

# Introduction to Machine Learning

Thomas G. Dietterich  
tgd@eecs.oregonstate.edu

# Outline

- What is Machine Learning?
- Introduction to Supervised Learning: Linear Methods
- Overfitting, Regularization, and the Bias-Variance Tradeoff
- Review and Summary

# Machine Learning: The Original Motivation

## Traditional Software Process

- Interview the experts
- Create an algorithm that automates their process

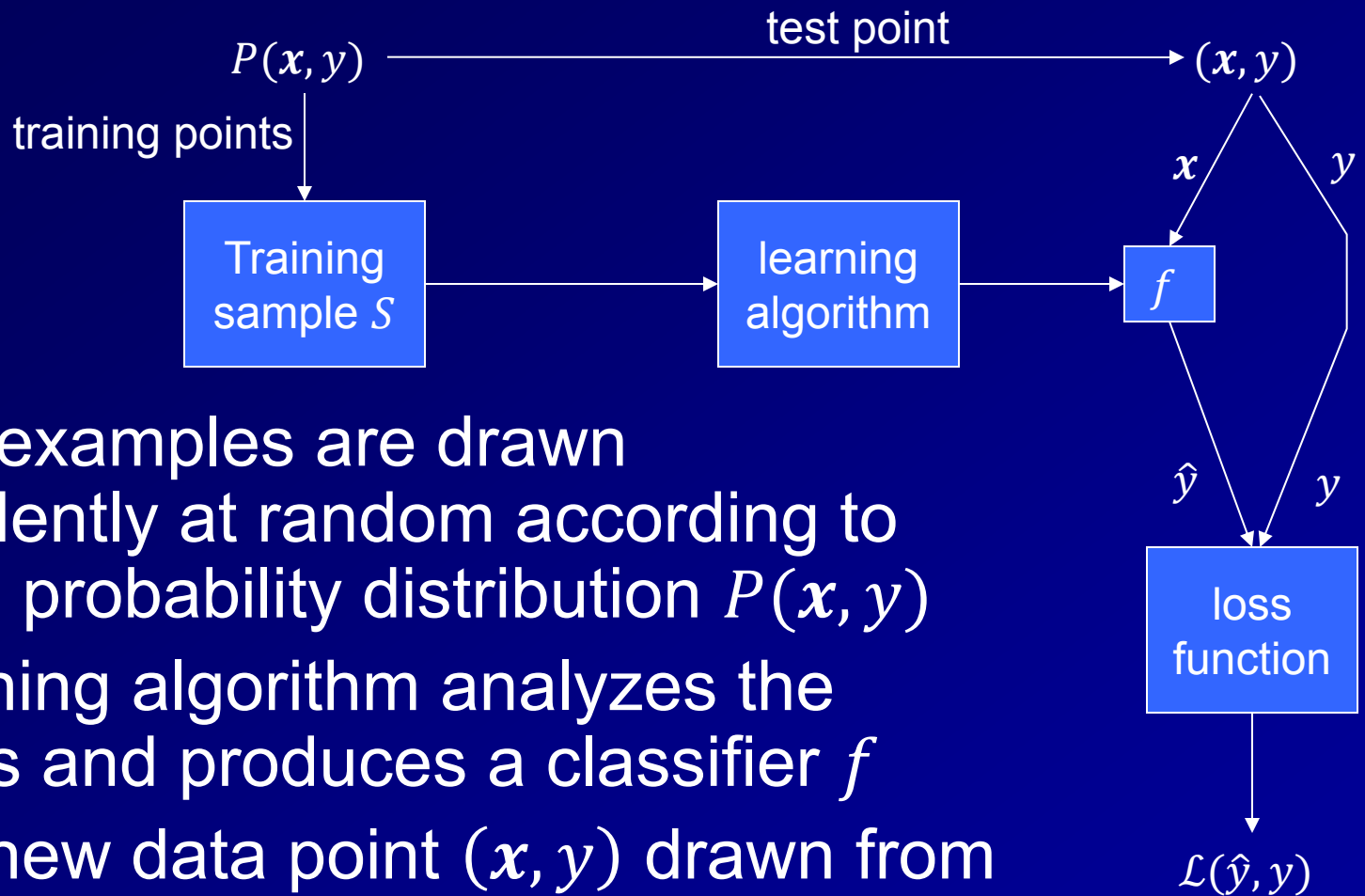
## Machine Learning Process

- Collect input-output examples from the experts
- Learn a function to map from the input to the output

# Supervised Learning

- Given: Training examples  $(x_i, f(x_i))$  for some unknown function  $f$ .
- Find: A good approximation to  $f$ .
  
- Example Applications
  - Handwriting recognition
    - $x$ : data from pen motion
    - $f(x)$ : letter of the alphabet
  - Disease Diagnosis
    - $x$ : properties of patient (symptoms, lab tests)
    - $f(x)$ : disease (or maybe, recommended therapy)
  - Face Recognition
    - $x$ : bitmap picture of person's face
    - $f(x)$ : name of person
  - Spam Detection
    - $x$ : email message
    - $f(x)$ : spam or not spam

# Formal Setting



- Training examples are drawn independently at random according to unknown probability distribution  $P(x, y)$
- The learning algorithm analyzes the examples and produces a classifier  $f$
- Given a new data point  $(x, y)$  drawn from  $P$ , the classifier is given  $x$  and predicts  $\hat{y} = f(x)$
- The loss  $\mathcal{L}(\hat{y}, y)$  is then measured
- Goal of the learning algorithm: Find the  $f$  that minimizes the *expected loss*

# Formal Version of Spam Detection

- $P(x, y)$ : distribution of email messages  $x$  and their true labels  $y$  (“spam” or “not spam”)
- training sample: a set of email messages that have been labeled by the user
- learning algorithm: what we study in MLSS!
- $f$ : the classifier output by the learning algorithm
- test point: A new email message  $x$  (with its true, but hidden, label  $y$ )
- loss function  $\mathcal{L}(\hat{y}, y)$  :

predicted label $\hat{y}$	true label $y$	
	spam	not spam
spam	0	10
not spam	1	0

# Three Main Approaches to Machine Learning

- Learn a classifier: a function  $f$ .
- Learn a conditional distribution: a conditional distribution  $P(y | \mathbf{x})$
- Learn the joint probability distribution:  $P(\mathbf{x}, y)$
- We will study one example of each method:
  - Learn a classifier: The Perceptron algorithm
  - Learn a conditional distribution: Logistic regression
  - Learn the joint distribution: Linear discriminant analysis

# Linear Threshold Units

$$h(\mathbf{x}) = \begin{cases} +1 & \text{if } w_1x_1 + \dots + w_nx_n \geq w_0 \\ -1 & \text{otherwise} \end{cases}$$

- We assume that each feature  $x_j$  and each weight  $w_j$  is a real
- We will study three different algorithms for learning linear threshold units:
  - Perceptron: function
  - Logistic Regression: conditional distribution
  - Linear Discriminant Analysis: joint distribution

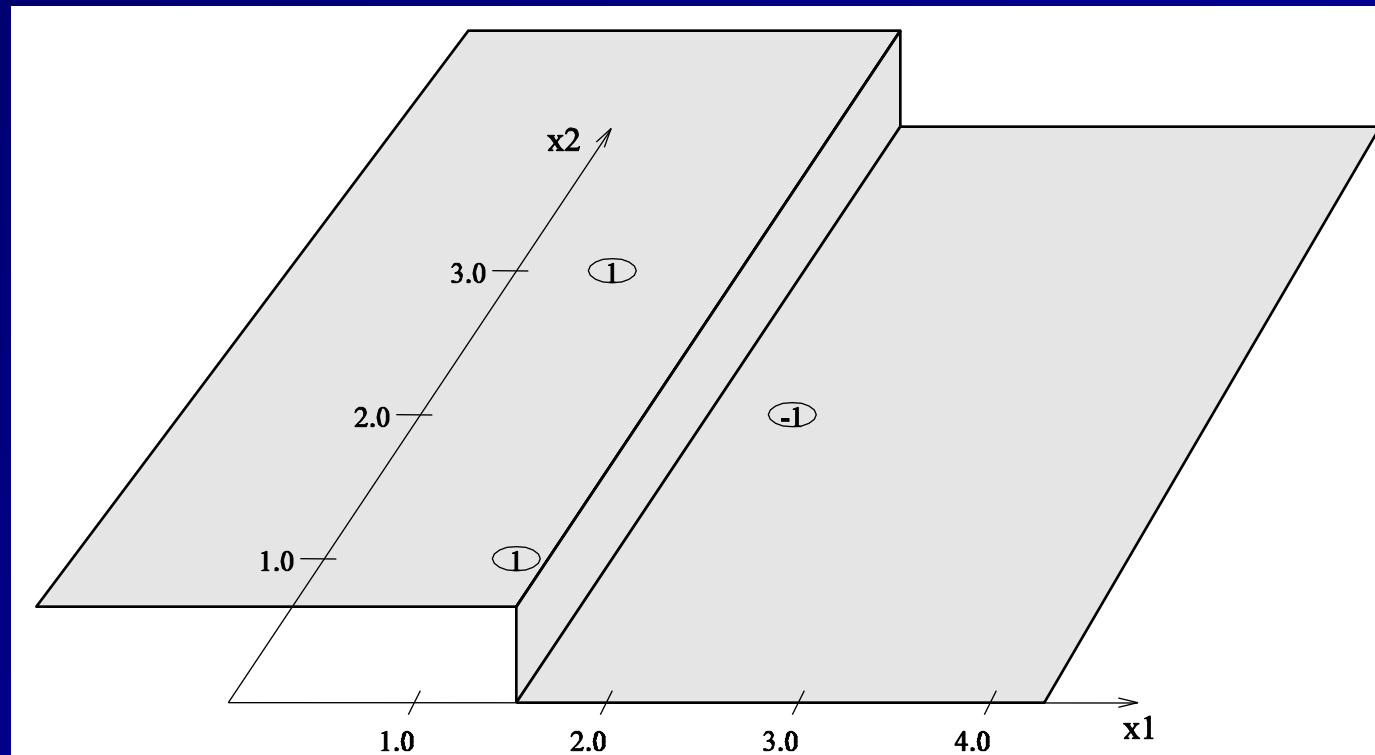


# A canonical representation

- Given a training example of the form  
 $(\langle x_1, x_2, x_3, x_4 \rangle, y)$
- transform it to  
 $(\langle -1, x_1, x_2, x_3, x_4 \rangle, y)$
- The parameter vector will then be  
 $(w_0, w_1, w_2, w_3, w_4)$
- We will call the *unthresholded* hypothesis  $u(\mathbf{x}, \mathbf{w})$   
 $u(\mathbf{x}, \mathbf{w}) = \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^\top \mathbf{x}$
- Each hypothesis can be written  
 $h(\mathbf{x}) = \text{sgn}(u(\mathbf{x}, \mathbf{w}))$
- Our goal is to find  $\mathbf{w}$

# Geometrical View

- Consider three training examples:  $(\langle 1.0, 1.0 \rangle, +1)$   
 $(\langle 0.5, 3.0 \rangle, +1)$   
 $(\langle 2.0, 2.0 \rangle, -1)$
- We want a classifier that looks like the following:



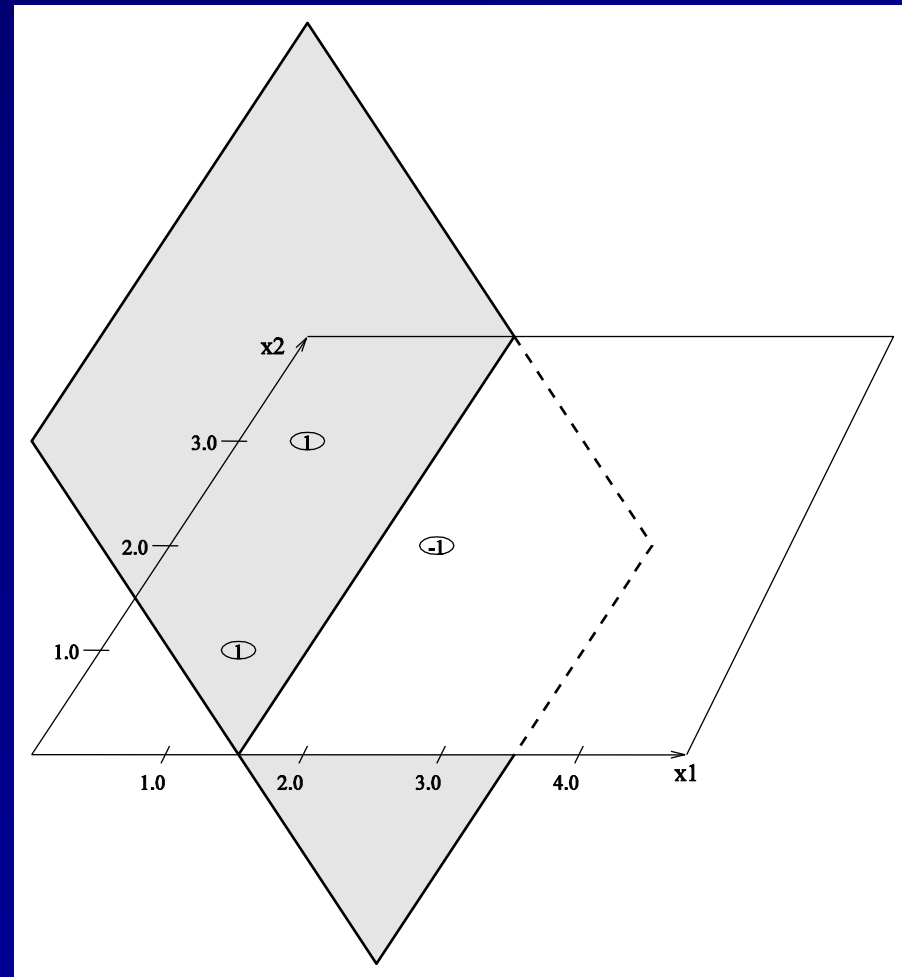
# The Unthresholded Discriminant Function is a Hyperplane

■ The equation

$$u(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

is a plane

$$\hat{y} = \begin{cases} +1 & \text{if } u(\mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$



# Machine Learning == Optimization

## ■ Given:

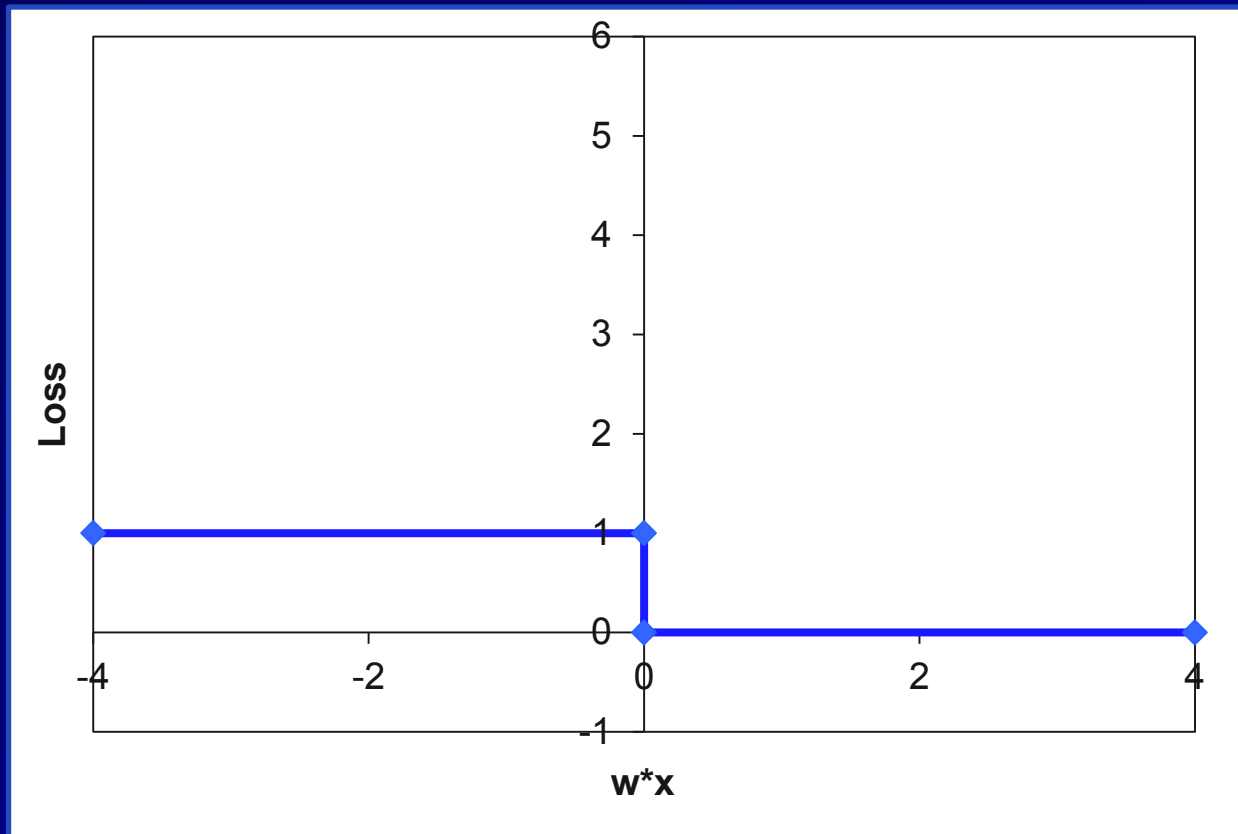
- A set of  $N$  training examples  
 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- A loss function  $\mathcal{L}$

## ■ Find:

- The weight vector  $\mathbf{w}$  that minimizes the expected loss on the training data

