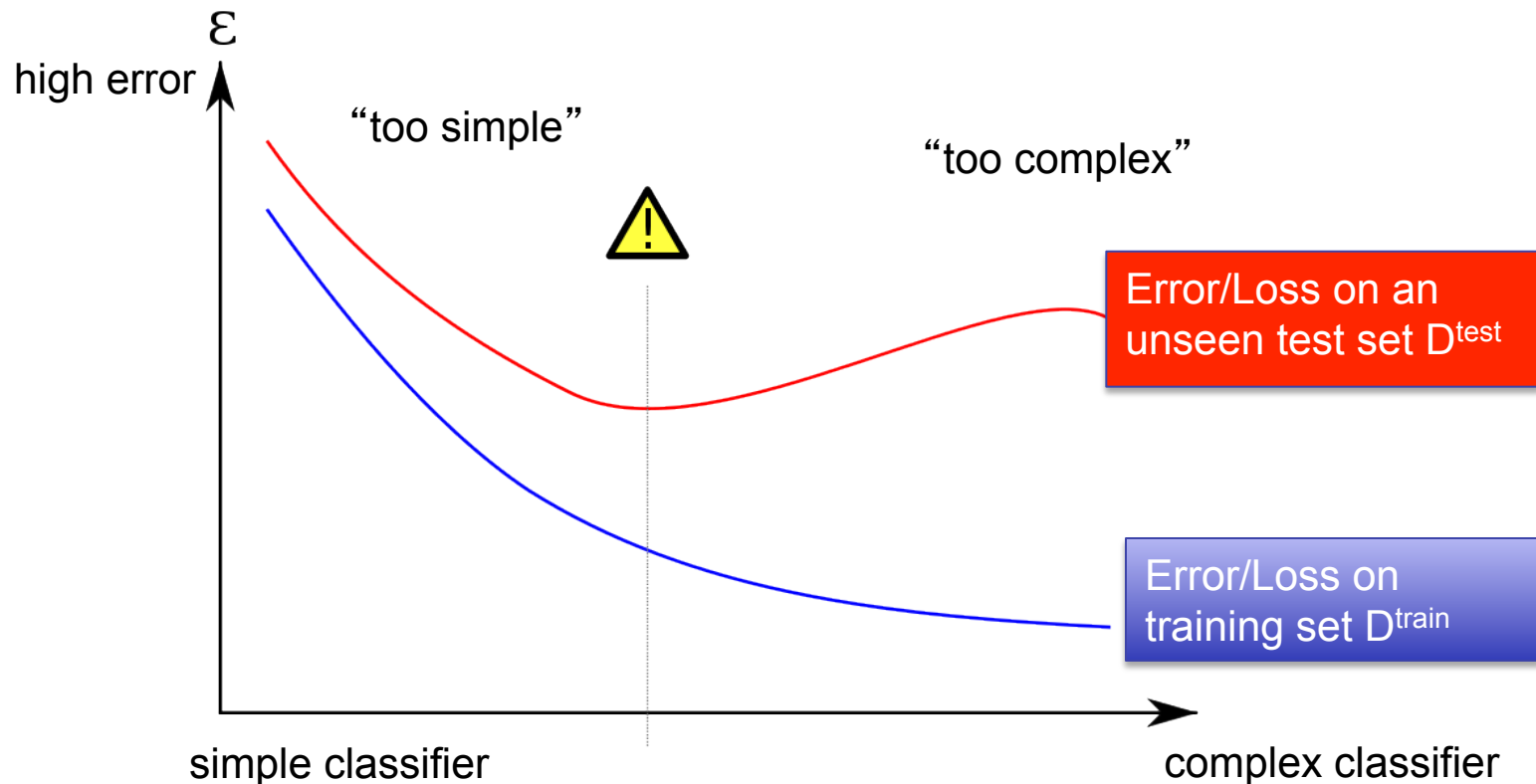


# Bias-Variance in Machine Learning

# Bias-Variance: Outline

- Underfitting/overfitting:
  - Why are complex hypotheses bad?
- Simple example of bias/variance
- Error as bias+variance for regression
  - brief comments on how it extends to classification
- Measuring bias, variance and error
- Bagging - a way to reduce variance
- Bias-variance for classification

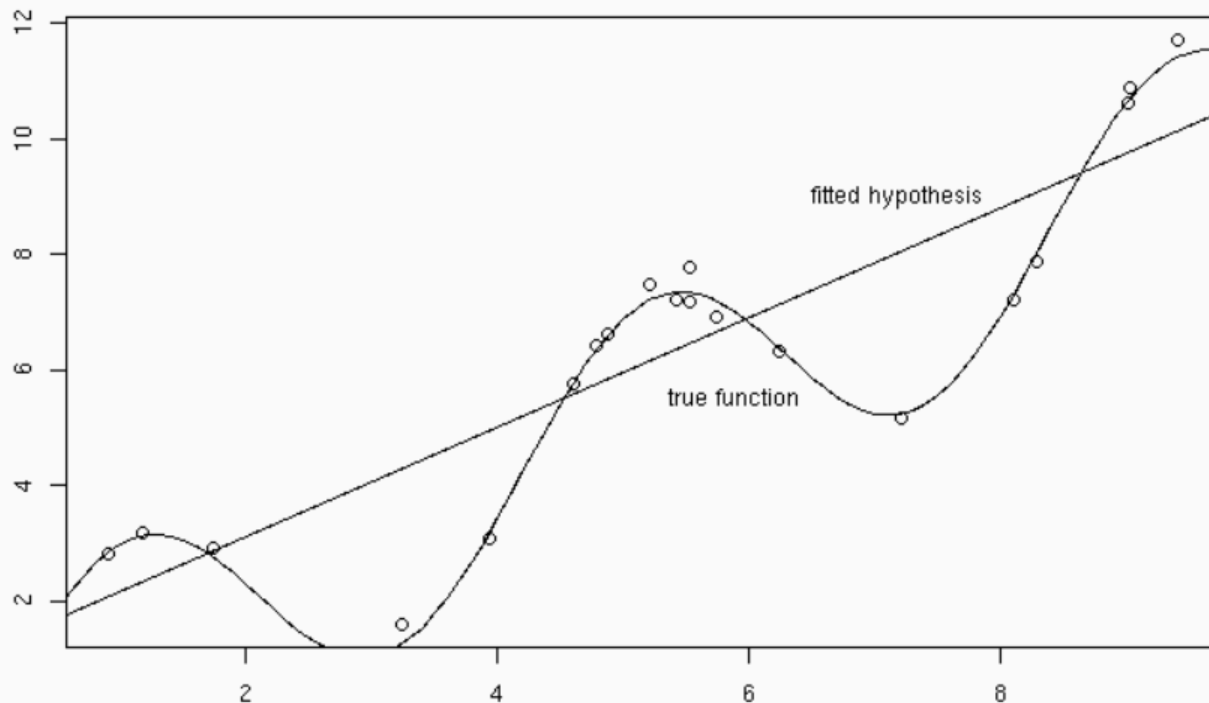
# Bias/Variance is a Way to Understand Overfitting and Underfitting



