4.9 | 2.8 | 43% |
| breast cancer | 5.9 | 3.7 | 37% |
| ionosphere | 11.2 | 7.9 | 29% |
| diabetes | 25.3 | 23.9 | 6% |
| glass | 30.4 | 23.6 | 22% |
| soybean | 8.6 | 6.8 | 21% |

When Will Bagging Improve Accuracy?

- Depends on the stability of the base-level classifiers
- A learner is **unstable** if a small change to the training set causes a large change in the output hypothesis
 - If small changes in D cause large changes in the output, then there will likely be an improvement in performance with bagging
- Bagging can help unstable procedures, but could hurt the performance of stable procedures
 - Decision trees are unstable
 - K-Nearest Neighbor and Naïve Bayes are stable classifiers

Random Forests

- Remember from couple of slides ago:
“In practice, models are correlated, so reduction is smaller than $1/N$ ”
- So, why not introduce some more randomness?
- Build large collection of de-correlated trees and average them.
- How to make the trees de-correlated?

Algorithm 15.1 *Random Forest for Regression or Classification.*

1. For $b = 1$ to B :
 - (a) Draw a bootstrap sample \mathbf{Z}^* of size N from the training data.
 - (b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.
 - i. Select m variables at random from the p variables.
 - ii. Pick the best variable/split-point among the m .
 - iii. Split the node into two daughter nodes.
2. Output the ensemble of trees $\{T_b\}_1^B$.

To make a prediction at a new point x :

Regression: $\hat{f}_{rf}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$.

Classification: Let $\hat{C}_b(x)$ be the class prediction of the b th random-forest tree. Then $\hat{C}_{rf}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$.

Random Forests

- Ensemble method specifically designed for decision tree classifiers
- Introduce two sources of randomness: “bagging” and “random input vectors”
- Bagging method: each tree is grown using a bootstrap sample of training data
- **Random vector method:** best split at each node is chosen from a random sample of m attributes instead of all attributes

Random Forest = Bagging + Random Vector Method

Readings

- **Pattern Recognition and Machine Learning by Christopher M. Bishop – Chapter 14**
- **Bagging Predictors by Leo Breiman (1996)**