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\text{sgn}(\mathbf{w}^\top \mathbf{x}_i), y_i)$$

# Problem: Step-wise Constant Loss Function

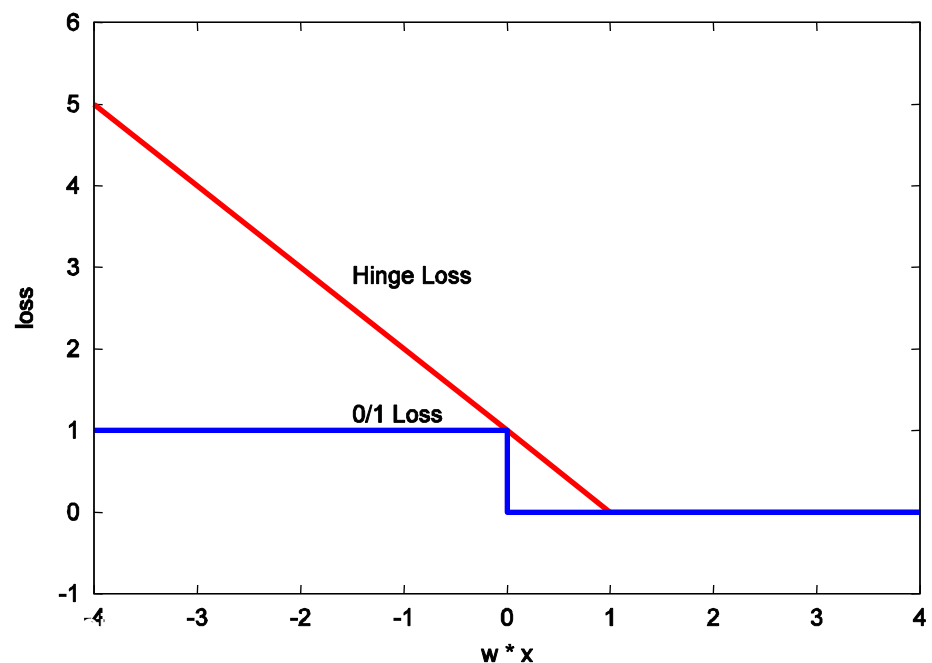


# Approximating the expected loss by a smooth function

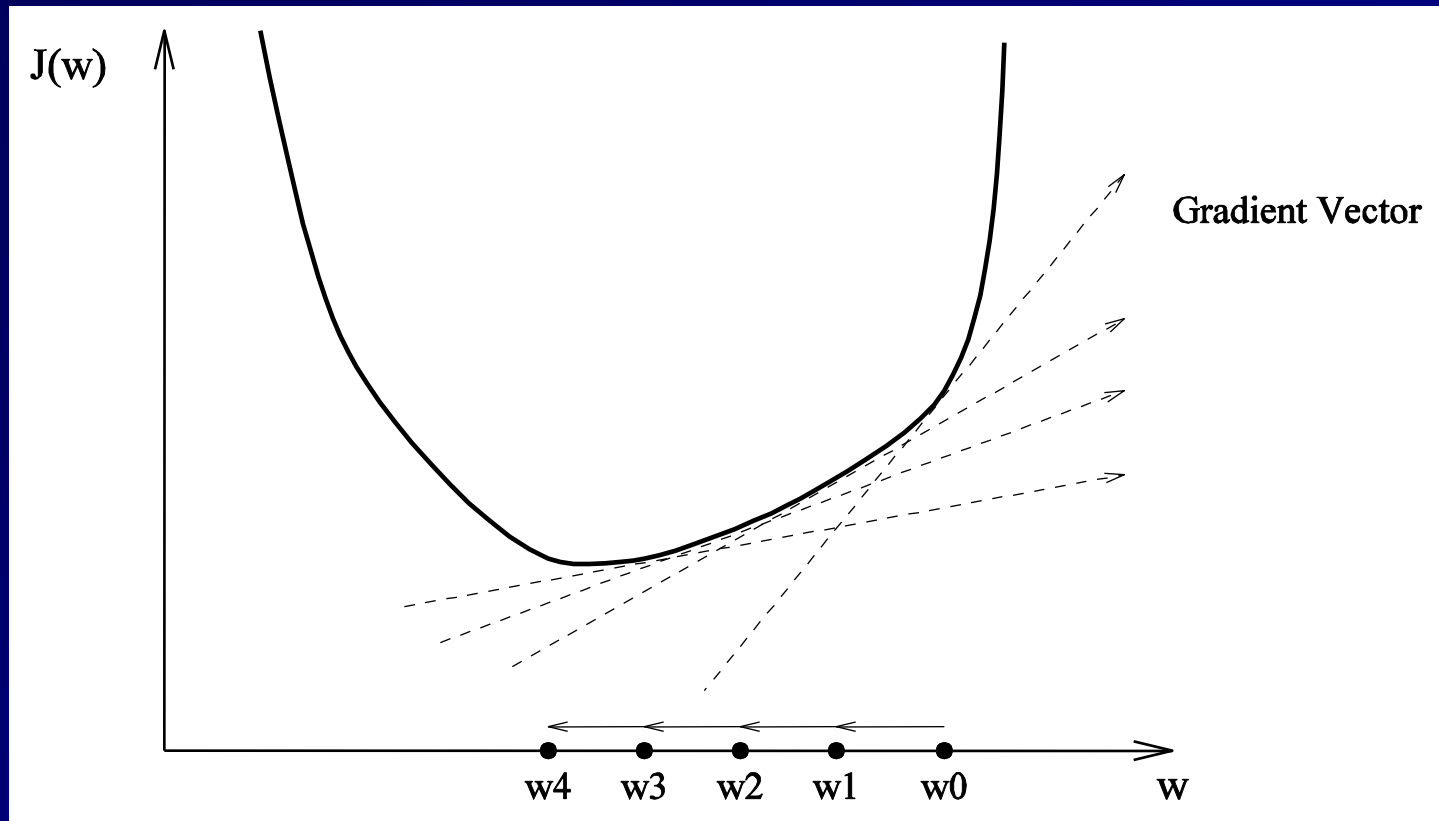
- Simplify the optimization problem by replacing the original objective function by a surrogate loss function. For example, consider the *hinge loss*:

$$\tilde{J}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x})$$

When  $y = 1$ :



# Minimizing $\tilde{J}$ by Gradient Descent Search



- Start with weight vector  $w^0$
- Compute gradient  $\nabla \tilde{J}(w^0) = \left( \frac{\partial \tilde{J}(w^0)}{\partial w_0}, \frac{\partial \tilde{J}(w^0)}{\partial w_1}, \dots, \frac{\partial \tilde{J}(w^0)}{\partial w_n} \right)$
- Compute  $w^1 = w^0 - \eta \nabla \tilde{J}(w^0)$   
where  $\eta$  is a “step size” parameter
- Repeat until convergence

# Gradient of the Hinge Loss

Let  $\tilde{J}_i(\mathbf{w}) = \max(0, -y_i \mathbf{w}^\top \mathbf{x})$

$$\frac{\partial \tilde{J}(\mathbf{w})}{\partial w_k} = \frac{\partial}{\partial w_k} \left( \frac{1}{N} \sum_{i=1}^N \tilde{J}_i(\mathbf{w}) \right) = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial w_k} \tilde{J}_i(\mathbf{w})$$

$$\sum_{i=1}^N \frac{\partial}{\partial w_k} \tilde{J}_i(\mathbf{w}) = \frac{\partial}{\partial w_k} \max \left( 0, -y_i \sum_j w_j x_{ij} \right)$$

$$= \begin{cases} 0 & \text{if } y_i \mathbf{w}^\top \mathbf{x} \geq 0 \\ -y_i x_{ik} & \text{otherwise} \end{cases}$$



# Batch Perceptron Algorithm

Input: training examples  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

$\mathbf{w} = (0, \dots, 0)$  // initial weight vector

Repeat until convergence

$\mathbf{g} = (0, \dots, 0)$  // initial gradient vector

    For  $i = 1$  to  $N$  do

        if  $(y_i \mathbf{w}^\top \mathbf{x}_i < 0)$  //  $\mathbf{x}_i$  is misclassified

            For  $j = 1$  to  $n$  do

$$g_j = g_j - y_i x_{ij}$$

$\mathbf{g} := \mathbf{g} / N$  // average gradient

$\mathbf{w} := \mathbf{w} - \eta \mathbf{g}$  // take a gradient step

# Online Perceptron Algorithm

$\mathbf{w} = (0, \dots, 0)$  be the initial weight vector

$\mathbf{g} = (0, \dots, 0)$  be the initial gradient vector

Repeat forever

    Accept training example  $(\mathbf{x}_i, y_i)$

    if  $(y_i \mathbf{w}^\top \mathbf{x}_i < 0)$  //  $\mathbf{x}_i$  is misclassified

        For  $j = 1$  to  $n$  do  $g_j = -y_i x_{ij}$

$\mathbf{w} := \mathbf{w} - \eta \mathbf{g}$  // take a gradient step

This is called stochastic gradient descent because the overall gradient is approximated by the gradient from each individual example

# Learning Rates and Convergence

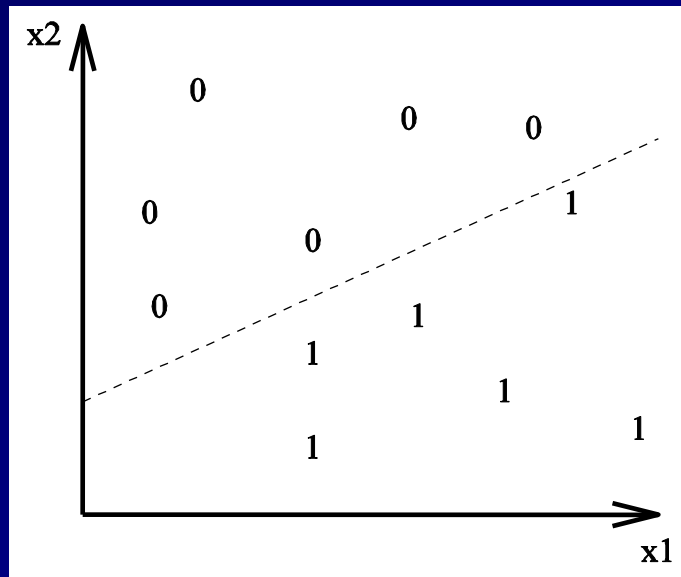
- The learning rate  $\eta$  must decrease to zero in order to guarantee convergence. The online case is known as the Robbins-Munro algorithm. It is guaranteed to converge under the following assumptions:

$$\begin{aligned} \lim_{t \rightarrow \infty} \eta_t &= 0 \\ \sum_{t=0}^{\infty} \eta_t &= \infty \\ \sum_{t=0}^{\infty} \eta_t^2 &< \infty \end{aligned}$$

- The learning rate is also called the step size. Some algorithms (e.g., Newton's method, conjugate gradient) choose the stepsize automatically and converge faster
- There is only one "basin" for linear threshold units, so a local minimum is the global minimum. Choosing a good starting point can make the algorithm converge faster

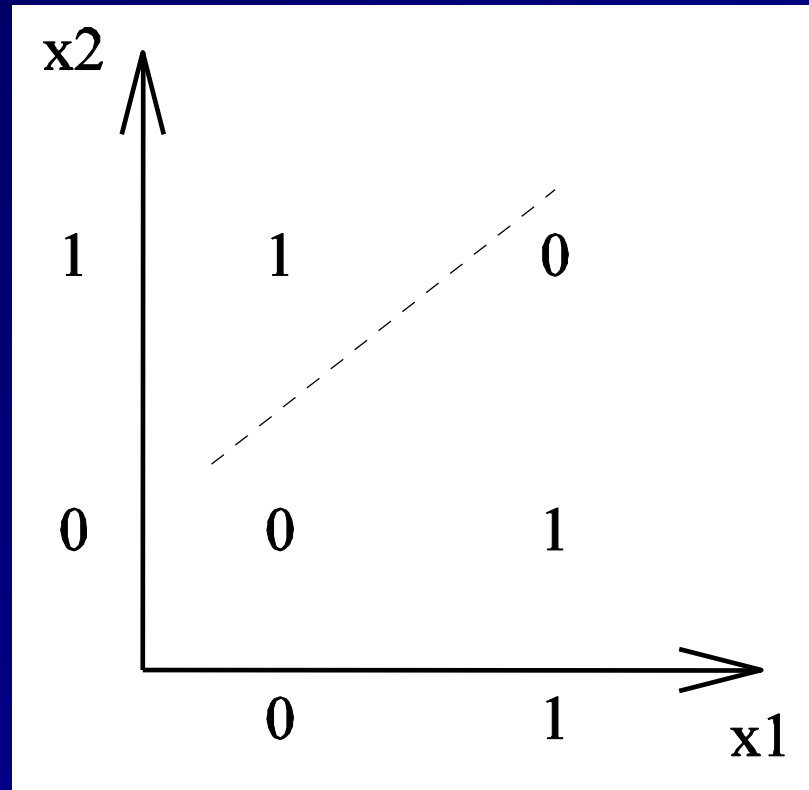
# Decision Boundaries

- A classifier can be viewed as partitioning the input space or feature space  $X$  into decision regions



- A linear threshold unit always produces a linear decision boundary. A set of points that can be separated by a linear decision boundary is said to be linearly separable.

# Exclusive-OR is Not Linearly Separable



# Review

- We adopted the discriminant function approach (no probabilistic model)
- We adopted the hinge loss as a surrogate for the 0/1 loss
- We formulated the optimization problem of minimizing the average hinge loss on the training data
- We solved this problem using gradient descent → Perceptron algorithm

# Logistic Regression

- Learn the conditional probability  $P(y | \mathbf{x})$
- Let  $p_y(\mathbf{x}; \mathbf{w})$  be our estimate of  $P(y | \mathbf{x})$ , where  $\mathbf{w}$  is a vector of adjustable parameters. Assume only two classes  $y = 0$  and  $y = 1$ , and

$$p_1(\mathbf{x}|\mathbf{w}) = \frac{\exp \mathbf{w}^\top \mathbf{x}}{1 + \exp \mathbf{w}^\top \mathbf{x}}$$

$$p_0(\mathbf{x}|\mathbf{w}) = 1 - p_1(\mathbf{x}|\mathbf{w})$$

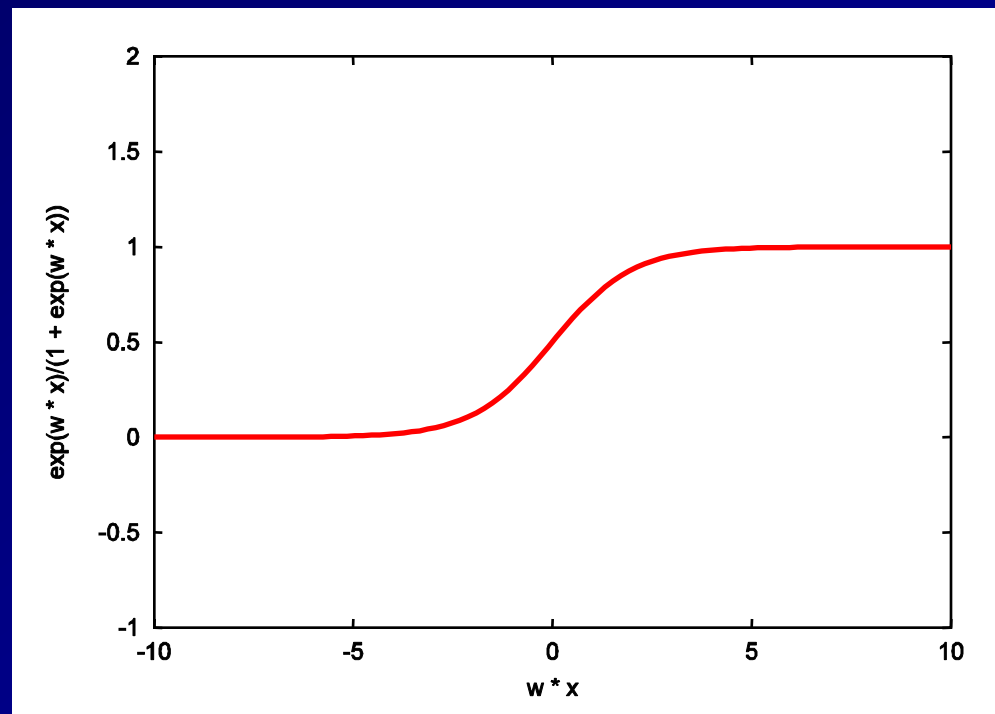
- It is easy to show that this is equivalent to

$$\log \frac{p_1(\mathbf{x}; \mathbf{w})}{p_0(\mathbf{x}; \mathbf{w})} = \mathbf{w}^\top \mathbf{x}$$

- In other words, the log odds of class 1 is a linear function of  $\mathbf{x}$ .