# Bias-Variance: An Example

# Example

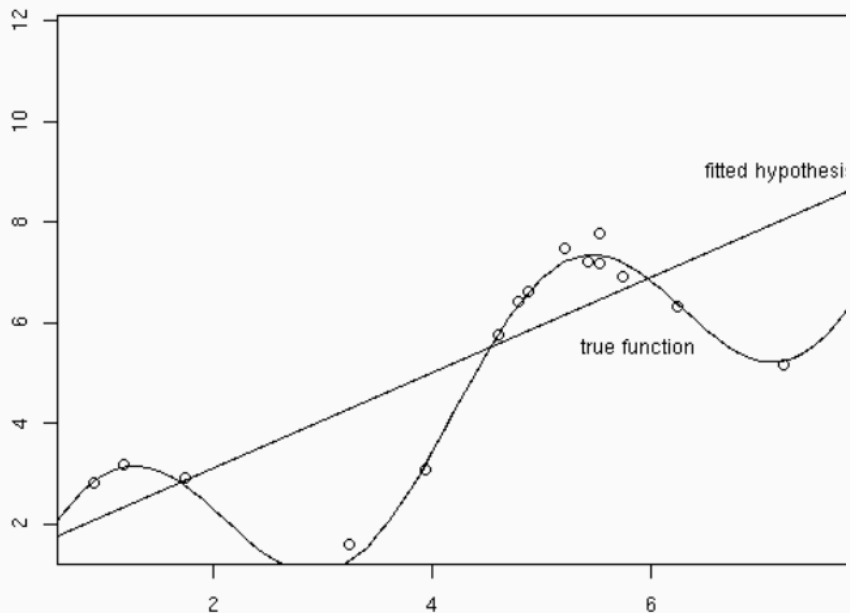
Tom Dietterich, Oregon St



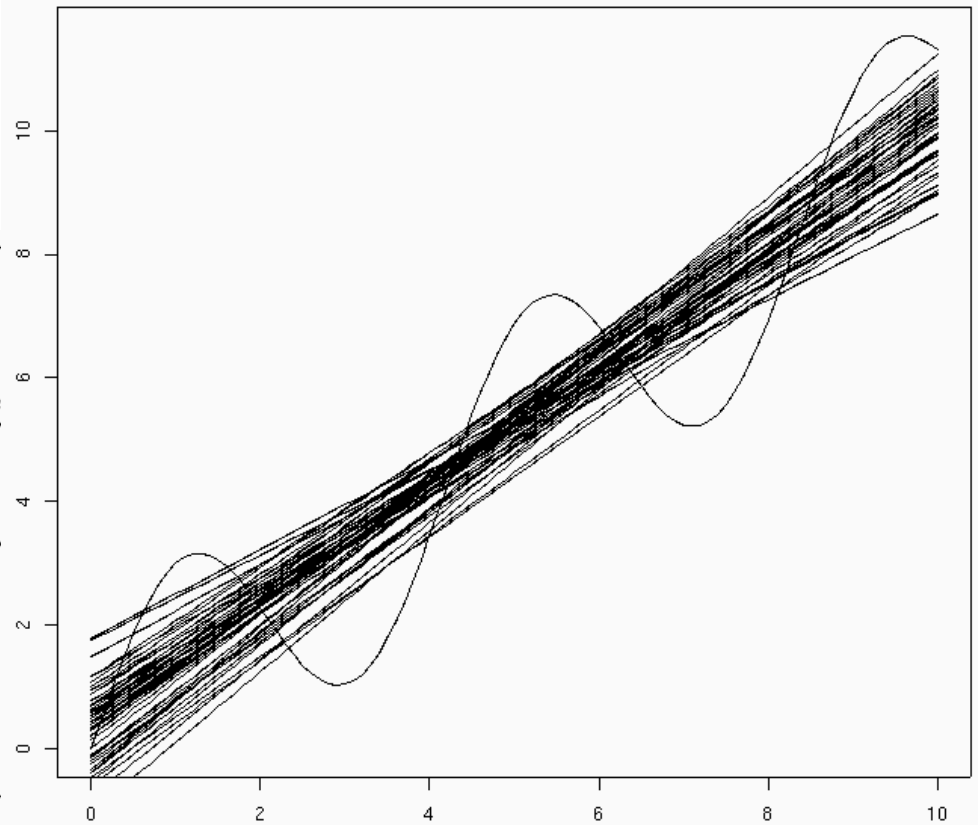
$$y = x + 2 \sin(1.5x) + N(0,0.2)$$

# Example

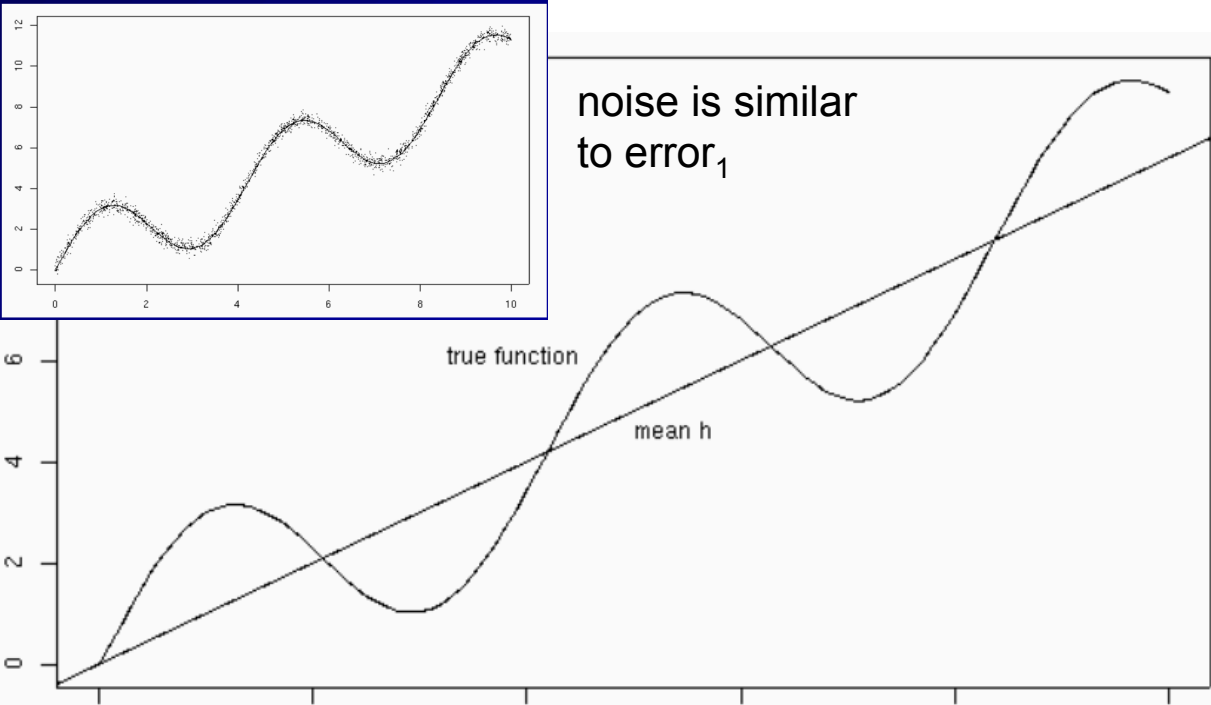
Tom Dietterich, Oregon St



$$y = x + 2 \sin(1.5x) + N(0,0.2)$$

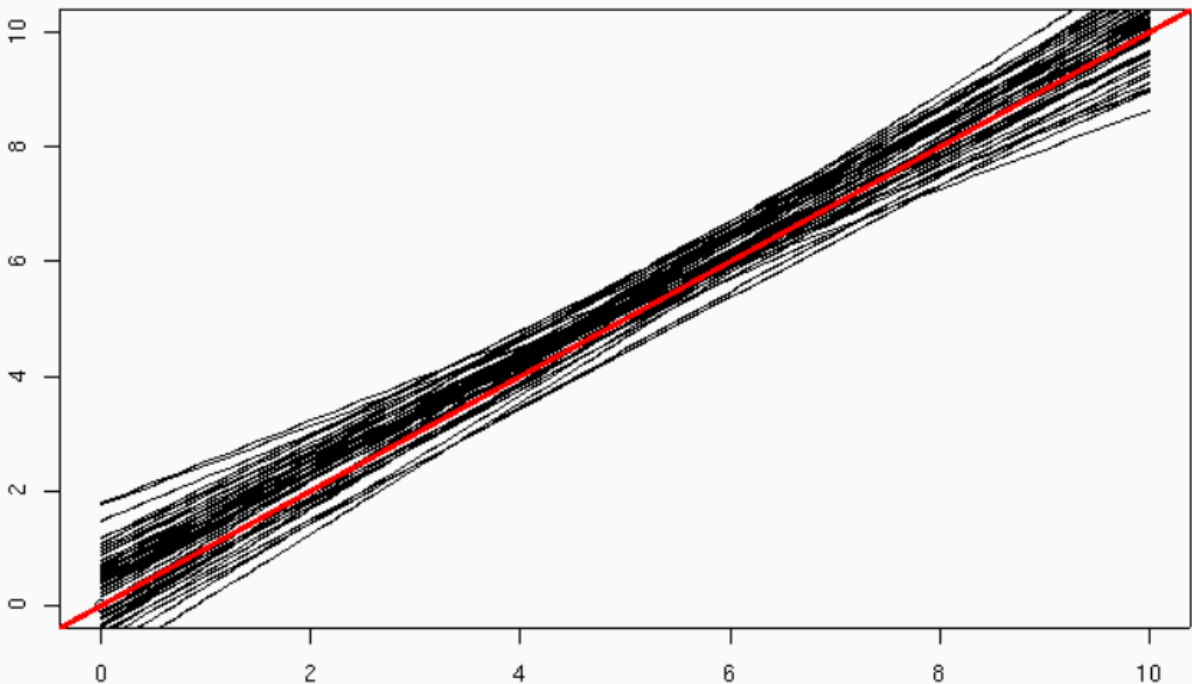


Same experiment, repeated:  
with 50 samples of 20 points each



The true function  $f$  can't be fit perfectly with hypotheses from our class  $H$  (lines)  $\rightarrow \text{Error}_1$