# Why the exp function?

- One reason: A linear function has a range from  $-\infty$  to  $+\infty$  and we need to force it to be positive and sum to 1 in order to be a probability:

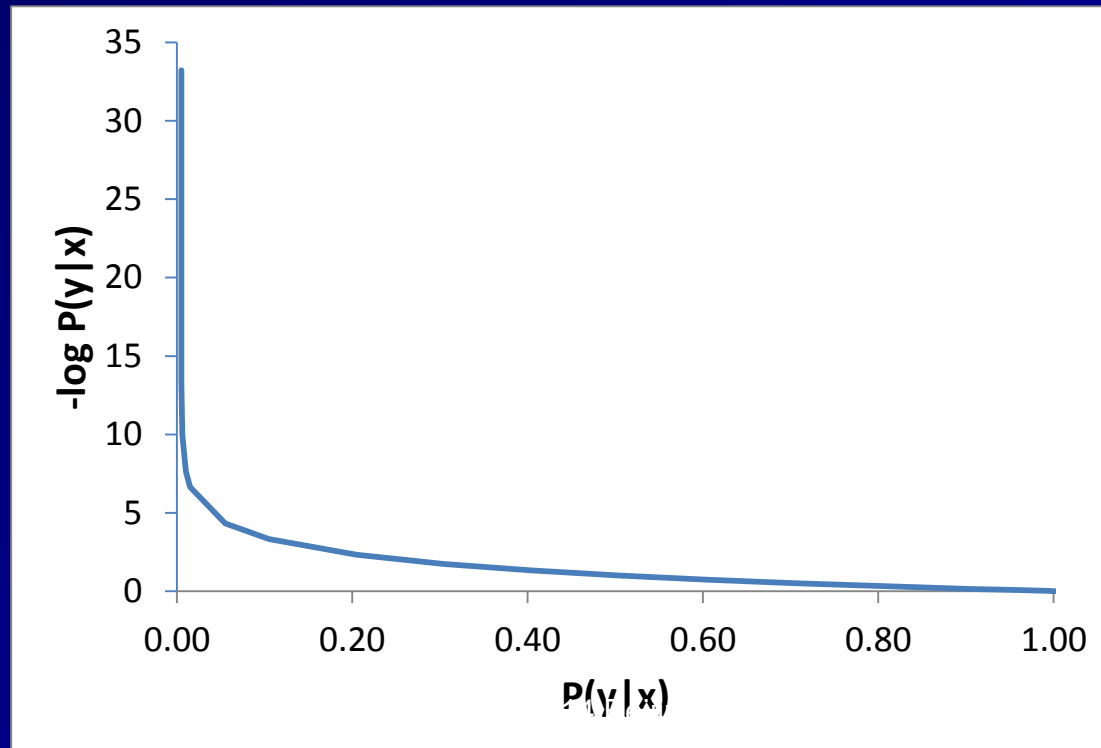




# Deriving a Learning Algorithm: Choosing the Loss Function

- For probabilistic models, we use the log loss:

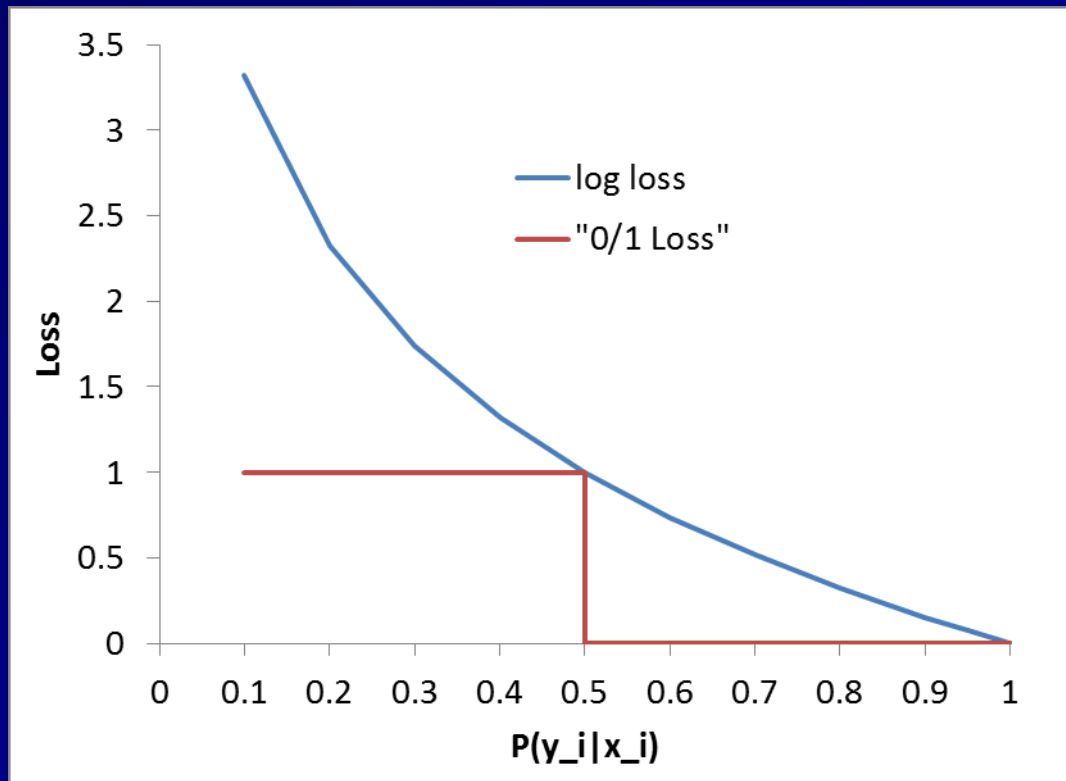
$$\mathcal{L}(\hat{P}(y|x), y) = \begin{cases} -\log \hat{P}(y = 1|x_i) & \text{if } y_i = 1 \\ -\log \hat{P}(y = 0|x_i) & \text{if } y_i = 0 \end{cases}$$



# Comparison with 0/1 Loss

- For probabilistic models, we use the log loss:

$$\mathcal{L}(\hat{P}(y|x), y) = \begin{cases} -\log \hat{P}(y = 1|x_i) & \text{if } y_i = 1 \\ -\log \hat{P}(y = 0|x_i) & \text{if } y_i = 0 \end{cases}$$



# Maximum Likelihood Fitting

■ To minimize the log loss, we should maximize  $\log \hat{P}(y_i|x_i)$

■ The *likelihood* of the data is:

$$\prod_i \hat{P}(y_i|x_i)$$

■ It is easier to work with the *log likelihood*:

$$\sum_i \log \hat{P}(y_i|x_i)$$

# Maximizing the log likelihood via gradient ascent

■ Rewrite the log likelihood in terms of  $p_1(\mathbf{x}_i; \mathbf{w})$ :

$$\ell(\mathbf{w}) = \sum_i y_i \log p_1(\mathbf{x}_i; \mathbf{w}) + (1 - y_i) \log(1 - p_1(\mathbf{x}_i; \mathbf{w}))$$

Differentiate with respect to  $w_j$ :

$$\frac{\partial \ell(y_i; \mathbf{x}_i, \mathbf{w})}{\partial w_j} = \frac{y_i}{p_1(\mathbf{x}_i, \mathbf{w})} \frac{\partial p_1(\mathbf{x}_i, \mathbf{w})}{\partial w_j} + \frac{1 - y_i}{1 - p_1(\mathbf{x}_i, \mathbf{w})} \left( - \frac{\partial p_1(\mathbf{x}_i, \mathbf{w})}{\partial w_j} \right)$$

Gather terms

$$\frac{\partial \ell(y_i; \mathbf{x}_i, \mathbf{w})}{\partial w_j} = \left[ \frac{y_i}{p_1(\mathbf{x}_i, \mathbf{w})} - \frac{1 - y_i}{1 - p_1(\mathbf{x}_i, \mathbf{w})} \right] \left( \frac{\partial p_1(\mathbf{x}_i, \mathbf{w})}{\partial w_j} \right)$$

# Maximizing the log likelihood via gradient ascent (2)

- Collect over common denominator:

$$\frac{\partial \ell(y_i; \mathbf{x}_i, \mathbf{w})}{\partial w_j} = \left[ \frac{y_i(1 - p_1(\mathbf{x}_i, \mathbf{w})) - (1 - y_i)p_1(\mathbf{x}_i, \mathbf{w})}{p_1(\mathbf{x}_i, \mathbf{w})(1 - p_1(\mathbf{x}_i, \mathbf{w}))} \right] \left( \frac{\partial p_1(\mathbf{x}_i, \mathbf{w})}{\partial w_j} \right)$$

- Simplify:

$$\frac{\partial \ell(y_i; \mathbf{x}_i, \mathbf{w})}{\partial w_j} = \left[ \frac{y_i - p_1(\mathbf{x}_i, \mathbf{w})}{p_1(\mathbf{x}_i, \mathbf{w})(1 - p_1(\mathbf{x}_i, \mathbf{w}))} \right] \left( \frac{\partial p_1(\mathbf{x}_i, \mathbf{w})}{\partial w_j} \right)$$

- Now we just need to compute

$$\frac{\partial p_1(\mathbf{x}_i, \mathbf{w})}{\partial w_j}$$

# Computing $\frac{\partial p_1(x_i, w)}{\partial w_j}$

■  $p_1$  can be written as

$$p_1(x_i; w) = \frac{1}{1 + \exp(w^\top x)}$$

■ From this, we obtain:

$$\begin{aligned} \frac{\partial p_1(x_i, w)}{\partial w_j} &= \frac{1}{(1 + \exp(-w^\top x_i))^2} \frac{\partial(1 + \exp(-w^\top x_i))}{\partial w_j} \\ &= -\frac{1}{(1 + \exp(-w^\top x_i))^2} \exp(-w^\top x_i) \frac{\partial}{\partial w_j} (-w^\top x_i) \\ &= -\frac{1}{(1 + \exp(-w^\top x_i))^2} \exp(-w^\top x_i) (-x_{ij}) \\ &= p_1(x_i; w)(1 - p_1(x_i; w))x_{ij} \end{aligned}$$

# Putting it together we have

$$\begin{aligned}\frac{\partial \ell(y_i; \mathbf{x}_i, \mathbf{w})}{\partial w_j} &= \left[ \frac{y_i - p_1(\mathbf{x}_i; \mathbf{w})}{p_1(\mathbf{x}_i; \mathbf{w})(1 - p_1(\mathbf{x}_i; \mathbf{w}))} \right] \left( \frac{\partial p_1(\mathbf{x}_i; \mathbf{w})}{\partial w_j} \right) \\ \frac{\partial p_1(\mathbf{x}_i; \mathbf{w})}{\partial w_j} &= p_1(\mathbf{x}_i; \mathbf{w})(1 - p_1(\mathbf{x}_i; \mathbf{w}))x_{ij} \\ \frac{\partial \ell(y_i; \mathbf{x}_i, \mathbf{w})}{\partial w_j} &= [y_i - p_1(\mathbf{x}_i; \mathbf{w})]x_{ij}\end{aligned}$$

■ The overall gradient is therefore

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i (y_i - p_1(\mathbf{x}_i; \mathbf{w}))x_{ij}$$

■ Note that the first term is the error on the probability scale

# Batch Gradient Ascent for Logistic Regression

Input: training examples  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

$\mathbf{w} = (0, \dots, 0)$  be the initial weight vector

Repeat until convergence

$\mathbf{g} = (0, \dots, 0)$  be the initial gradient vector

For  $i = 1$  to  $N$  do

$$p_i = 1 / (1 + \exp(-\mathbf{w}^T \mathbf{x}_i))$$

$$\text{error}_i := y_i - p_i$$

For  $j = 1$  to  $n$  do

$$g_j = g_j + \text{error}_i \cdot x_{ij}$$

$$\mathbf{g} := \mathbf{g} / N \quad // \text{average gradient}$$

$$\mathbf{w} := \mathbf{w} + \eta \mathbf{g} \quad // \text{take a gradient step}$$

- An online gradient ascent algorithm can be constructed, of course
- Most statistical packages use a second-order (Newton-Raphson) algorithm for faster convergence



# Logistic Regression Implements a Linear Discriminant Function

- In the 2-class 0/1 loss function case, we should predict  $\hat{y} = 1$  if  $\hat{P}(y = 1|\mathbf{x}; \mathbf{w}) > 0.5$

$$\frac{P(y = 1|\mathbf{x}; \mathbf{w})}{P(y = 0|\mathbf{x}; \mathbf{w})} > 1$$

- Take log of both sides

$$\log \frac{P(y = 1|\mathbf{x}; \mathbf{w})}{P(y = 0|\mathbf{x}; \mathbf{w})} > 0$$

- or

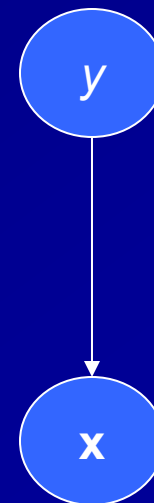
$$\mathbf{w}^T \mathbf{x} > 0$$

# Review

- We adopted the conditional probability approach:  $P(y|\mathbf{x})$
- We adopted the log loss as a surrogate for 0/1 loss
- We formulated the optimization problem of maximizing the average log likelihood on the training data
- We solved this problem using gradient ascent

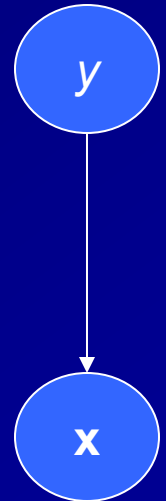
# The Joint Probability Approach: Linear Discriminant Analysis

- Learn  $P(\mathbf{x}, y)$ . This is called the generative approach, because we can think of  $P(\mathbf{x}, y)$  as a model of how the data is generated.
  - For example, if we factor the joint distribution into the form
$$P(\mathbf{x}, y) = P(y) P(\mathbf{x} | y)$$
  - Generative “story”
    - draw  $y \sim P(y)$  choose a class
    - draw  $\mathbf{x} \sim P(\mathbf{x} | y)$  generate the features for  $\mathbf{x}$
  - This can be represented as a probabilistic graphical model



# Linear Discriminant Analysis (2)

- $P(y)$  is a discrete multinomial distribution
  - example:  $P(y = 0) = 0.31$ ,  $P(y = 1) = 0.69$  will generate 31% negative examples and 69% positive examples
- For LDA, we assume that  $P(\mathbf{x} | y)$  is a multivariate normal distribution with mean  $\boldsymbol{\mu}_k$  and covariance matrix  $\Sigma$



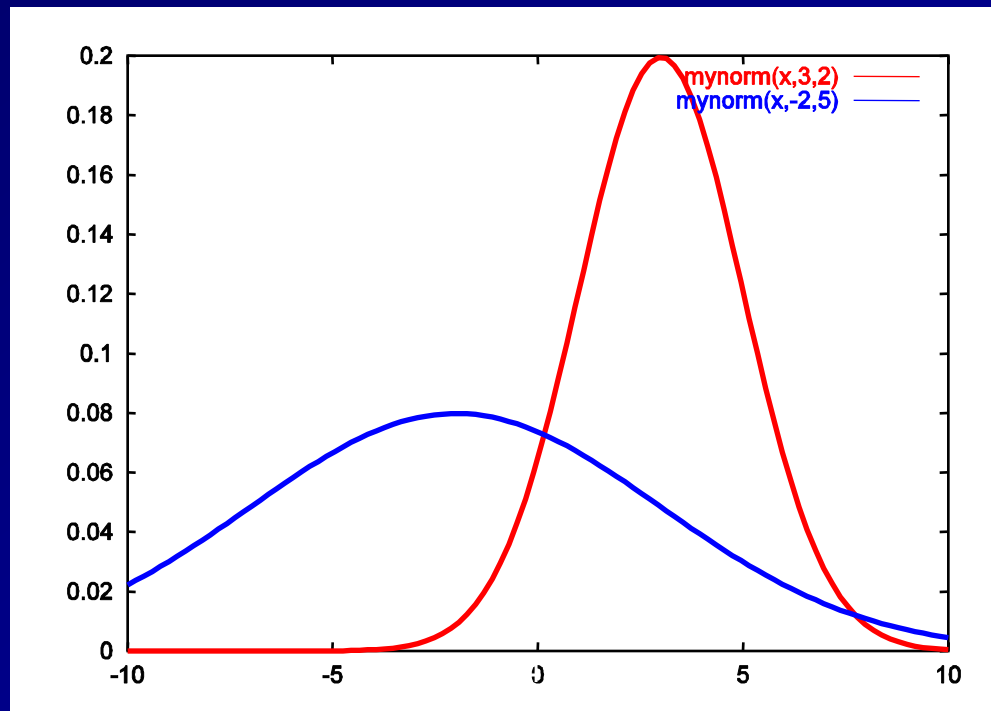
$$P(\mathbf{x}|y = k) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}_k]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_k]\right)$$

# Multivariate Normal Distributions: A tutorial

- Recall that the univariate normal (Gaussian) distribution has the formula

$$p(x) = \frac{1}{(2\pi)^{\frac{1}{2}}\sigma} \exp\left[-\frac{1}{2}\frac{(x - \mu)^2}{\sigma^2}\right]$$

- where  $\mu$  is the mean and  $\sigma^2$  is the variance
- Graphically, it looks like this:



# The Multivariate Gaussian

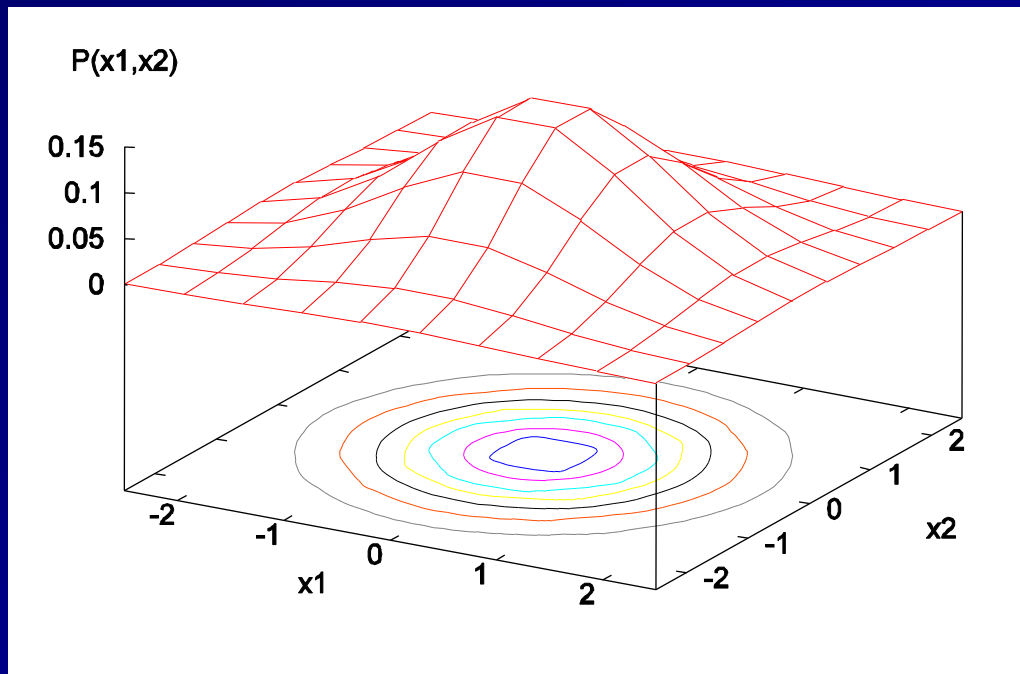
- A 2-dimensional Gaussian is defined by a mean vector  $\mu = (\mu_1, \mu_2)$  and a covariance matrix

$$\Sigma = \begin{bmatrix} \sigma_{1,1}^2 & \sigma_{1,2}^2 \\ \sigma_{2,1}^2 & \sigma_{2,2}^2 \end{bmatrix}$$

- where  $\sigma_{i,j}^2 = E[(x_i - \mu_i)(x_j - \mu_j)]$  is the variance (if  $i = j$ ) or co-variance (if  $i \neq j$ ).  
 $\Sigma$  is symmetric and positive-definite

# The Multivariate Gaussian (2)

- If  $\Sigma$  is the identity matrix  $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  and  $\mu = (0, 0)$ , we get the 2-D standard normal distribution:

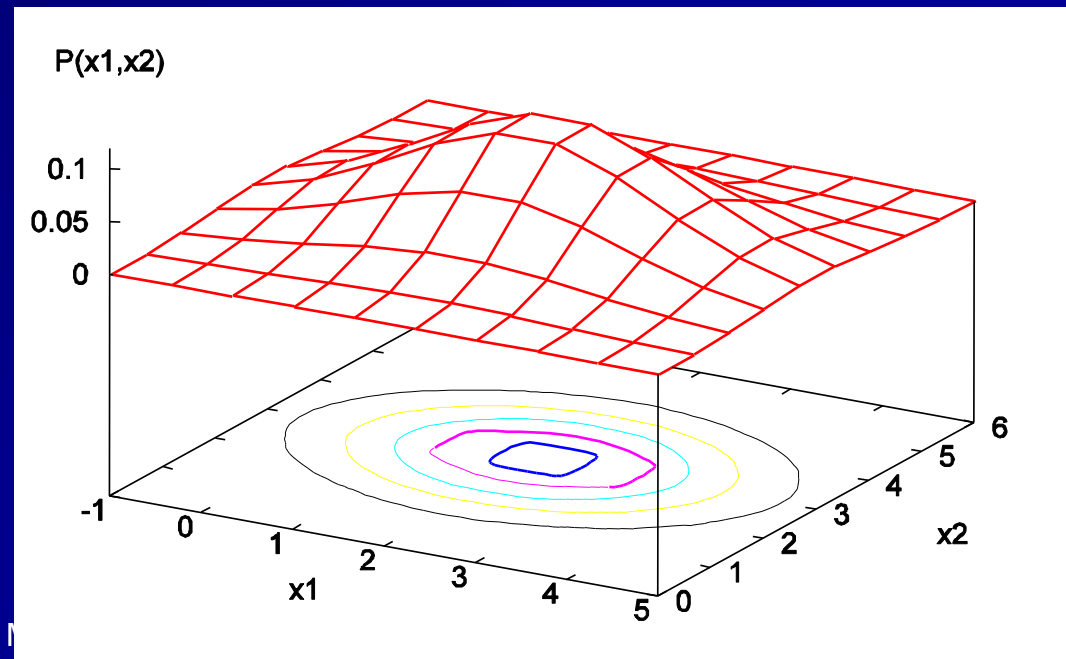


# The Multivariate Gaussian (3)

- If  $\Sigma$  is a diagonal matrix, then  $x_1$ , and  $x_2$  are independent random variables, and lines of equal probability are ellipses parallel to the coordinate axes. For example, when

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \text{ and}$$

$$\mu = (2,3) \text{ we obtain}$$



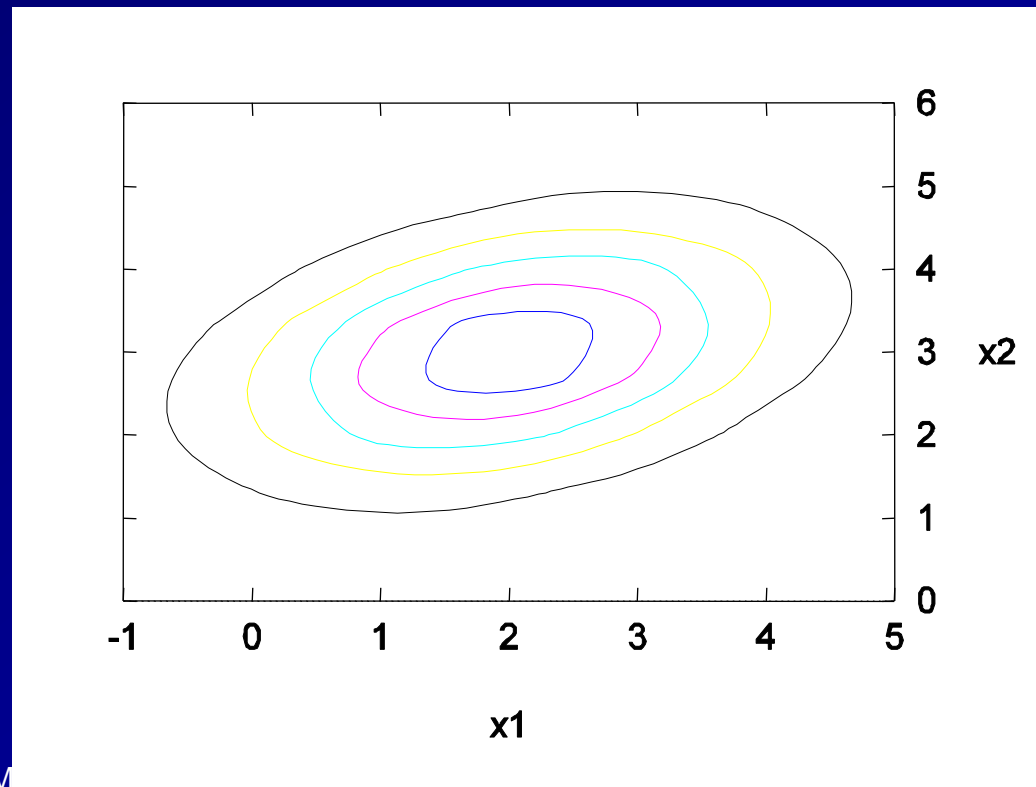


# The Multivariate Gaussian (4)

- Finally, if  $\Sigma$  is an arbitrary matrix, then  $x_1$  and  $x_2$  are dependent, and lines of equal probability are ellipses tilted relative to the coordinate axes. For example, when

$$\Sigma = \begin{bmatrix} 2 & 0.5 \\ 0.5 & 1 \end{bmatrix} \text{ and}$$

$$\mu = (2,3) \text{ we obtain}$$



# Estimating a Multivariate Gaussian

- Given a set of  $N$  data points  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , we can compute the maximum likelihood estimate for the multivariate Gaussian distribution as follows:

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_i \mathbf{x}_i \qquad \hat{\boldsymbol{\Sigma}} = \frac{1}{N} \sum_i (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top$$

- Note that the dot product in the second equation is an outer product. The outer product of two vectors is a matrix:

$$\mathbf{x}\mathbf{y}^\top = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} [y_1 \quad y_2 \quad y_3] = \begin{bmatrix} x_1y_1 & x_1y_2 & x_1y_3 \\ x_2y_1 & x_2y_2 & x_2y_3 \\ x_3y_1 & x_3y_2 & x_3y_3 \end{bmatrix}$$

- For comparison, the usual dot product is written as  $\mathbf{x}^\top \mathbf{y}$

# The LDA Model

- Linear discriminant analysis assumes that the joint distribution has the form

$$P(\mathbf{x}, y) = P(y) \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}_y]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_y]\right)$$

where each  $\boldsymbol{\mu}_y$  is the mean of a multivariate Gaussian for examples belonging to class  $y$  and  $\Sigma$  is a single covariance matrix shared by all classes.

# Fitting the LDA Model

- It is easy to learn the LDA model in a single pass through the data:
  - Let  $\hat{\pi}_k$  be our estimate of  $P(y = k)$
  - Let  $N_k$  be the number of training examples belonging to class  $k$ .

$$\hat{\pi}_k = \frac{N_k}{N}$$

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{N_k} \sum_{\{i:y_i=k\}} \mathbf{x}_i$$

$$\hat{\Sigma} = \frac{1}{N} \sum_i (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{y_i})(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{y_i})^\top$$

- Note that each  $\mathbf{x}_i$  is subtracted from its corresponding  $\boldsymbol{\mu}_{y_i}$  prior to taking the outer product. This gives us the “pooled” estimate of  $\Sigma$
- This is known as the Method of Moments

# LDA learns an LTU

- Just as with Logistic Regression, we should classify  $x$  into class 1 if  $\hat{P}(y = 1|\mathbf{x}) > 0.5$
- Our model contains  $P(y)$  and  $P(\mathbf{x}|y)$ , so we need to perform probabilistic inference to obtain the condition probability:

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} = \frac{\text{Norm}(\mathbf{x}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma})\pi_y}{\sum_k \text{Norm}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma})\pi_k}$$

- As before, we re-express this as

$$\frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} > 1$$

- The denominators cancel, and we have

$$\frac{\text{Norm}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma})\pi_1}{\text{Norm}(\mathbf{x}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma})\pi_0} > 1$$

# LDA Learns an LTU (2)

- Substitute the formula for *Normal*:

$$\frac{\pi_1 \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}_1]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_1]\right)}{\pi_0 \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}_0]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_0]\right)} > 1$$

- Cancel terms, take logs

$$\log \frac{\pi_1}{\pi_0} + [\mathbf{x} - \boldsymbol{\mu}_1]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_1] - [\mathbf{x} - \boldsymbol{\mu}_0]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_0] > 0$$

- With a bit more work

$$\log \frac{\pi_1}{\pi_0} + 2\mathbf{x}^\top \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0^\top \Sigma^{-1} \boldsymbol{\mu}_0 > 0$$

# LDA Learns an LTU (3)

$$\log \frac{\pi_1}{\pi_0} + 2\mathbf{x}^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0^\top \Sigma^{-1} \boldsymbol{\mu}_0 > 0$$

■ Define

$$\mathbf{w} = 2\Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$$
$$c = \log \frac{\pi_1}{\pi_0} + \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0^\top \Sigma^{-1} \boldsymbol{\mu}_0$$

■ Then LDA will classify into  $\hat{y} = 1$  iff

$$\mathbf{x}^\top \mathbf{w} > -c$$

■ which is a linear threshold unit

# Two Geometric Views of LDA

## View 1: Mahalanobis Distance

- The quantity  $D_M(\mathbf{x}, \mathbf{u})^2 = [\mathbf{x} - \mathbf{u}]^\top \Sigma^{-1} [\mathbf{x} - \mathbf{u}]$  is known as the (squared) Mahalanobis distance between  $\mathbf{x}$  and  $\mathbf{u}$ . We can think of the matrix  $\Sigma^{-1}$  as a linear distortion/rotation of the coordinate system that converts the standard Euclidean distance into the Mahalanobis distance

- Note that

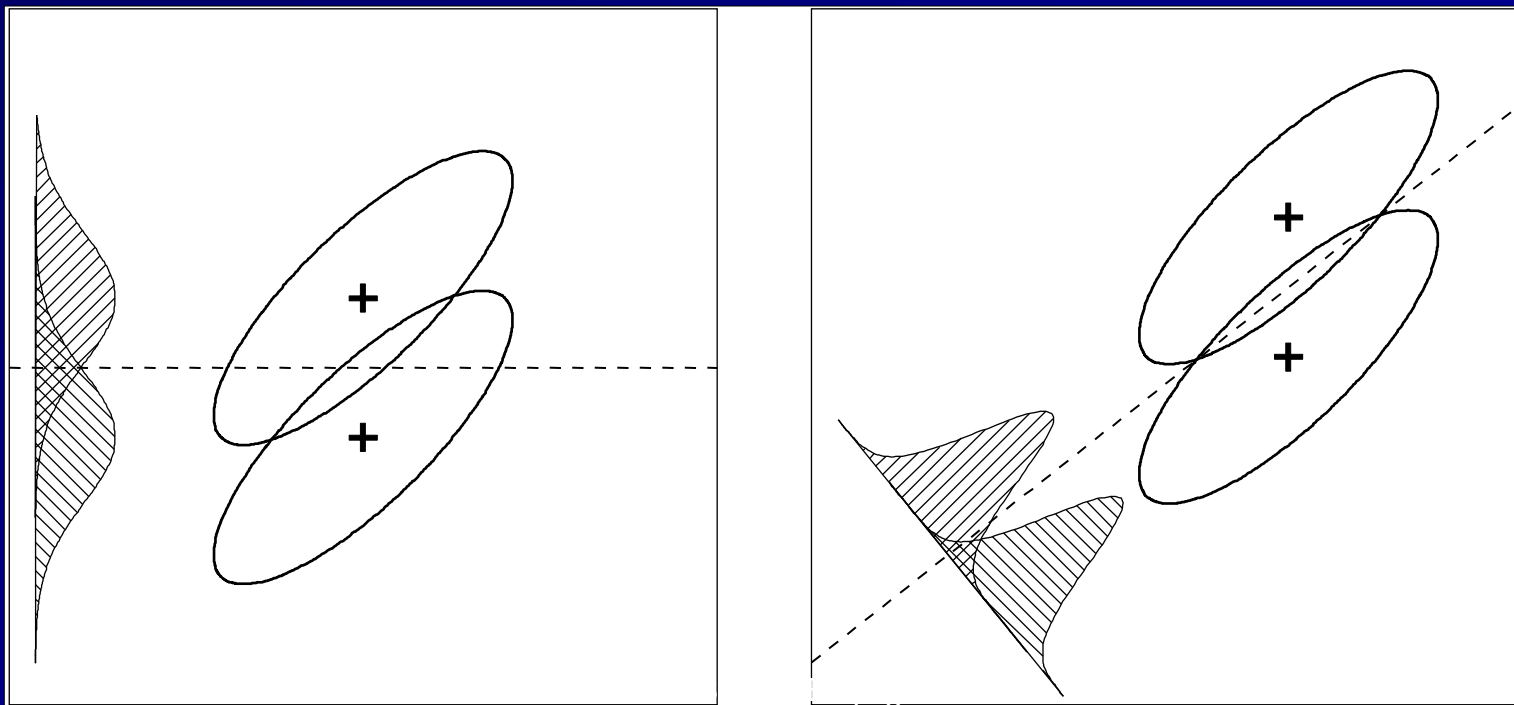
$$\log P(y = k | \mathbf{x}) \propto \log \pi_k - \frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}_k]^\top \Sigma^{-1} [\mathbf{x} - \boldsymbol{\mu}_k]$$
$$\log P(y = k | \mathbf{x}) \propto \log \pi_k - \frac{1}{2} D_M(\mathbf{x}, \boldsymbol{\mu}_k)^2$$

- Therefore, we can view LDA as computing  $D_M(\mathbf{x}, \boldsymbol{\mu}_0)^2$  and  $D_M(\mathbf{x}, \boldsymbol{\mu}_1)^2$  and then classifying  $\mathbf{x}$  according to which mean  $\boldsymbol{\mu}_0$  or  $\boldsymbol{\mu}_1$  is closer in Mahalanobis distance (corrected by  $\log \pi_k$ )



# View 2: Most Informative Low-Dimensional Projection

- LDA can also be viewed as finding a hyperplane of dimension  $K - 1$  such that  $x$  and the  $\{\mu_k\}$  are projected down into this hyperplane and then  $x$  is classified to the nearest  $\mu_k$  using Euclidean distance inside this hyperplane



# Generalizations of LDA

## ■ General Gaussian Classifier

- Instead of assuming that all classes share the same  $\Sigma$ , we can allow each class  $k$  to have its own  $\Sigma_k$ . In this case, the resulting classifier will be a quadratic threshold unit (instead of an LTU)

## ■ Naïve Gaussian Classifier

- Allow each class to have its own  $\Sigma_k$ , but require that each  $\Sigma_k$  be diagonal. This means that *within* each class, any pair of features  $x_{j_1}$  and  $x_{j_2}$  will be assumed to be statistically independent. The resulting classifier is still a quadratic threshold unit (but with a restricted form)

# Review

## Linear Discriminant Analysis

- We adopted the joint probability approach:  $P(x, y)$
- We adopted the log loss as a surrogate for 0/1 loss
- We fit the model directly, via the method of moments

# Comparing Perceptron, Logistic Regression, and LDA

- How should we choose among these three algorithms?
- There are several trade-offs
- There is a big debate in the machine learning community!

# Issues in the Debate

- Statistical Efficiency. If the generative model  $P(x, y)$  is correct, then LDA usually gives the highest accuracy, particularly when the amount of training data is small. If the model is correct, LDA requires 30% less data than Logistic Regression in theory
- Computational Efficiency. Generative models typically are the easiest to learn. In our example, the LDA parameters can be computed directly from the data without using gradient descent.