Fix: *more* expressive set of hypotheses  $H$



We don't get the best hypothesis from  $H$  because of noise/small sample size  $\rightarrow \text{Error}_2$

Fix: *less* expressive set of hypotheses  $H$

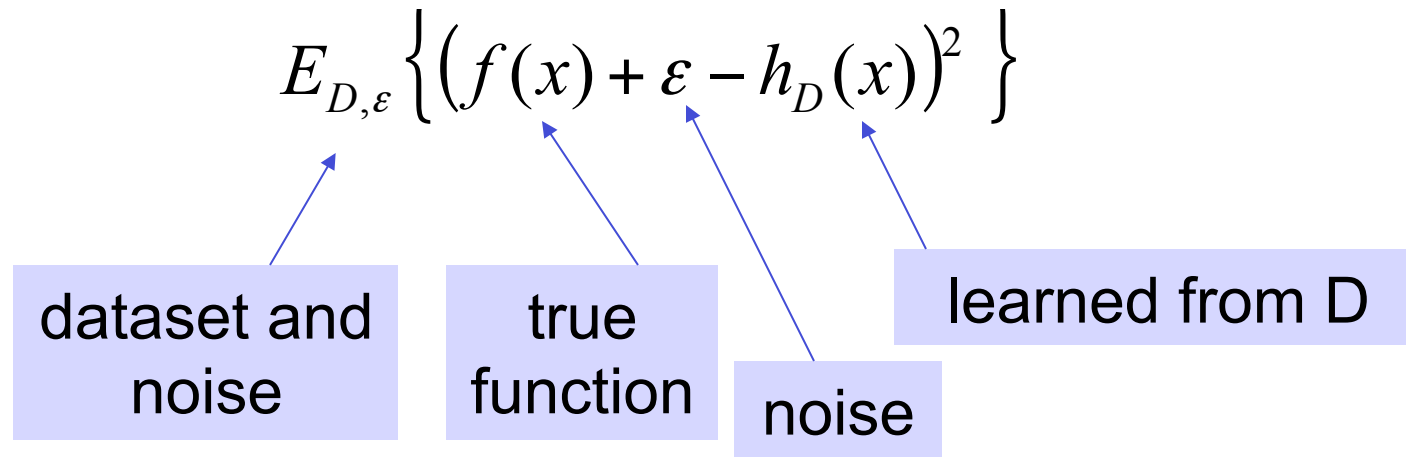
# Bias-Variance Decomposition: Regression



# Bias and variance for regression

- For regression, we can easily **decompose** the error of the learned model into two parts: bias (error 1) and variance (error 2)
  - Bias: the class of models **can't** fit the data.
    - **Fix:** a *more expressive* model class.
  - Variance: the class of models **could** fit the data, but doesn't because it's hard to fit.
    - **Fix:** a *less expressive* model class.

# Bias – Variance decomposition of error



Fix test case  $x$ , then do this experiment:

1. Draw size  $n$  sample  $D=(x_1, y_1), \dots, (x_n, y_n)$
2. Train linear regressor  $h_D$  using  $D$
3. Draw one test example  $(x, f(x)+\epsilon)$
4. Measure squared error of  $h_D$  on that example  $x$

What's the expected error?

# Bias – Variance decomposition of error

Notation - to simplify this

$$f \equiv f(x) + \varepsilon \qquad \hat{y} = \hat{y}_D \equiv h_D(x)$$
$$E_{D,\varepsilon} \left\{ \left( \overbrace{f(x) + \varepsilon}^{\text{true function}} - \overbrace{h_D(x)}^{\text{learned from D}} \right)^2 \right\}$$

dataset and noise

true function

noise

learned from D

$$h \equiv E_D \{h_D(x)\}$$

long-term expectation of learner's prediction on this  $x$  averaged over many data sets  $D$

# Bias – Variance decomposition of error

$$E_{D,\varepsilon} \left\{ (f - \hat{y})^2 \right\}$$

$$= E \left\{ ([f - h] + [h - \hat{y}])^2 \right\}$$

$$= E \left\{ [f - h]^2 + [h - \hat{y}]^2 + 2[f - h][h - \hat{y}] \right\}$$

$$= E \left\{ [f - h]^2 + [h - \hat{y}]^2 + 2[fh - f\hat{y} - h^2 + h\hat{y}] \right\}$$

$$= E[(f - h)^2] + E[(h - \hat{y})^2] + 2(E[fh] - E[f\hat{y}] - E[h^2] + E[h\hat{y}])$$

$$h \equiv E_D \{h_D(x)\}$$

$$\hat{y} = \hat{y}_D \equiv h_D(x)$$

$$f \equiv f(x) + \varepsilon$$

$$E_{D,\varepsilon} \left\{ (f(x) + \varepsilon) * E_D \{h_D(x)\} \right\}$$

$$= E_{D,\varepsilon} \left\{ (f(x) + \varepsilon) * h_D(x) \right\}$$

$$E_{D,\varepsilon} \left\{ E_D \{h_D(x)\} * E_D \{h_D(x)\} \right\}$$

$$= E_{D,\varepsilon} \left\{ E_D \{h_D(x)\} * h_D(x) \right\}$$

# Bias – Variance decomposition of error

$$\begin{aligned} & E_{D,\varepsilon} \left\{ (f - \hat{y})^2 \right\} \\ &= E \left\{ ([f - h] + [h - \hat{y}])^2 \right\} \\ &= E \left\{ [f - h]^2 + [h - \hat{y}]^2 + 2[f - h][h - \hat{y}] \right\} \\ &= E[(f - h)^2] + E[(h - \hat{y})^2] \end{aligned}$$