# Issues in the Debate

- Robustness to changing loss functions. Both generative and conditional probability models allow the loss function to be changed at run time without re-learning. Perceptron requires re-training the classifier when the loss function changes.
  - Probabilistic modelling separates model learning from making predictions or decisions
  - Suppose the cost of a false positive SPAM prediction is 10 whereas a false negative is 1
  - The cost of classifying  $x$  as  $y = 1$  is  $P(y = 0|x) \times 10$
  - The cost of classifying  $x$  as  $y = 0$  is  $P(y = 1|x) \times 1$
  - So we can choose the value of  $y$  that minimizes the expected cost

# Issues in the Debate

## ■ Vapnik's Principle

- If your goal is to minimize 0/1 loss, then you should do that directly, rather than first solving a harder problem (probability estimation)
- This is what Perceptron does
- Other algorithms that follow this principle
  - Support Vector Machines
  - Decision Trees
  - Neural Networks

# Issues in the Debate

- Robustness to model assumptions. The generative model usually performs poorly when the assumptions are violated. For example, if  $P(\mathbf{x}|y)$  is very non-Gaussian, then LDA won't work well. Logistic Regression is more robust to model assumptions, and Perceptron is even more robust.

Consequently, making the generative approach work often requires more detailed modeling of  $P(\mathbf{x}|y)$ .

- Robustness to missing values and noise. In many applications, some of the features  $x_{ij}$  may be missing or corrupted in some of the training examples. Generative models typically provide better ways of handling this than non-generative models.



# Questions?

# Some Questions to Think About

- Machine learning has many powerful non-linear classifiers, why didn't you discuss those?
  - Linear methods are surprisingly powerful, especially in computer vision and natural language processing where the number of features is very large
- You have showed how to fit these models to training data, but that doesn't guarantee that they will make good predictions on new data points
  - Excellent question! That is the subject of Part 2

# Break

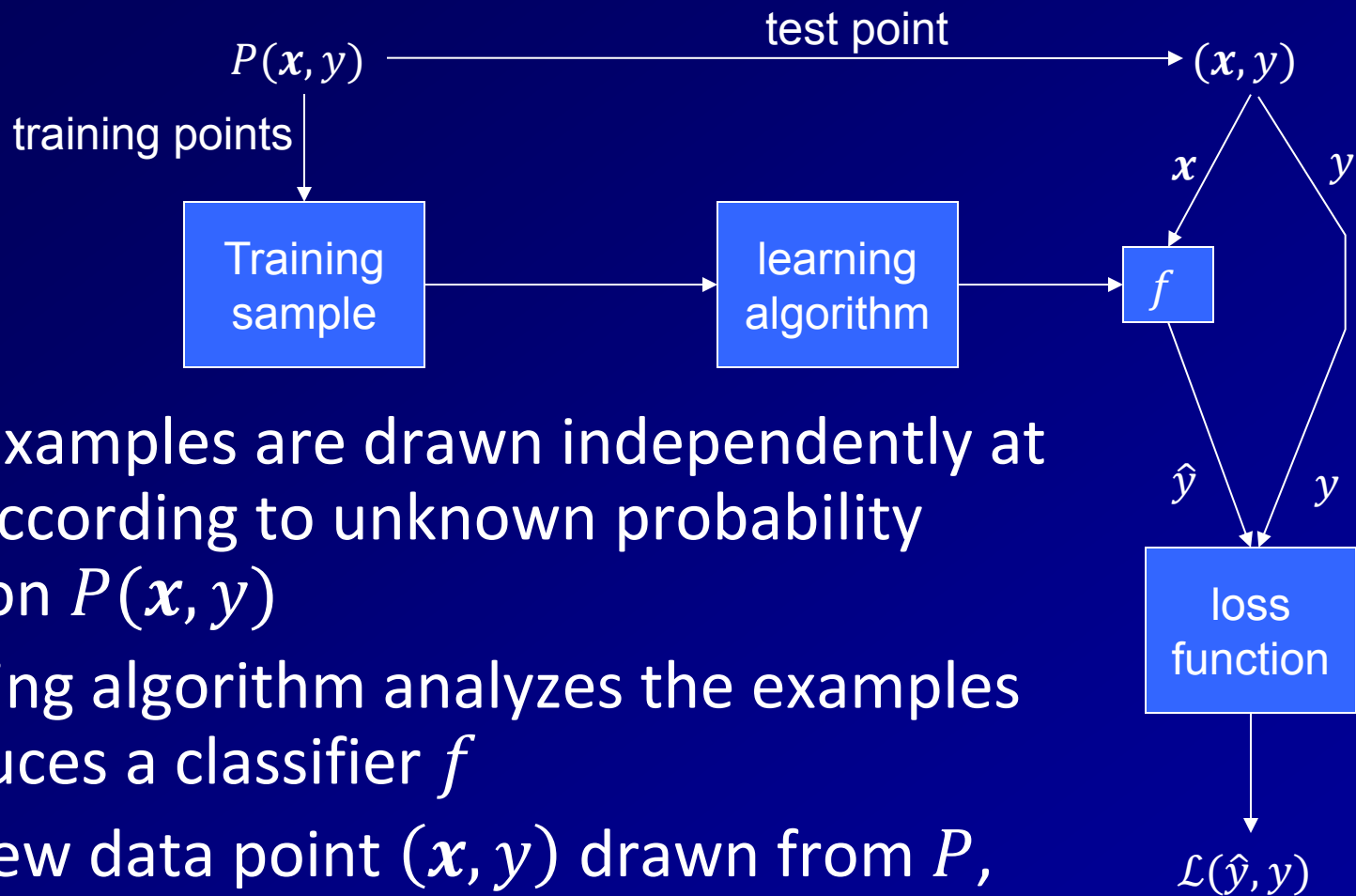
# Introduction to Machine Learning Part 2

Thomas G. Dietterich  
tgd@eecs.oregonstate.edu

# Outline

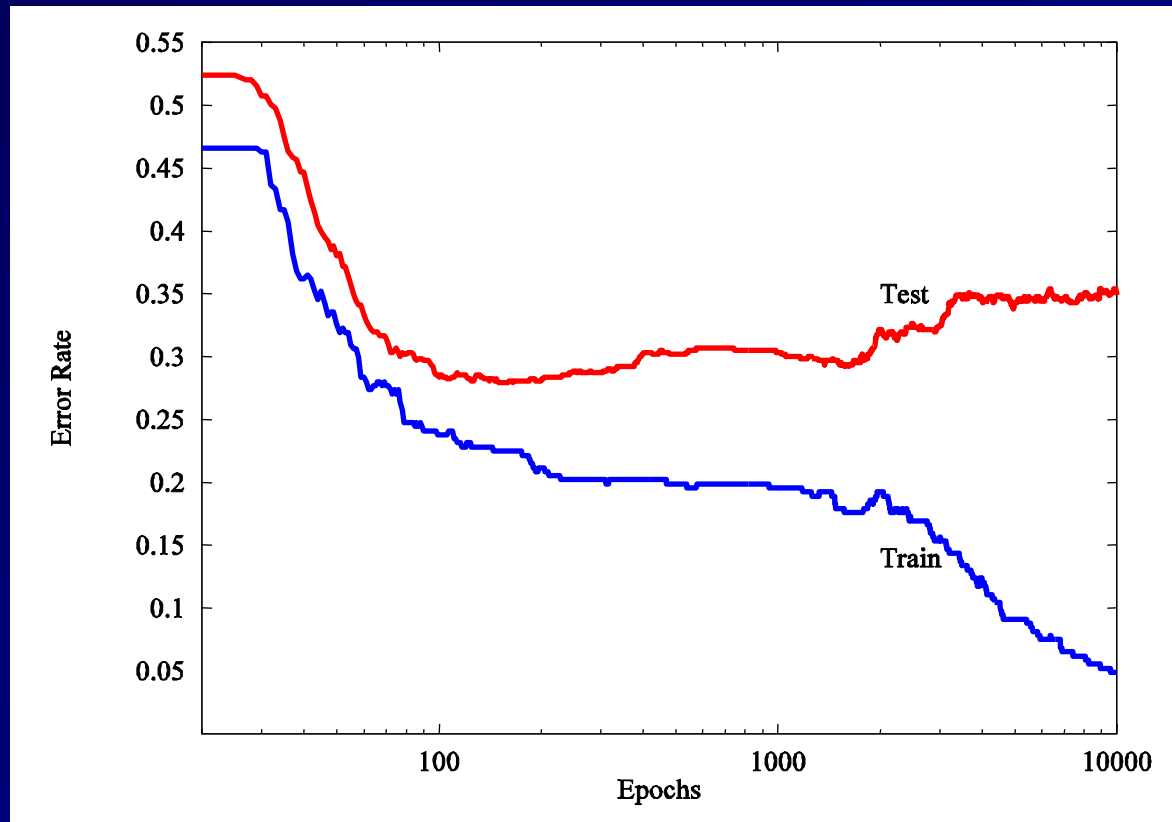
- What is Machine Learning?
- Introduction to Supervised Learning: Linear Methods
- Overfitting, Regularization, and the Bias-Variance Tradeoff
- Review and Summary

# Formal Setting



- Training examples are drawn independently at random according to unknown probability distribution  $P(x, y)$
- The learning algorithm analyzes the examples and produces a classifier  $f$
- Given a new data point  $(x, y)$  drawn from  $P$ , the classifier is given  $x$  and predicts  $\hat{y} = f(x)$
- The loss  $\mathcal{L}(\hat{y}, y)$  is then measured
- Goal of the learning algorithm: Find the  $f$  that minimizes the *expected loss* on new points

# The Problem of Overfitting



- Model: Neural Network
- Epoch: One batch gradient descent step
- After about 180 epochs, error on the test data starts to increase even though the model continues to become more accurate on the training data

Lesson:

It is not enough to minimize the  
loss on the training data  
if our goal is to optimize accuracy  
on *new* data points



# Another Example

- True function:

$$y = f(x_1) = -(x_1 - 6)^2 + 5$$

- Training data:

$$x_1 \sim \text{unif}(1,6)$$

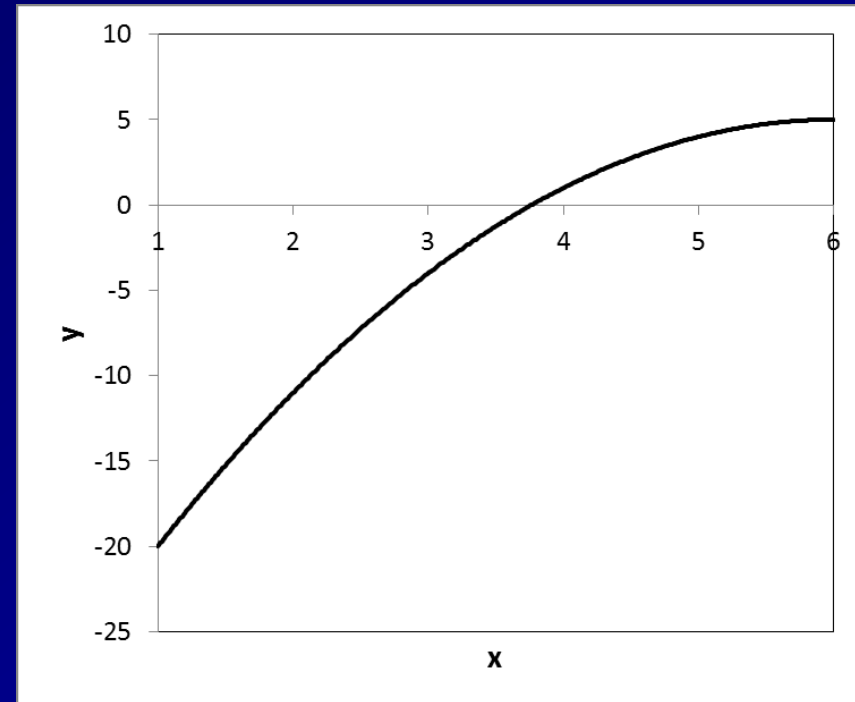
$$x_2 \sim \text{unif}(-6,6)$$

$$y = f(x_1) + \text{norm}(0,10)$$

- Two models:

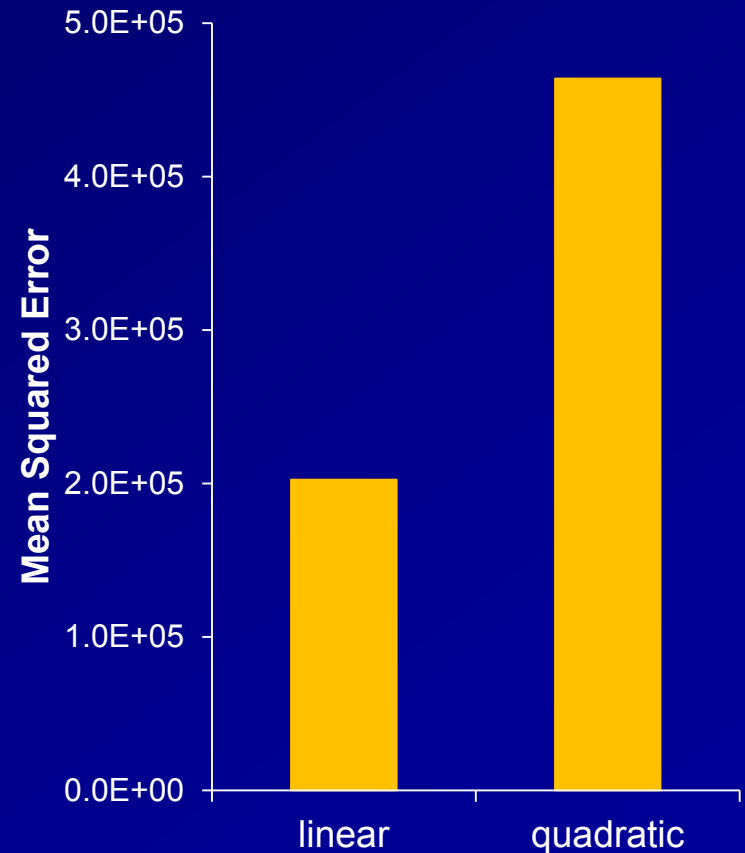
$$\text{Model 1: } y = w_0 + w_1x_1 + w_2x_2$$

$$\text{Model 2: } y = w_0 + w_1x_1 + w_{12}x_1^2 + w_2x_2 + w_{22}x_2^2$$



# Small, Noisy Training Set

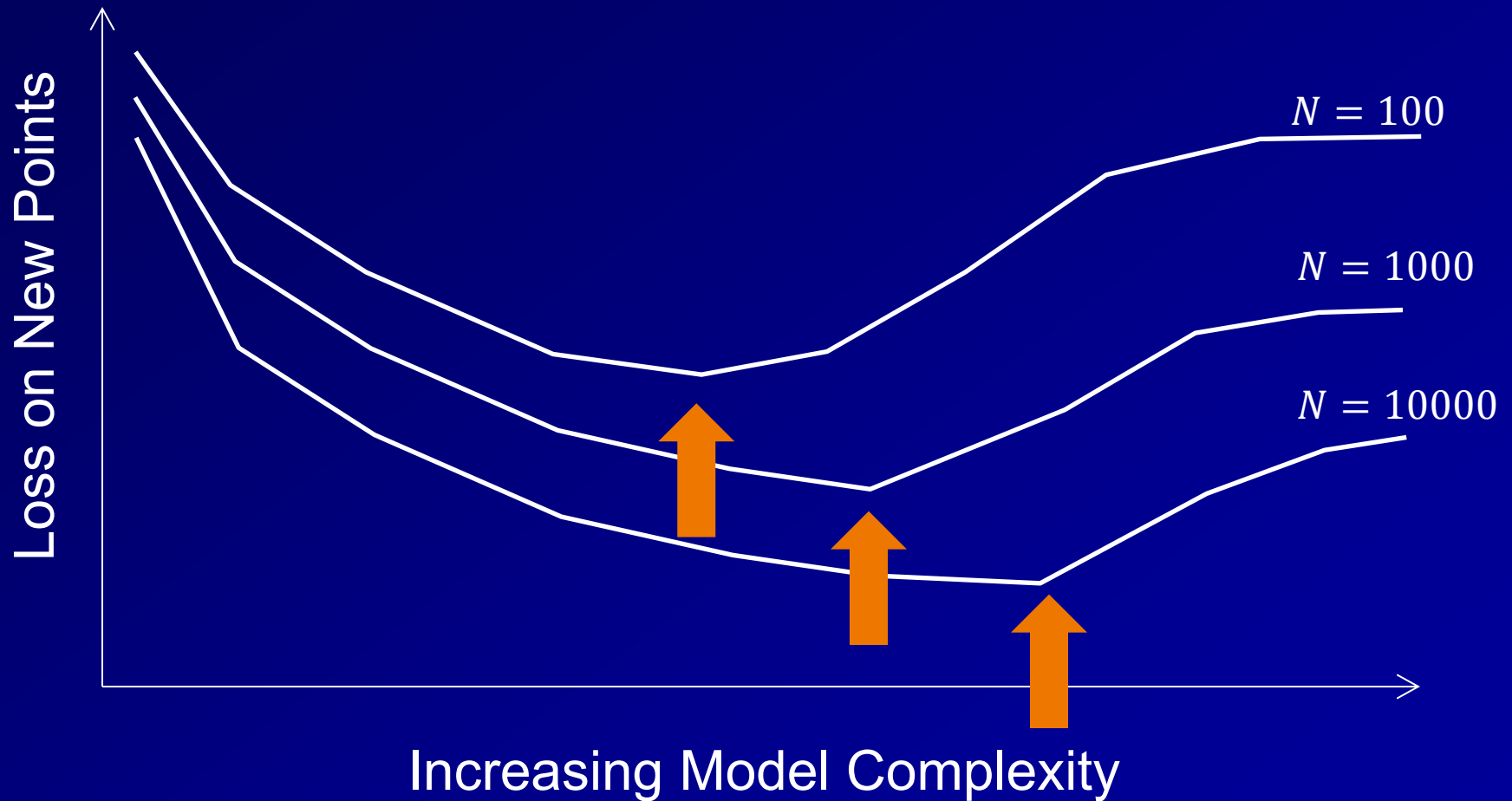
- 3 training examples
- Linear model is much more accurate even though
  - It cannot represent the true model
  - The quadratic model can represent the true model



# The Three-Way Tradeoff

- There is a tradeoff between
  - amount of data
  - complexity of the model fit to the data
  - accuracy of the model on new data points

# Three-way Tradeoff



# What is Model Complexity?

- Parametric models: Complexity = number of weights
  - Linear model is less complex than a quadratic model
  - Extra features → extra complexity
- More subtle:
  - Weights that are zero don't contribute to complexity
  - Small weights contribute less complexity than large weights

# Controlling Model Complexity Via Regularization

- Regularization: Penalize the magnitude of the weights in the model
- Example: “square” penalty

$$J(\mathbf{w}) = \frac{1}{N} \sum_i \mathcal{L}(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_j w_j^2$$

The severity of the penalty is controlled by  $\lambda > 0$

Adds a term of  $2\lambda \sum_j w_j$  to the gradient descent

In neural networks, this is called “weight decay”

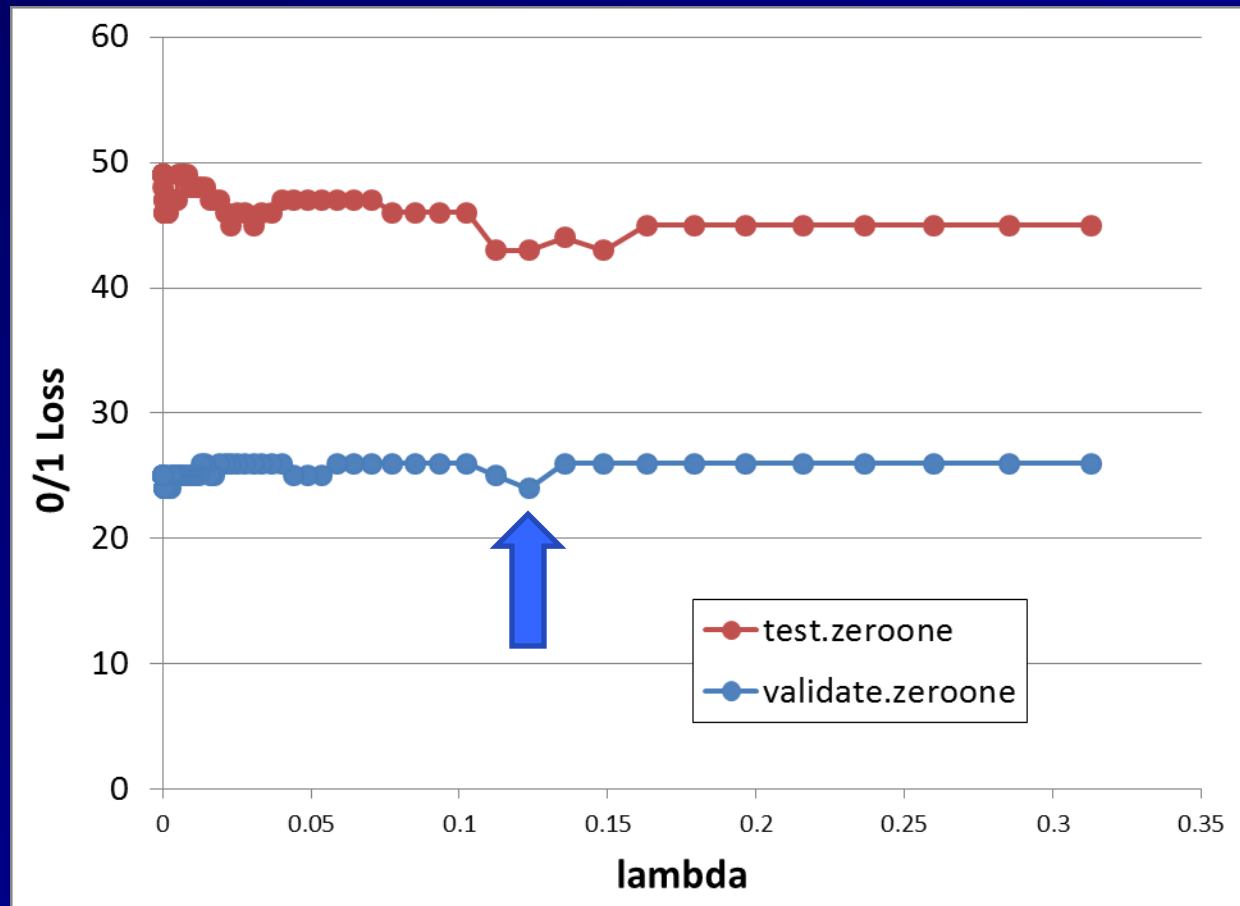
# How to choose $\lambda$ ?

- Simple Holdout Method
- Subdivide training data  $S$  into two subsets:  $S_{train}$  and  $S_{holdout}$ .
- Choose a set of candidate  $\lambda$  values =  $\lambda_0, \lambda_1, \dots, \lambda_R$
- Minimize the penalized loss on  $S_{train}$ , measure the prediction loss on  $S_{holdout}$
- Choose the  $\lambda$  value that gives the smallest prediction loss

# Simple Holdout Example

Logistic regression  
with regularization  
penalty

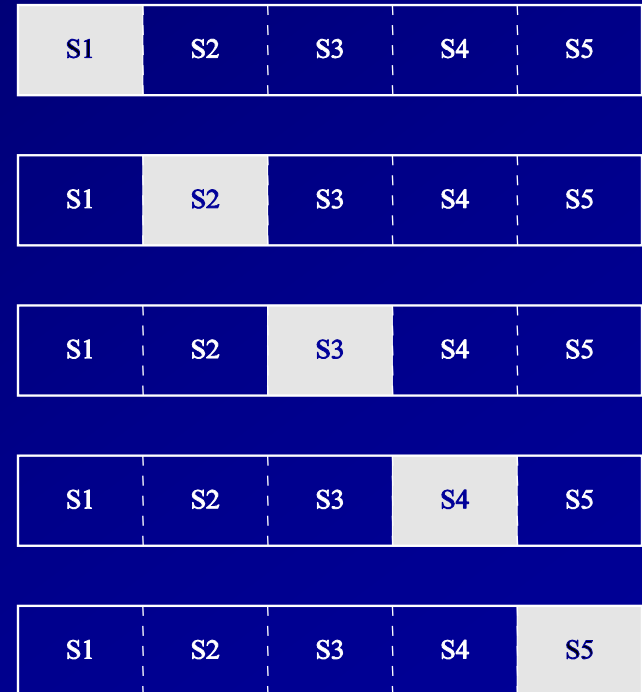
Overfitting was  
not too severe in  
this problem





# $k$ -fold Cross-Validation to determine $\lambda$

- To evaluate a value for  $\lambda$
- Randomly divide  $S$  into  $k$  equal-sized subsets
- Run learning algorithm  $k$  times, each time use one subset for  $S_{eval}$  and the rest for  $S_{train}$
- Compute the average loss on  $S_{eval}$
- Choose the  $\lambda$  value with minimum loss



# A Bayesian Perspective

- A fully-generative Bayesian story:

- Generate the weights:  $\mathbf{w} \sim P(\mathbf{w})$

- For each data point  $i$

  - Generate the class label:  $y_i \sim P(y_i)$

  - Generate the features:  $\mathbf{x}_i \sim P(\mathbf{x}_i | y_i, \mathbf{w})$

  - Assemble the data set:  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

# Inferring the most likely $\mathbf{w}$

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} P(\mathbf{w}|S)$$

$$P(\mathbf{w}|S) = P(\mathbf{w}) \prod_i P(y_i) P(\mathbf{x}_i|y_i, \mathbf{w})$$

$$\log P(\mathbf{w}|S) = \log P(\mathbf{w}) + \sum_i \log P(y_i) + \log P(\mathbf{x}_i|y_i, \mathbf{w})$$

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \ell(S, \mathbf{w}) + \log P(\mathbf{w})$$

The regularization penalty is the same as the log prior on the weights

This provides a way of incorporating different penalties on different weights in a model, based on prior knowledge

# Controlling Complexity by Early Stopping

## ■ Incrementally add complexity

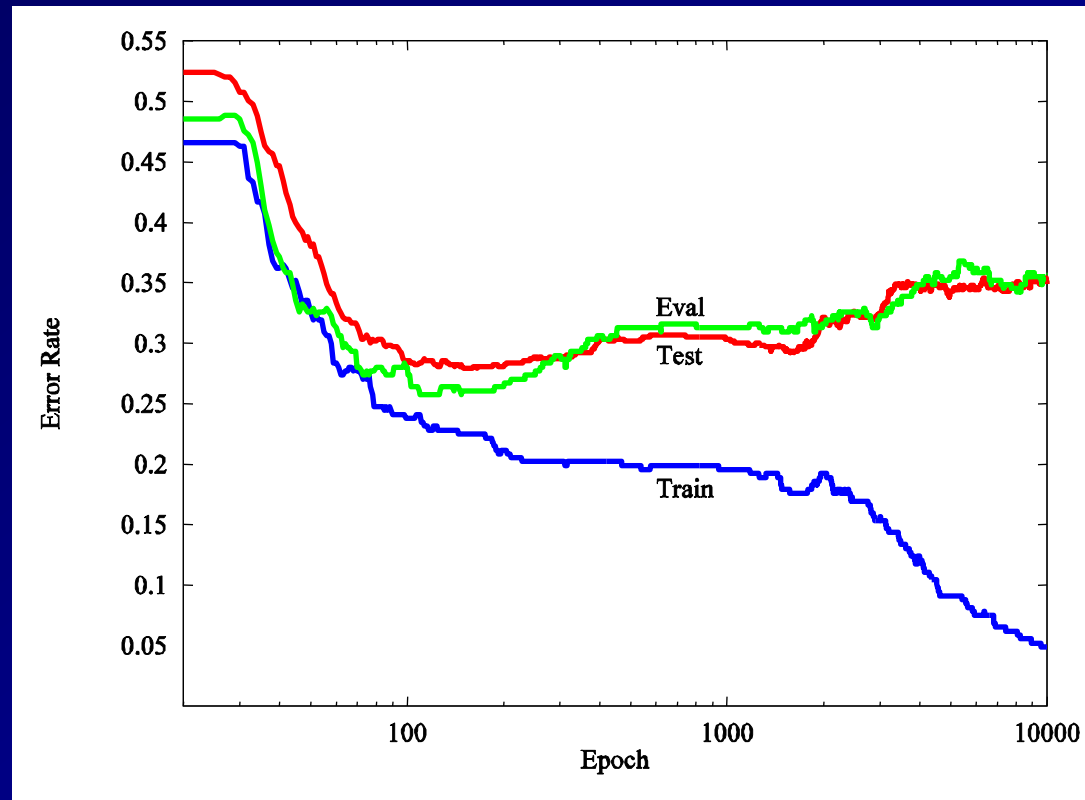
- monitor performance on  $S_{eval}$
- stop when performance drops

## ■ Examples:

- incrementally adding variables in regression
- growing a decision tree
- early stopping in stochastic gradient descent

# Early Stopping

- Subdivide  $S$  into  $S_{train}$  and  $S_{eval}$
- Initialize weights to 0
- Perform gradient descent on  $S_{train}$ , which gradually causes the weights to grow (in magnitude)
- Measure the loss on  $S_{eval}$  after each gradient step
- Stop when the loss starts to increase



# Other Methods for Controlling Complexity

- Model pruning
  - over-fit a model, then prune/shrink
- Adding noise
  - to the inputs or to intermediate quantities
  - makes it harder to overfit
- Ensembles (see below)

# Bias-Variance Analysis

- An alternative view of over-fitting and complexity control
- Bias: Systematic error in the model
  - typically caused by an inability to express the full complexity of the data
- Variance: Variability in the fitted model
  - typically caused by having a model that is too complex for the amount of data
- The total error of the model can be partitioned into the sum of a bias term and a variance term
- As we increase  $\lambda$  we increase bias but reduce variance

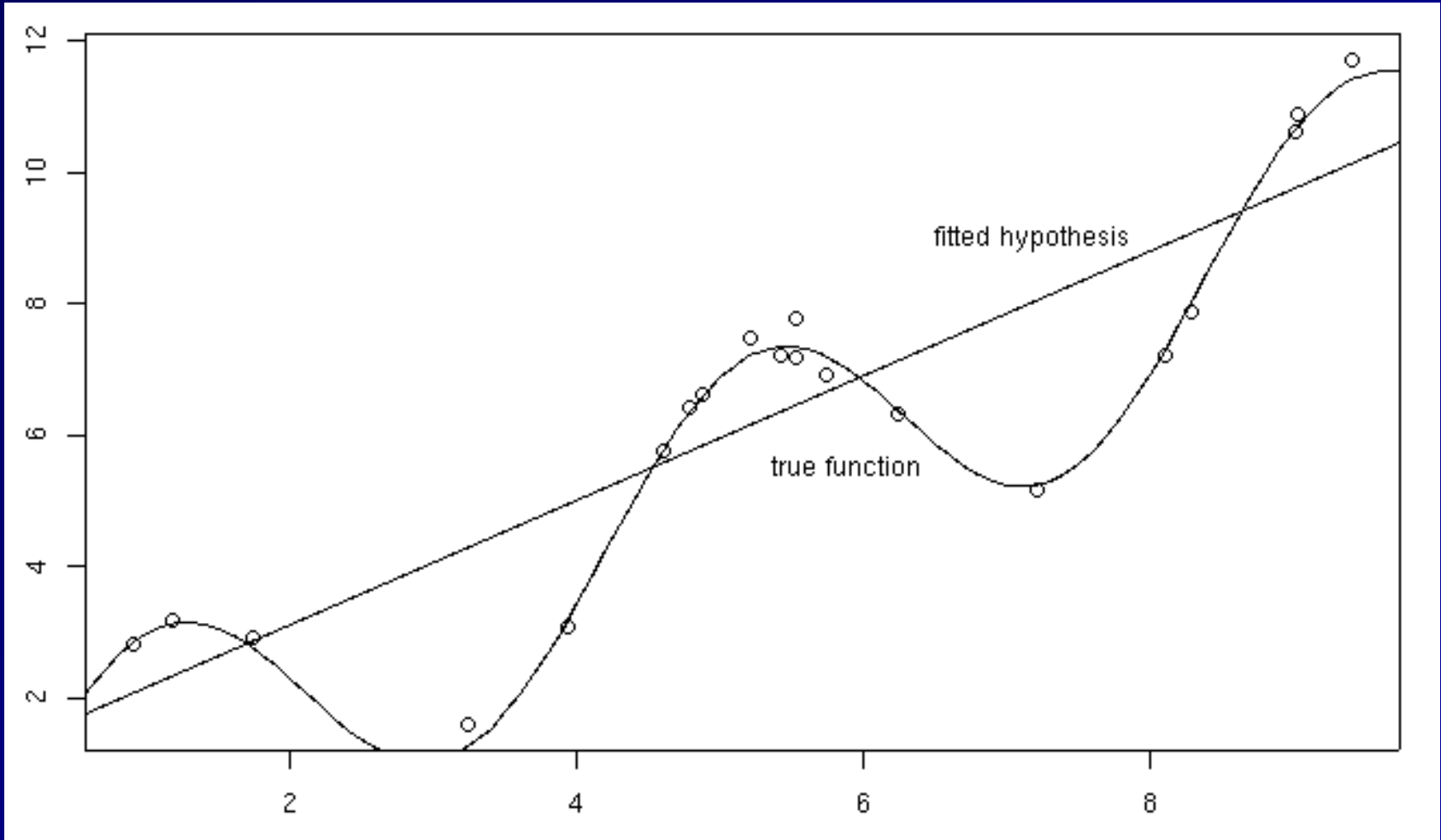
# Bias-Variance Analysis of Regression

- True function:  $y = f(\mathbf{x}) + \epsilon$ 
  - where  $\epsilon \sim \text{Norm}(0, \sigma^2)$
- We are given a set  $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$
- We fit a model  $h(\mathbf{x})$  to minimize the square loss

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2$$
$$J(h) = \sum_i (\hat{y}_i - y_i)^2$$