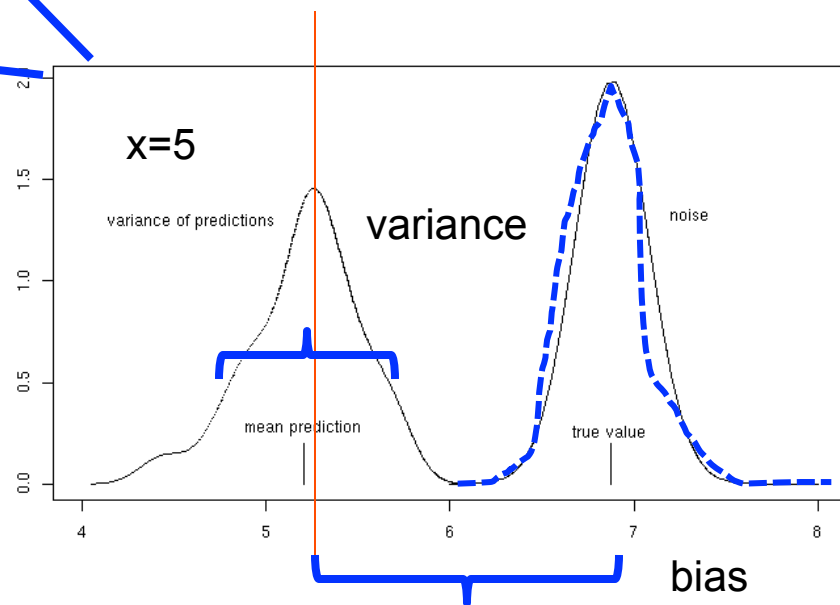
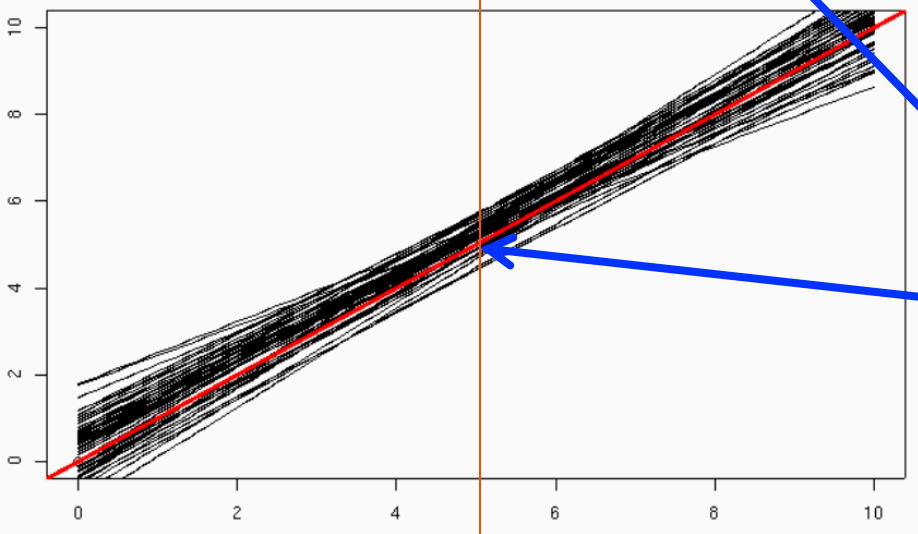
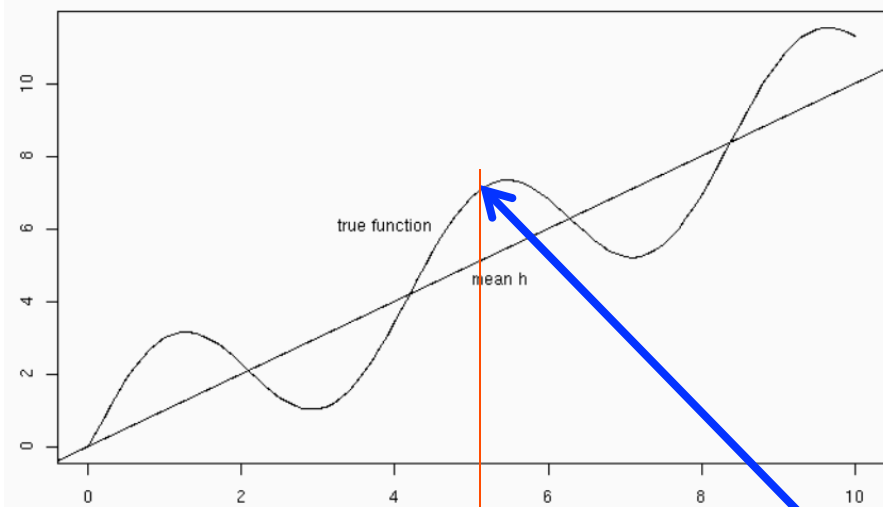
$$\begin{aligned} h &\equiv E_D \{h_D(x)\} \\ \hat{y} &= \hat{y}_D \equiv h_D(x) \\ f &\equiv f(x) + \varepsilon \end{aligned}$$

Squared difference between best possible prediction for  $x$ ,  $f(x)$ , and our “long-term” expectation for what the learner will do if we averaged over many datasets  $D$ ,  $E_D[h_D(x)]$

BIAS<sup>2</sup>

VARIANCE

Squared difference btwn our long-term expectation for the learners performance,  $E_D[h_D(x)]$ , and what we expect in a representative run on a dataset  $D$  ( $\hat{y}$ )



# Bias-variance decomposition

- This is something real that you can (approximately) measure experimentally
  - if you have synthetic data
- Different learners and model classes have different *tradeoffs*
  - large bias/small variance: few features, highly regularized, highly pruned decision trees, large-k k-NN...
  - small bias/high variance: many features, less regularization, unpruned trees, small-k k-NN...

# Bias-Variance Decomposition: Classification



# A generalization of bias-variance decomposition to other loss functions

- “Arbitrary” real-valued loss  $L(y, y')$

But  $L(y, y') = L(y', y)$ ,  $L(y, y) = 0$ ,  
and  $L(y, y') \neq 0$  if  $y \neq y'$

- Define “optimal prediction”:

$$y^* = \operatorname{argmin}_y L(t, y')$$

- Define “main prediction of learner”

$$y_m = y_{m,D} = \operatorname{argmin}_y E_D\{L(y, y')\}$$

$$m = |D|$$

- Define “bias of learner”:

$$\operatorname{Bias}(x) = L(y^*, y_m)$$

- Define “variance of learner”

$$\operatorname{Var}(x) = E_D[L(y_m, y)]$$

- Define “noise for  $x$ ”:

$$N(x) = E_t[L(t, y^*)]$$

Claim:

$$E_{D,t}[L(t, y)] = c_1 N(x) + \operatorname{Bias}(x) + c_2 \operatorname{Var}(x)$$

where

$$c_1 = \Pr_D[y = y^*] - 1$$

$$c_2 = 1 \text{ if } y_m = y^*, -1 \text{ else}$$

For 0/1 loss, the *main prediction* is the most common class predicted by  $h_D(x)$ , weighting  $h$ 's by  $\Pr(D)$

# Bias and variance

- For **classification**, we can also decompose the error of a learned classifier into two terms: bias and variance
  - Bias: the class of models **can't** fit the data.
  - **Fix**: a *more expressive* model class.
  - Variance: the class of models **could** fit the data, but doesn't because it's hard to fit.
  - **Fix**: a *less expressive* model class.

# Bias-Variance Decomposition: Measuring

# Bias-variance decomposition

- This is something real that you can (approximately) measure experimentally
  - if you have synthetic data
  - ...or if you're clever
  - You need to somehow approximate  $E_D\{h_D(x)\}$
  - I.e., construct many variants of the dataset  $D$

# Background: “Bootstrap” sampling

- **Input:** dataset  $D$
- **Output:** many variants of  $D$ :  $D_1, \dots, D_T$
- For  $t=1, \dots, T$ :
  - $D_t = \{ \}$
  - For  $i=1 \dots |D|$ :
    - Pick  $(\mathbf{x}, y)$  uniformly at random from  $D$  (i.e., **with** replacement) and add it to  $D_t$
    - Some examples never get picked ( $\sim 37\%$ )
    - Some are picked 2x, 3x, ....

# Measuring Bias-Variance with “Bootstrap” sampling

- Create  $B$  bootstrap variants of  $D$  (approximate many draws of  $D$ )
- For each bootstrap dataset
  - $T_b$  is the dataset;  $U_b$  are the “out of bag” examples
  - Train a hypothesis  $h_b$  on  $T_b$
  - Test  $h_b$  on each  $x$  in  $U_b$
- Now for each  $(\mathbf{x}, y)$  example we have many predictions  $h_1(x), h_2(x), \dots$  so we can estimate (ignoring noise)
  - **variance**: ordinary variance of  $h_1(x), \dots, h_n(x)$
  - **bias**:  $\text{average}(h_1(x), \dots, h_n(x)) - y$

# Applying Bias-Variance Analysis

- By measuring the bias and variance on a problem, we can determine how to improve our model
  - If bias is high, we need to allow our model to be more complex
  - If variance is high, we need to reduce the complexity of the model
- Bias-variance analysis also suggests a way to reduce variance: *bagging* (later)

# Bagging



# Bootstrap Aggregation (Bagging)

- Use the **bootstrap** to create  $B$  variants of  $D$
- Learn a classifier from **each variant**
- **Vote** the learned classifiers to predict on a test example

# Bagging (bootstrap aggregation)

- Breaking it down:

- input: dataset  $D$  and YFCL
- output: a classifier  $h_{D-BAG}$

Note that you can use *any* learner you like!