Example:  $N = 20$  points  
 $y = x + 2 \sin(1.5x) + \text{Norm}(0,0.2)$



# Bias-Variance Analysis

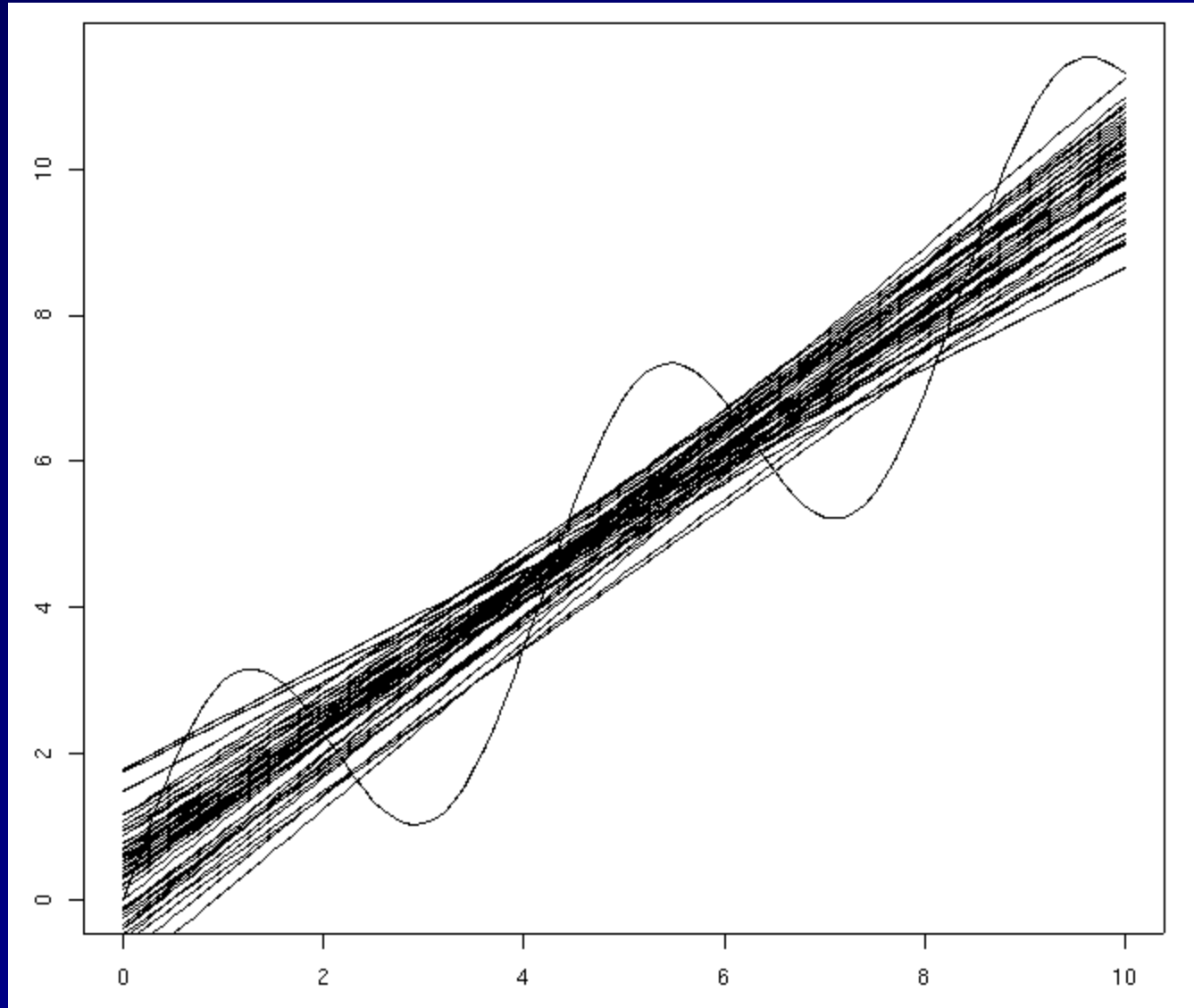
- Given a new data point  $(x^*, y^*)$  (with predicted value  $h(x^*)$ ), we would like to decompose our error

$$(y^* - h(x^*))^2$$

# Classical Statistical Analysis

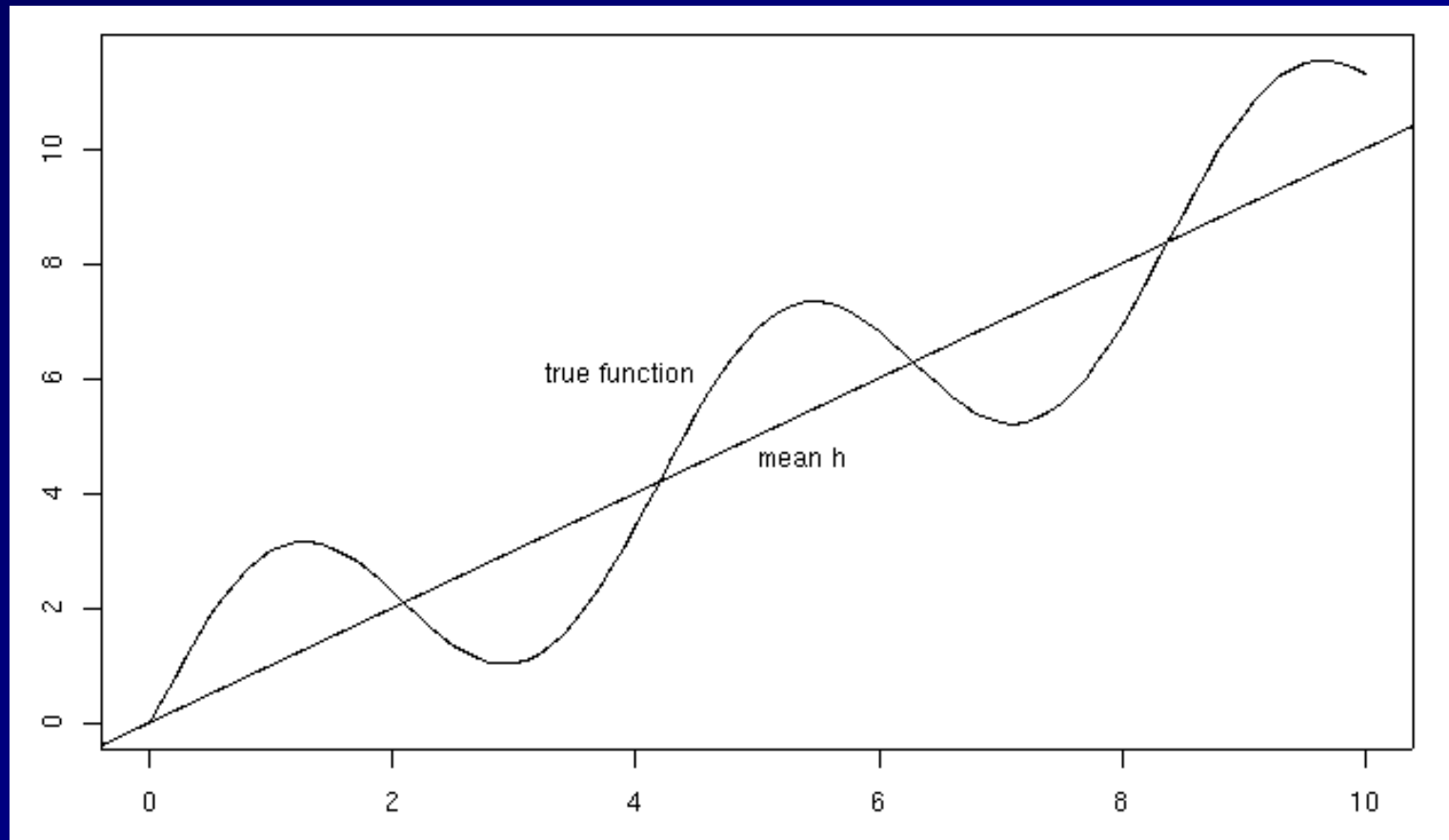
- Imagine that our particular training sample  $S$  is drawn from some population of possible training samples according to  $P(S)$ .
- We fit  $h$  to  $S$  ( $h$  is a random quantity)
- Compute  $\mathbb{E}_P[(y^* - h(x^*))^2]$
- Decompose this into “bias”, “variance”, and “noise”

# 50 fits (20 examples each)



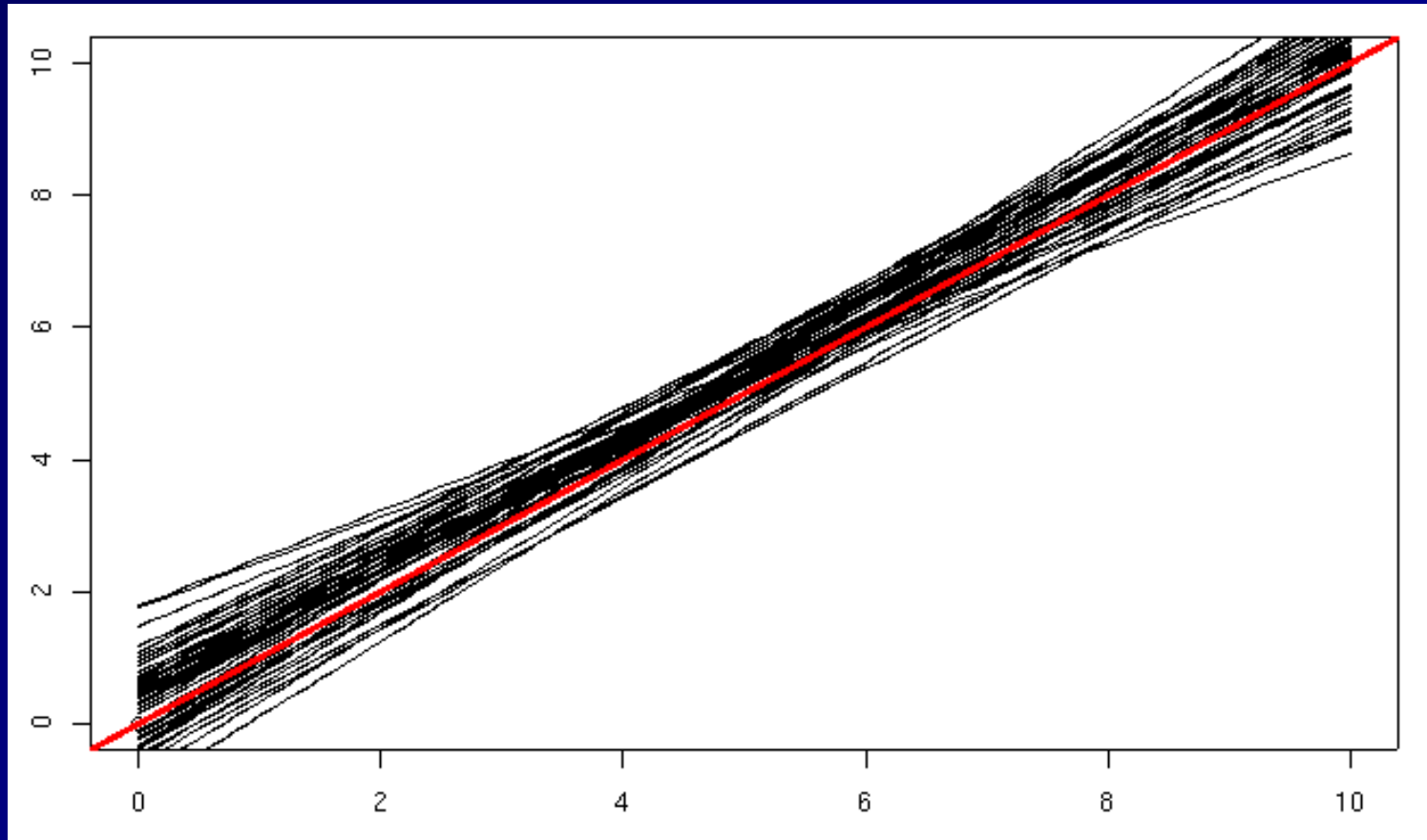
# Bias

- The difference between the average predicted value  $\mathbb{E}_P[h(x^*)]$  and the value of the true function  $f(x^*)$ :  $\mathbb{E}_P[h(x^*)] - f(x^*)$



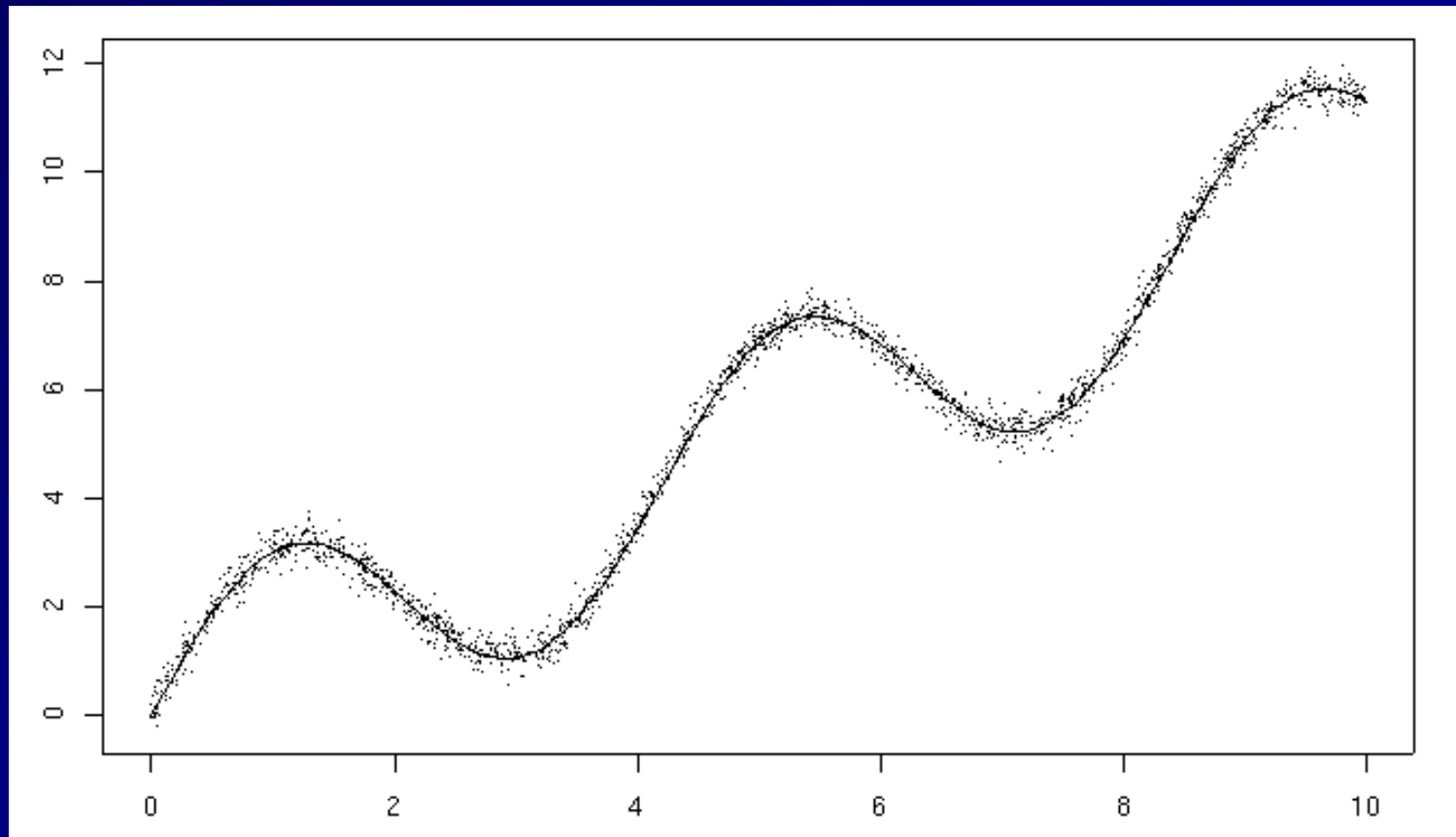
# Variance

- The variance of  $h(x^*)$  around its average value  $\mathbb{E}_P[h(x^*)]$ :  $\mathbb{E}_P[(h(x^*) - \mathbb{E}_P[h(x^*)])^2]$



# Noise

- The variation of  $y^*$  around its true average value  $f(x^*)$ :  $(y^* - f(x^*))^2$