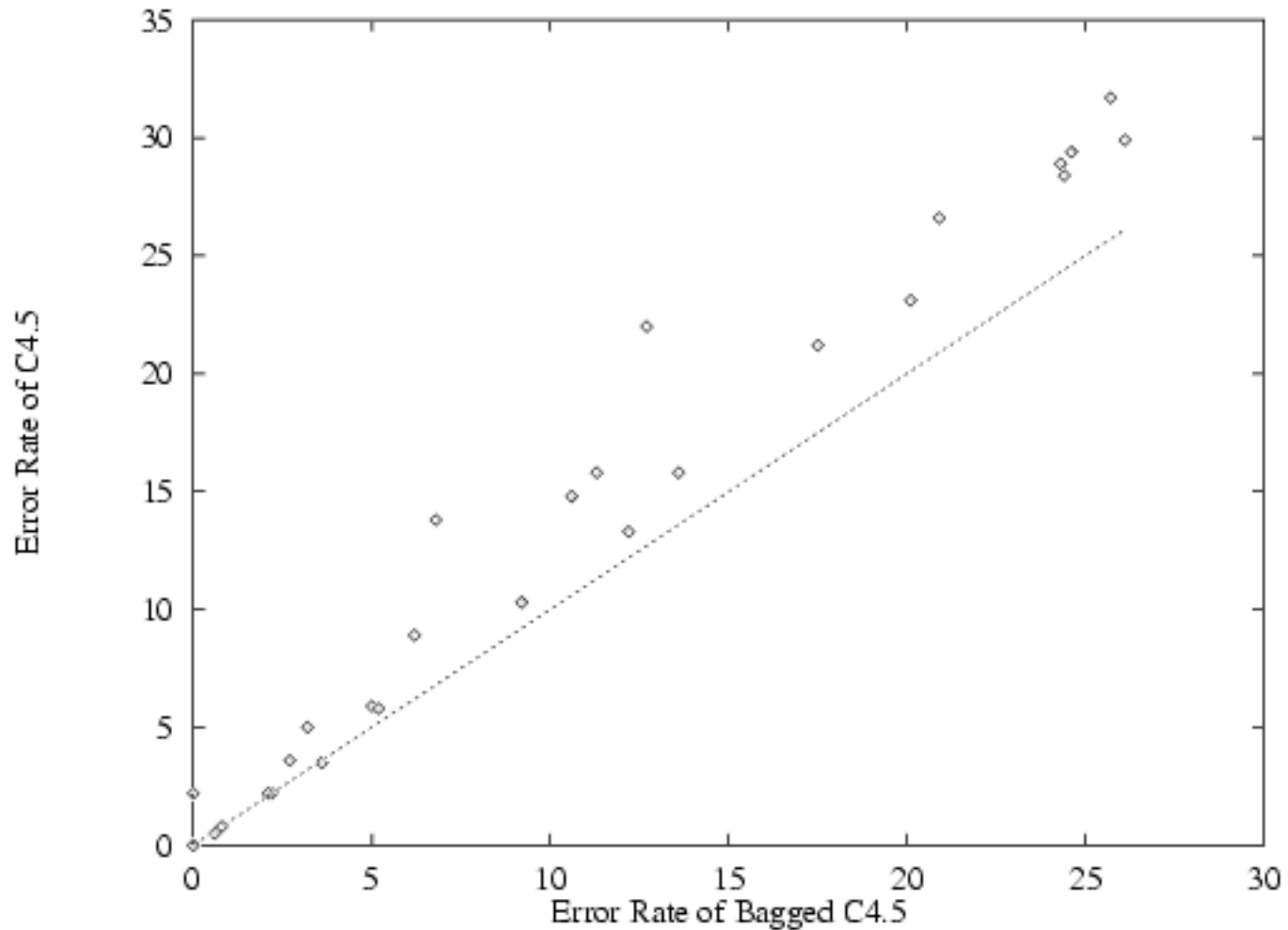
- use bootstrap to construct variants  $D_1, \dots, D_T$
- for  $t=1, \dots, T$ : train YFCL on  $D_t$  to get  $h_t$

You can also test  $h_t$  on the “out of bag” examples

- to classify  $x$  with  $h_{D-BAG}$ 
  - classify  $x$  with  $h_1, \dots, h_T$  and predict the most frequently predicted class for  $x$  (majority vote)

# Experiments

Freund and Schapire

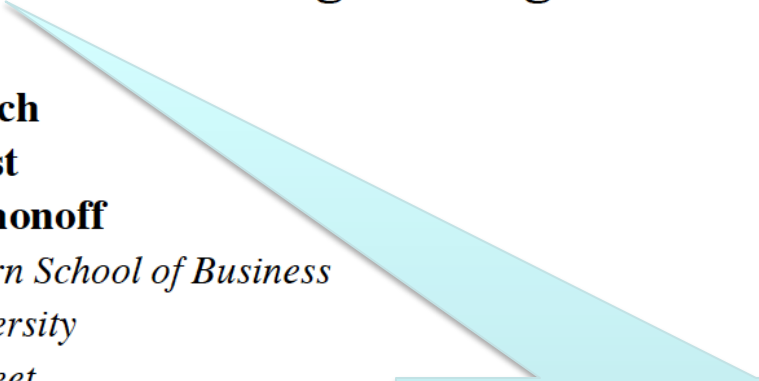


# Tree Induction vs. Logistic Regression: A Learning-Curve Analysis

**Claudia Perlich**  
**Foster Provost**  
**Jeffrey S. Simonoff**

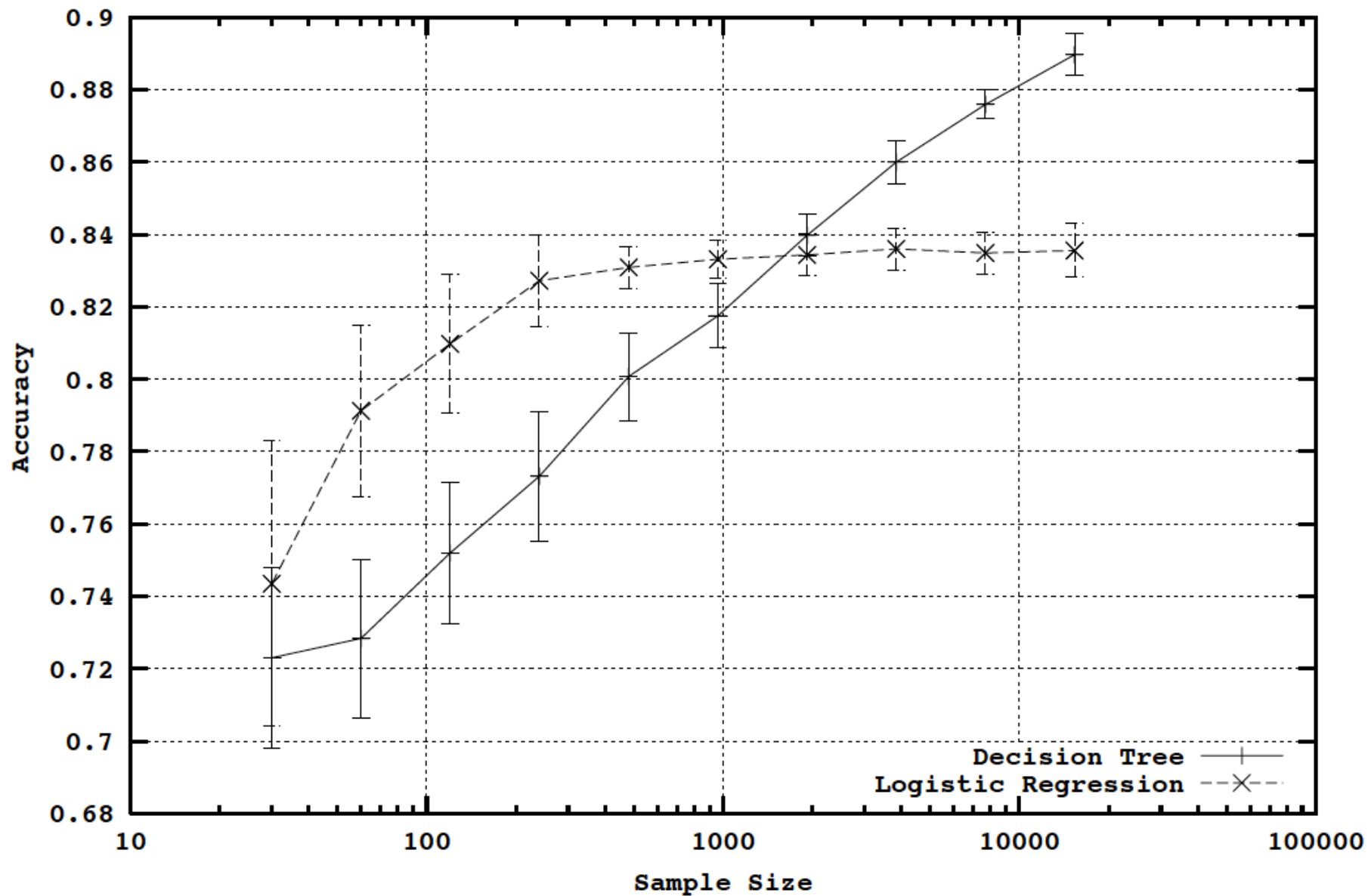
*Leonard N. Stern School of Business*  
*New York University*  
*44 West 4th Street*  
*New York, NY 10012*

CERLICH@STERN.NYU.EDU  
FPROVOST@STERN.NYU.EDU  
JSIMONOF@STERN.NYU.EDU



Bagged, minimally pruned decision trees

Learning Curve of Californian Housing Data



Generally, bagged  
decision trees  
outperform the linear  
classifier eventually if  
the data is large  
enough and clean  
enough.

Data set	Winner AUR	Winner Acc	Max-AUR	Result
Nurse	none	none	1	Indistinguishable
Mushrooms	none	none	1	Indistinguishable
Optdigit	none	none	0.99	Indistinguishable
Letter-V	C4	C4	0.99	C4 dominates
Letter-A	C4	C4	0.99	C4 crosses
Intrusion	C4	C4	0.99	C4 dominates
DNA	C4	C4	0.99	C4 dominates
Coverttype	C4	C4	0.99	C4 crosses
Telecom	C4	C4	0.98	C4 dominates
Pendigit	C4	C4	0.98	C4 dominates
Pageblock	C4	C4	0.98	C4 crosses
CarEval	none	C4	0.98	C4 crosses
Spam	C4	C4	0.97	C4 dominates
Chess	C4	C4	0.95	C4 dominates
CalHous	C4	C4	0.95	C4 crosses
Ailerons	none	C4	0.95	C4 crosses
Firm	LR	LR	0.93	LR crosses
Credit	C4	C4	0.93	C4 dominates
Adult	LR	C4	0.9	Mixed
Connects	C4	none	0.87	C4 crosses
Move	C4	C4	0.85	C4 dominates
Downsize	C4	C4	0.85	C4 crosses
Coding	C4	C4	0.85	C4 crosses
German	LR	LR	0.8	LR dominates
Diabetes	LR	LR	0.8	LR dominates
Bookbinder	LR	LR	0.8	LR crosses
Bacteria	none	C4	0.79	C4 crosses
Yeast	none	none	0.78	Indistinguishable
Patent	C4	C4	0.75	C4 crosses
Contra	none	none	0.73	Indistinguishable
IntShop	LR	LR	0.7	LR crosses
IntCensor	LR	LR	0.7	LR dominates
Insurance	none	none	0.7	Indistinguishable
IntPriv	LR	none	0.66	LR crosses
Mailing	LR	none	0.61	LR dominates
Abalone	LR	LR	0.56	LR dominates

# Bagging (bootstrap aggregation)

- Experimentally:
  - especially with minimal pruning: decision trees have low *bias* but high *variance*.
  - bagging usually improves performance for decision trees and similar methods
  - It reduces *variance* without increasing the bias (much).

# More detail on bias-variance and bagging for classification



# A generalization of bias-variance decomposition to other loss functions

- “Arbitrary” real-valued loss  $L(y, y')$

But  $L(y, y') = L(y', y)$ ,  $L(y, y) = 0$ ,  
and  $L(y, y') \neq 0$  if  $y \neq y'$

- Define “optimal prediction”:

$$y^* = \operatorname{argmin}_y L(t, y')$$

- Define “main prediction of learner”

$$y_m = y_{m,D} = \operatorname{argmin}_y E_D\{L(y, y')\}$$

$$m = |D|$$

- Define “bias of learner”:

$$\operatorname{Bias}(x) = L(y^*, y_m)$$

- Define “variance of learner”

$$\operatorname{Var}(x) = E_D[L(y_m, y)]$$

- Define “noise for  $x$ ”:

$$N(x) = E_t[L(t, y^*)]$$

Claim:

$$E_{D,t}[L(t, y)] = c_1 N(x) + \operatorname{Bias}(x) + c_2 \operatorname{Var}(x)$$

where

$$c_1 = \Pr_D[y = y^*] - 1$$

$$c_2 = 1 \text{ if } y_m = y^*, -1 \text{ else}$$

For 0/1 loss, the *main prediction* is the most common class predicted by  $h_D(x)$ , weighting  $h$ 's by  $\Pr(D)$

# More detail on Domingos' s model

- *Noisy channel*:  $y_i = \text{noise}(f(x_i))$ 
  - $f(x_i)$  is true label of  $x_i$
  - Noise  $\text{noise}(.)$  may change  $y \rightarrow y'$
- $h=h_D$  is learned hypothesis
  - from  $D=\{(x_1, y_1), \dots, (x_m, y_m)\}$
- for test case  $(x^*, y^*)$ , and predicted label  $h(x^*)$ , loss is  $L(h(x^*), y^*)$ 
  - For instance,  $L(h(x^*), y^*) = 1$  if error, else 0

# More detail on Domingos' s model

- We want to decompose  $E_{D,P}\{L(h(x^*),y^*)\}$  where  $m$  is size of  $D$ ,  $(x^*,y^*)\sim P$
- *Main prediction of learner* is  $y_m(x^*)$ 
  - $y_m(x^*) = \operatorname{argmin}_{y'} E_{D,P}\{L(h(x^*),y')\}$
  - $y_m(x^*) =$  “most common”  $h_D(x^*)$  among all possible  $D$ ' s, weighted by  $\Pr(D)$
- *Bias* is  $B(x^*) = L(y_m(x^*) , f(x^*))$
- *Variance* is  $V(x^*) = E_{D,P}\{L(h_D(x^*) , y_m(x^*) )$
- *Noise* is  $N(x^*)= L(y^*, f(x^*))$

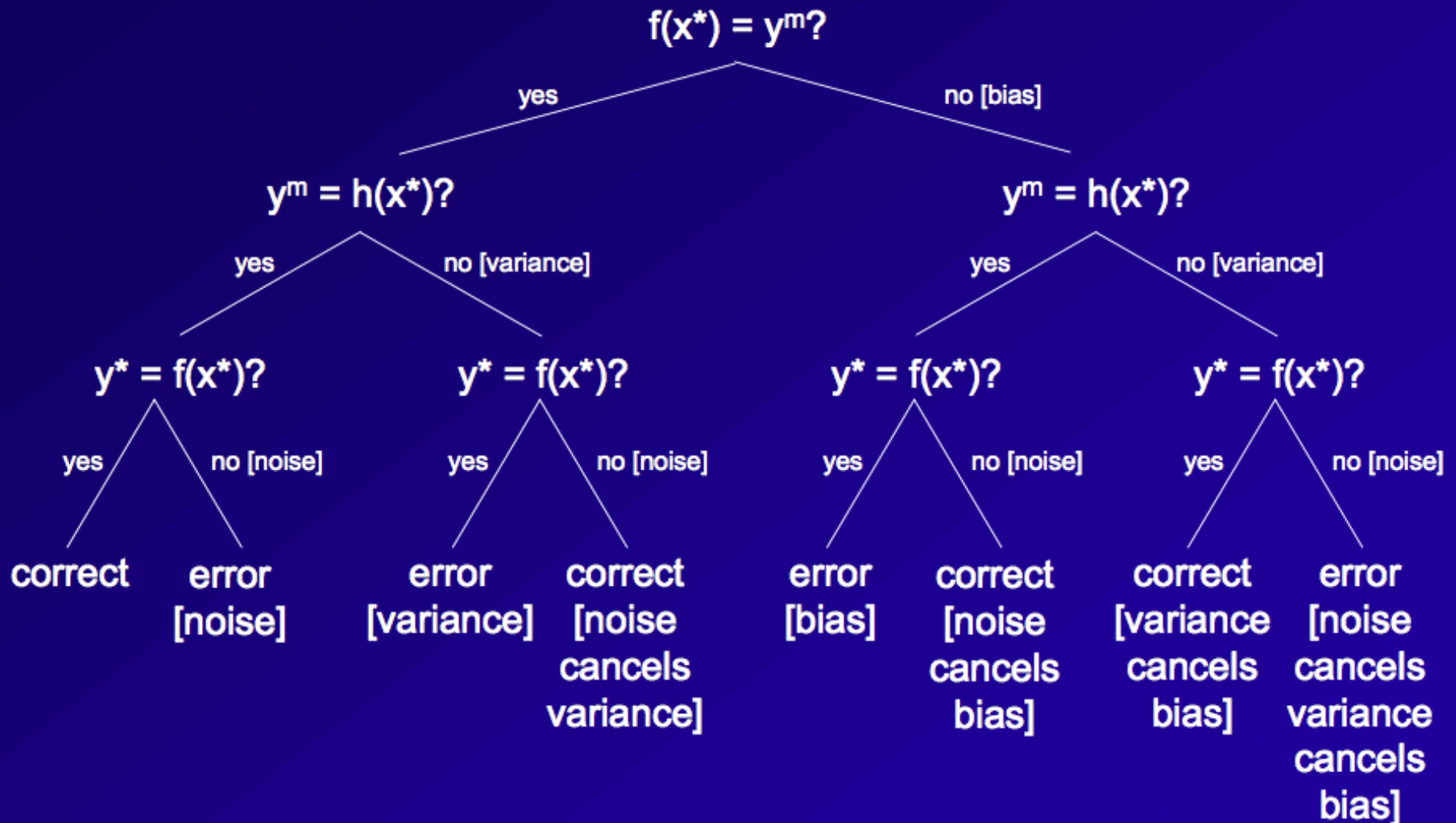
# More detail on Domingos' s model

- We want to decompose  $E_{D,P}\{L(h(x^*), y^*)\}$
- *Main prediction of learner* is  $y_m(x^*)$ 
  - “most common”  $h_D(x^*)$  over  $D$ 's for 0/1 loss
- *Bias* is  $B(x^*) = L(y_m(x^*), f(x^*))$ 
  - main prediction vs true label
- *Variance* is  $V(x^*) = E_{D,P}\{L(h_D(x^*), y_m(x^*))\}$ 
  - this hypothesis vs main prediction
- *Noise* is  $N(x^*) = L(y^*, f(x^*))$ 
  - true label vs observed label

# More detail on Domingos' s model

- We will decompose  $E_{D,P}\{L(h(x^*), y^*)\}$  into
  - *Bias* is  $B(x^*) = L(y_m(x^*), f(x^*))$ 
    - main prediction vs true label
    - this is 0/1, *not* a random variable
  - *Variance* is  $V(x^*) = E_{D,P}\{L(h_D(x^*), y_m(x^*))\}$ 
    - this hypothesis vs main prediction
  - *Noise* is  $N(x^*) = L(y^*, f(x^*))$ 
    - true label vs observed label

# Case analysis of error



# Analysis of error: unbiased case

■ Let  $P(y^* \neq f(x^*)) = N(x^*) = \tau$

■ Let  $P(y^m \neq h(x^*)) = V(x^*) = \sigma$

■ If  $(f(x^*) = y^m)$ , then we suffer a loss if exactly one of these events occurs:

$$\begin{aligned} L(h(x^*), y^*) &= \tau(1-\sigma) + \sigma(1-\tau) \\ &= \tau + \sigma - 2\tau\sigma \\ &= N(x^*) + V(x^*) - 2 N(x^*) V(x^*) \end{aligned}$$

Variance  
but no  
noise

Main  
prediction  
is correct

Noise but  
no  
variance

# Analysis of error: biased case

- Let  $P(y^* \neq f(x^*)) = N(x^*) = \tau$
- Let  $P(y^m \neq h(x^*)) = V(x^*) = \sigma$
- If  $(f(x^*) \neq y^m)$ , then we suffer a loss if either both or neither of these events occurs:

$$\begin{aligned} L(h(x^*), y^*) &= \tau\sigma + (1-\sigma)(1-\tau) \\ &= 1 - (\tau + \sigma - 2\tau\sigma) \\ &= B(x^*) - [N(x^*) + V(x^*) - 2 N(x^*) V(x^*)] \end{aligned}$$

No noise,  
no  
variance

Main  
prediction  
is wrong

Noise  
and  
variance



# Analysis of error: overall

$$E[ L(h(x^*), y^*) ] =$$

$$\text{if } B(x^*) = 1: B(x^*) - [N(x^*) + V(x^*) - 2 N(x^*) V(x^*)]$$

$$\text{if } B(x^*) = 0: B(x^*) + [N(x^*) + V(x^*) - 2 N(x^*) V(x^*)]$$

Hopefully we'll be in this case more often, if we've chosen a good classifier

Interaction terms are usually small

# Analysis of error: without noise

which is hard to estimate anyway

$$\begin{aligned} E[ L(h(x^*), y^*) ] = \\ \text{if } B(x^*) = 1: B(x^*) - V(x^*) \\ \text{if } B(x^*) = 0: B(x^*) + V(x^*) \end{aligned}$$

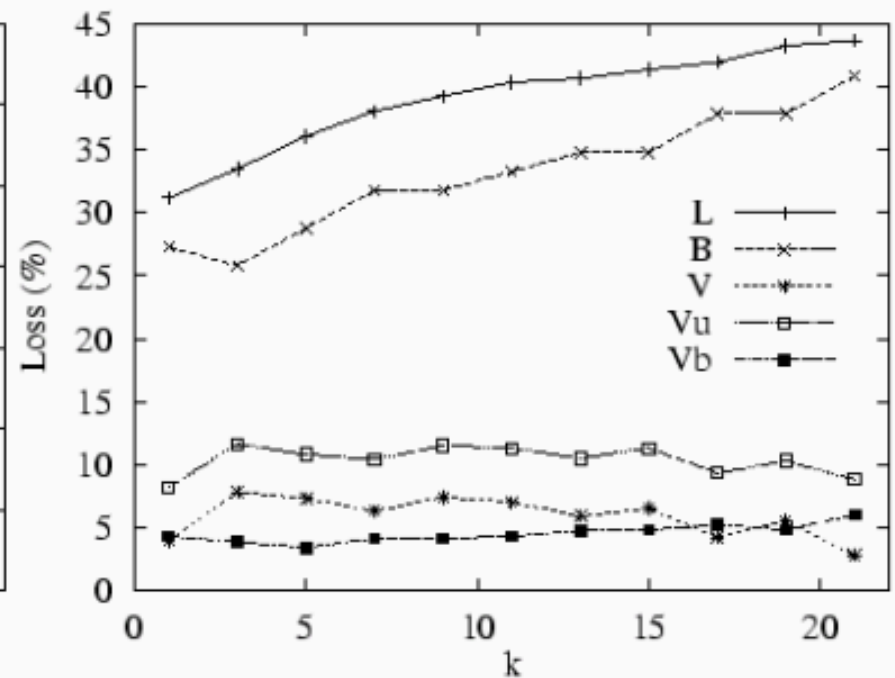
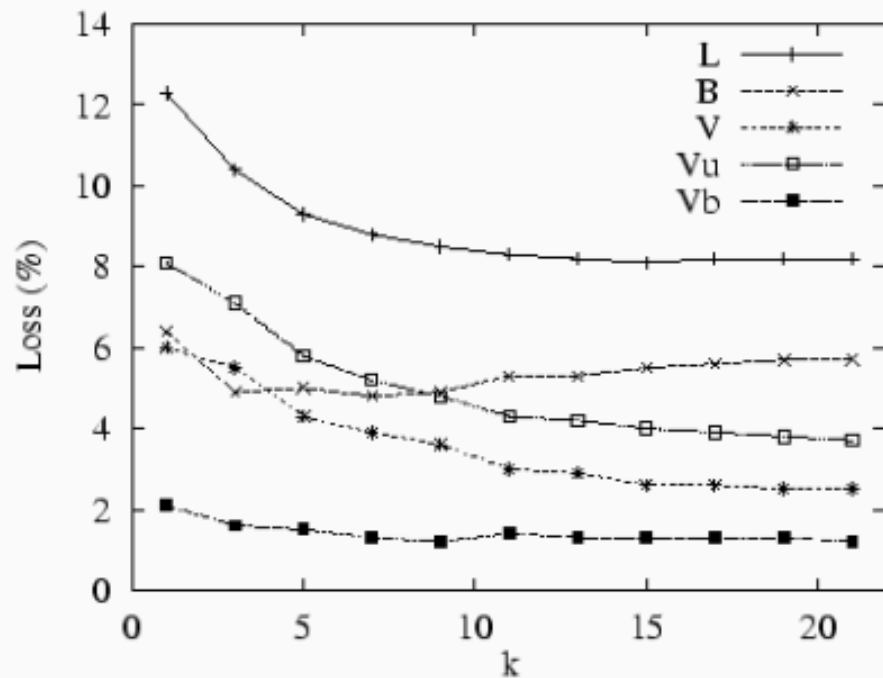
$V_b$

$V_u$

As with regression, we can experimentally approximately measure bias and variance with bootstrap replicates

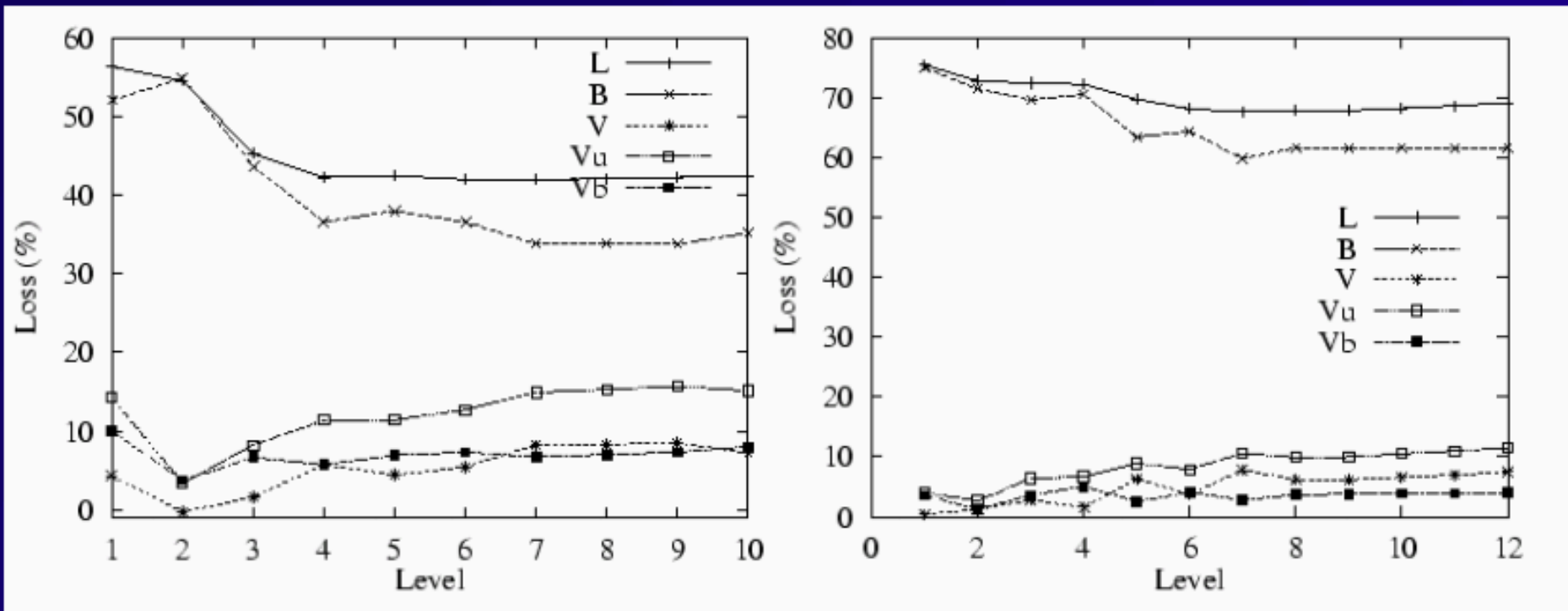
Typically break variance down into *biased variance*,  $V_b$ , and *unbiased variance*,  $V_u$ .

# K-NN Experiments



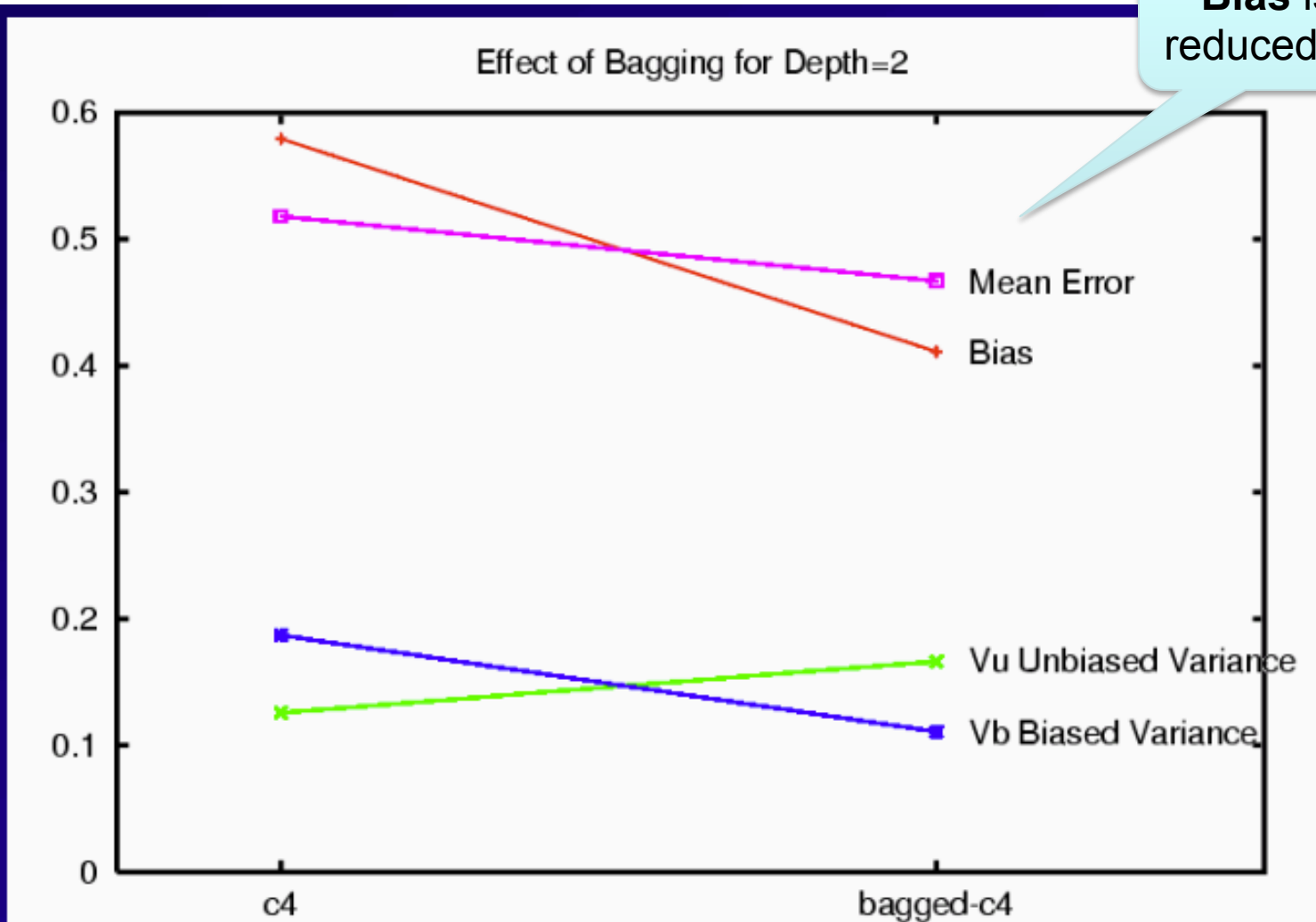
- Chess (left): Increasing K primarily reduces Vu
- Audiology (right): Increasing K primarily increases B.

# Tree Experiments



■ Glass (left), Primary tumor (right): deeper trees have lower B, higher Vu

# Tree “stump” experiments (depth 2)



# Large tree experiments (depth 10)