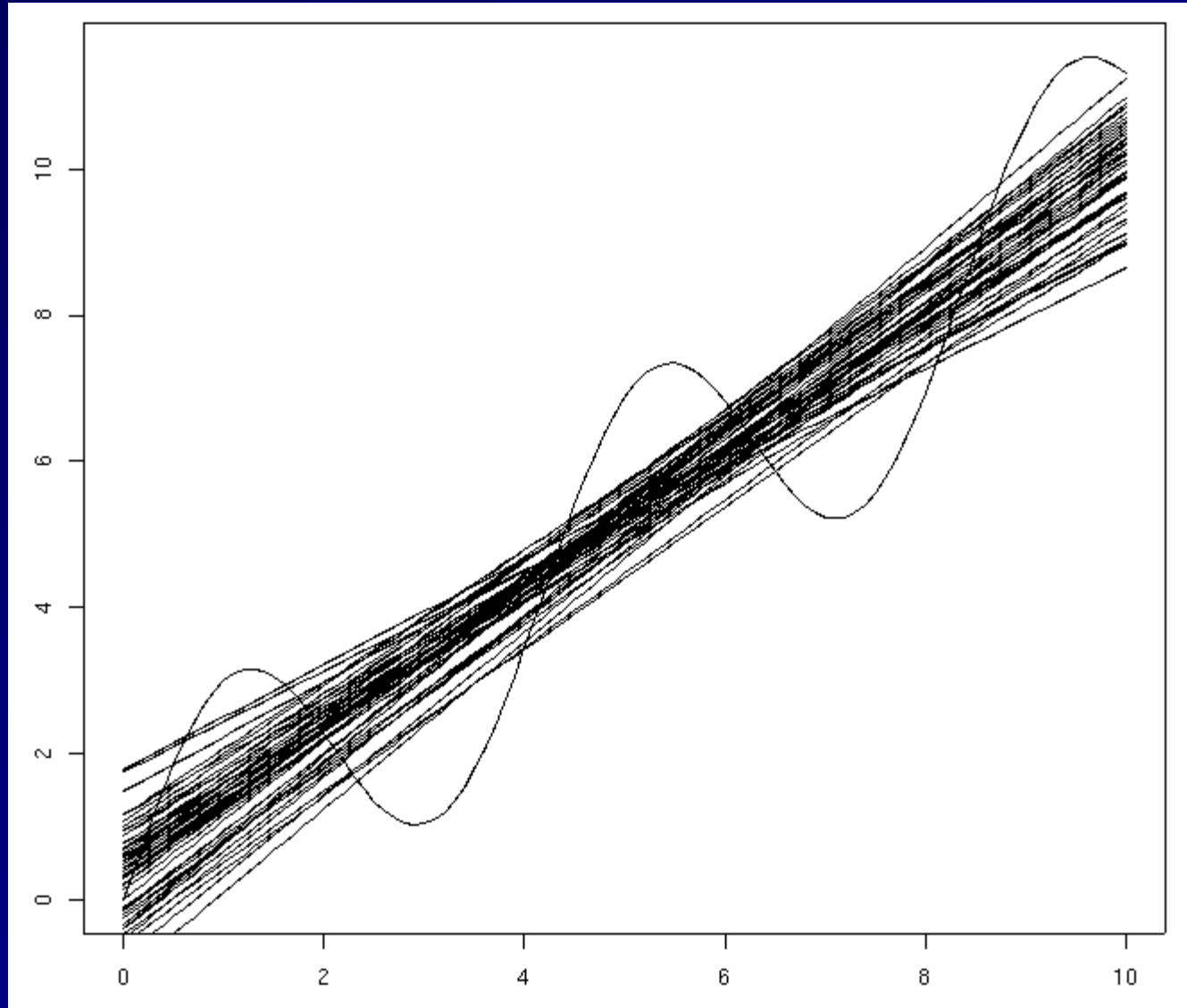


# The Bias-Variance Decomposition

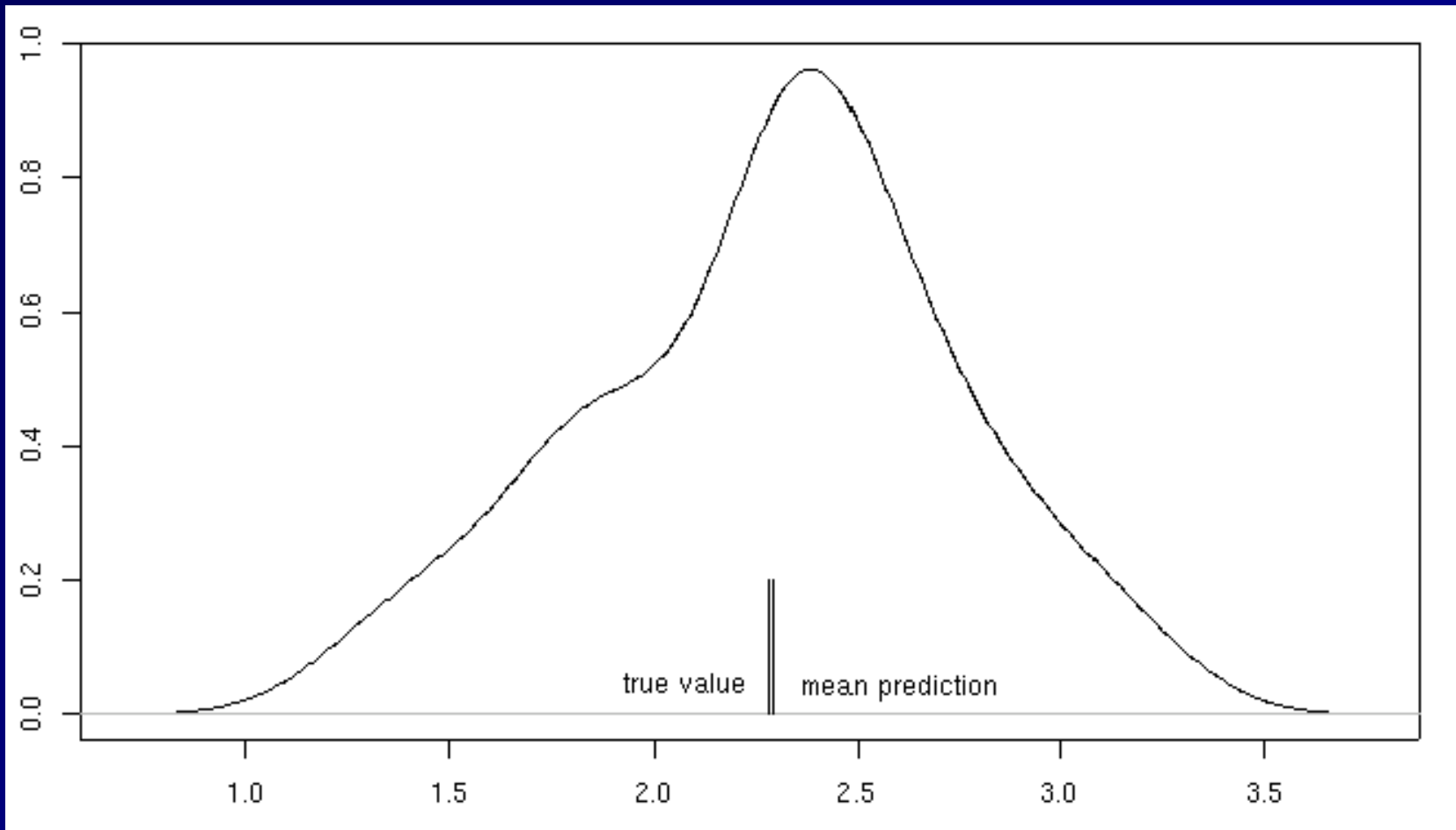
$$\begin{aligned}\mathbb{E}_h[(y^* - h(x^*))^2] &= \mathbb{E}_h[(h(x^*) - \mathbb{E}_h[h(x^*)])^2] && \text{variance} \\ &+ (\mathbb{E}_h[h(x^*)] - f(x^*))^2 && \text{squared bias} \\ &+ (y^* - f(x^*))^2 && \text{noise}\end{aligned}$$



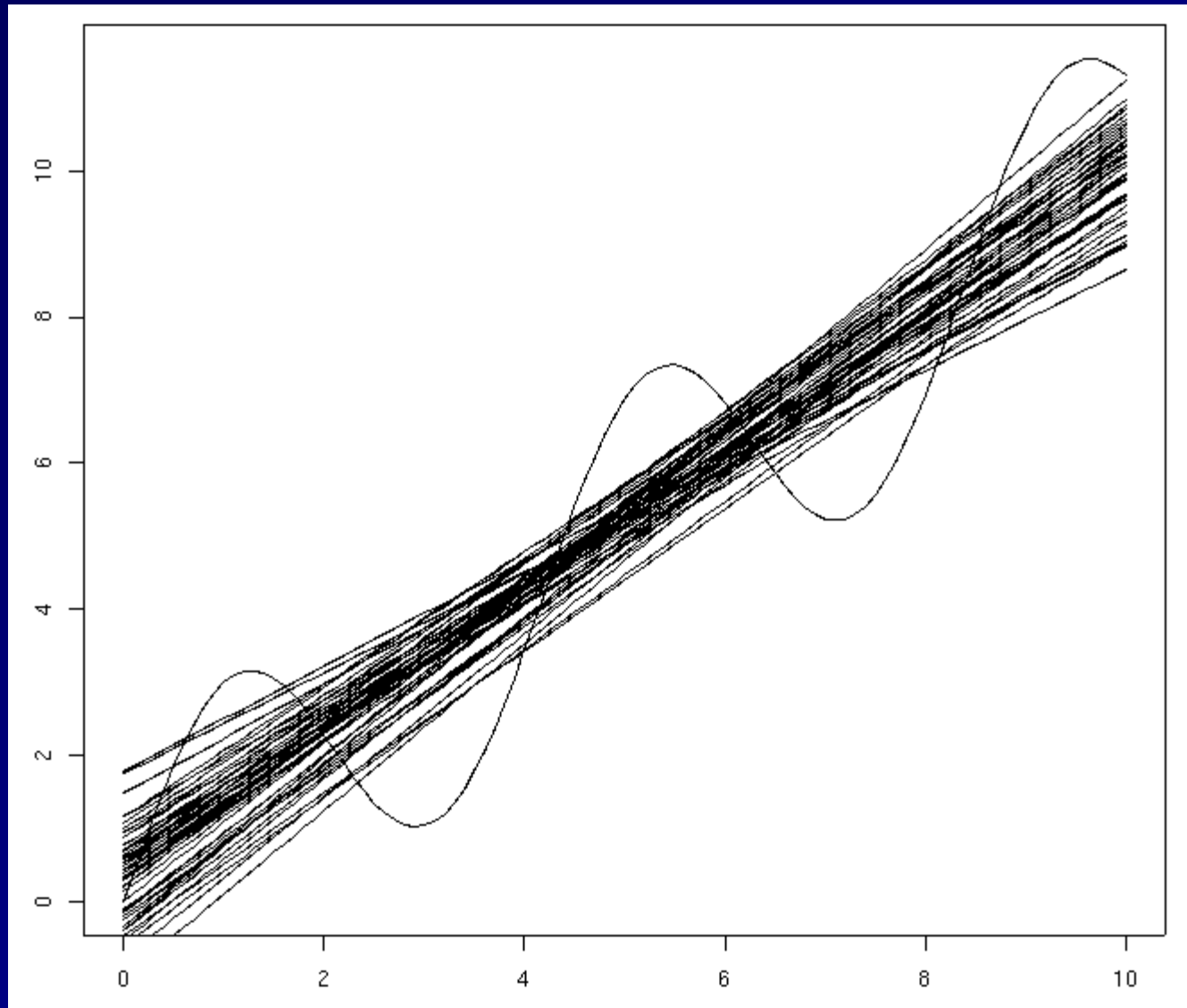
# 50 fits (20 examples each)



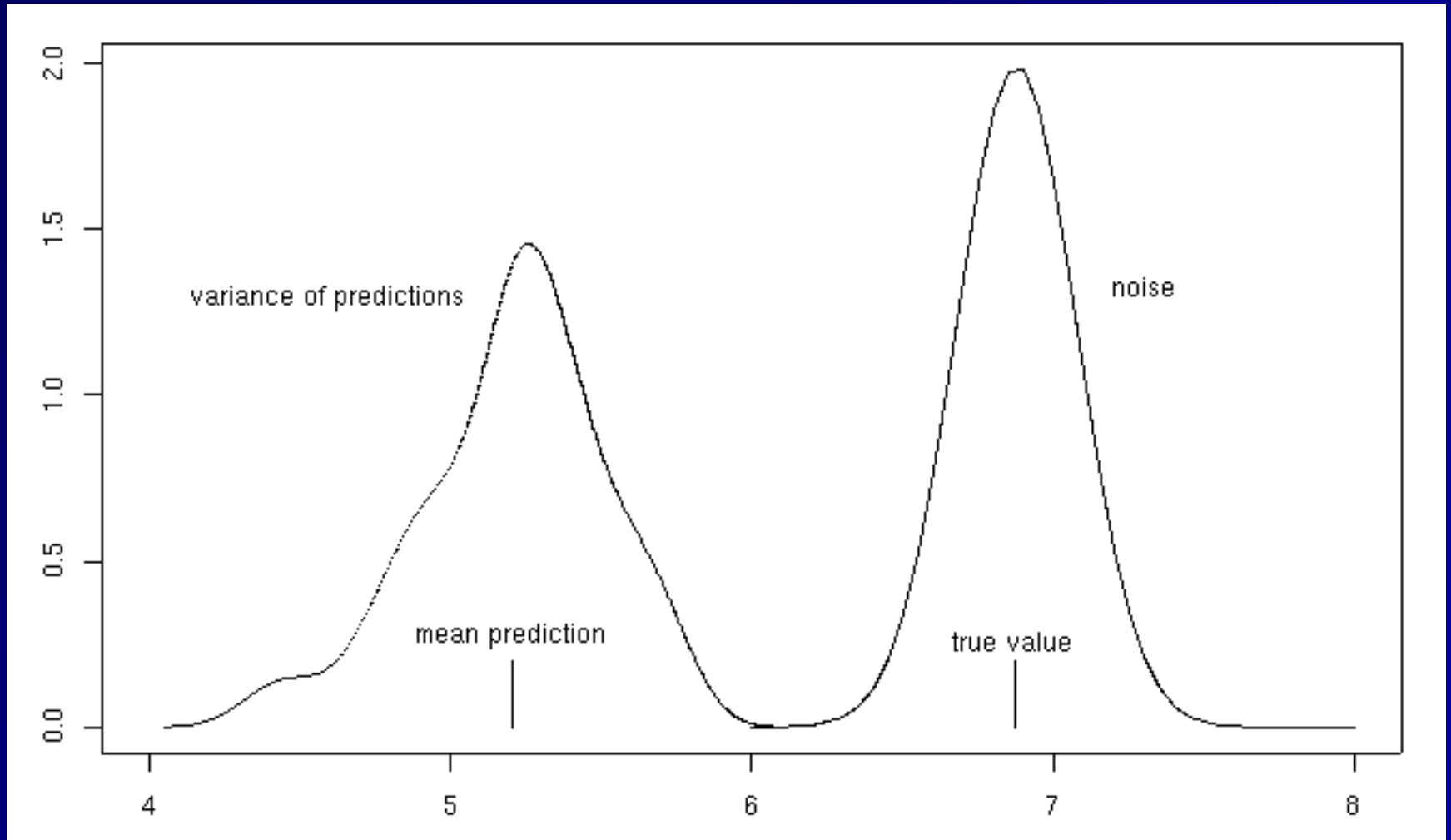
# Distribution of predictions at $x=2.0$



50 fits (20 examples each)



# Distribution of predictions at $x=5.0$



# Derivation

Let  $Z$  be a random variable with distribution  $P(Z)$

Let  $\bar{Z} = \mathbb{E}[Z]$  be the expected value of  $Z$

Lemma:  $\mathbb{E} \left[ (Z - \bar{Z})^2 \right] = \mathbb{E}[Z^2] - \bar{Z}^2$

Corollary:  $\mathbb{E}[Z^2] = \mathbb{E} \left[ (Z - \bar{Z})^2 \right] + \bar{Z}^2$

# Proof

$$\begin{aligned}\mathbb{E} \left[ (Z - \bar{Z})^2 \right] &= \mathbb{E} \left[ Z^2 - 2Z\bar{Z} + \bar{Z}^2 \right] \\ &= \mathbb{E}[Z^2] - 2\bar{Z}\mathbb{E}[Z] + \bar{Z}^2 \\ &= \mathbb{E}[Z^2] - 2\bar{Z}\bar{Z} + \bar{Z}^2 \\ &= \mathbb{E}[Z^2] - \bar{Z}^2\end{aligned}$$

# Derivation of the Decomposition

Expand the quadratic:

$$\mathbb{E}_P[(h(x^*) - y^*)^2] = \mathbb{E}_P[h(x^*)^2 - 2h(x^*)y^* + y^{*2}]$$

Push the expectation inside

$$= \mathbb{E}_P[h(x^*)^2] - 2\mathbb{E}_P[h(x^*)]\mathbb{E}_P[y^*] + \mathbb{E}_P[y^{*2}]$$

Apply the lemma twice

$$= \mathbb{E}_P[(h(x^*) - \mathbb{E}_P[h(x^*)])^2] + \mathbb{E}_P[h(x^*)]^2 - 2\mathbb{E}_P[h(x^*)]f(x^*) \\ + \mathbb{E}_P[(y^* - f(x^*))^2] + f(x^*)^2$$

Collapse the quadratic to get the squared bias term

$$= \mathbb{E}_P[(h(x^*) - \mathbb{E}_P[h(x^*)])^2] + (\mathbb{E}_P[h(x^*)] - f(x^*))^2 \\ + \mathbb{E}_P[(y^* - f(x^*))^2]$$

Note that we are also taking expectations wrt the noise  $\epsilon$

# Measuring Bias and Variance

- In practice (unlike in theory), we have only ONE training set  $S$ .
- We can simulate multiple training sets by bootstrap replicates  
 $S' = \{x \mid x \text{ is drawn at random with replacement from } S\}$  and  $|S'| = |S|$ .



# Procedure for Measuring Bias and Variance

- Construct  $B$  bootstrap replicates of  $S$  (e.g.,  $B = 200$ ):  $S_1, \dots, S_B$
- Apply learning algorithm to each replicate  $S_b$  to obtain hypothesis  $h_b$
- Let  $T_b = S \setminus S_b$  be the data points that do not appear in  $S_b$  (out of bag points)
- Compute predicted value  $h_b(x)$  for each  $x$  in  $T_b$

# Estimating Bias and Variance (continued)

- For each data point  $x$ , we will now have the observed corresponding value  $y$  and several predictions  $y_1, \dots, y_K$ .
- Compute the average prediction  $\bar{y} = \frac{1}{K} \sum_k y_k$ .
- Estimate bias as  $\bar{y} - y$
- Estimate variance as  $\frac{1}{K-1} \sum_k (y_k - \bar{y})^2$
- Assume noise is 0

# Approximations in this Procedure

- Bootstrap replicates are not real fresh data
- We ignore the noise
  - If we have multiple data points with the same  $x$  value, then we can estimate the noise
  - We can also estimate noise by pooling  $y$  values from nearby  $x$  values

# 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

# Ensemble Learning Methods

- Given training sample  $S$
- Generate multiple hypotheses,  $h_1, h_2, \dots, h_L$ .
- Optionally: determining corresponding weights  $\alpha_1, \alpha_2, \dots, \alpha_L$
- Classify new points according to

$$\sum_{\ell} \alpha_{\ell} h_{\ell}(x) > \theta$$

“weighted majority vote”

# Bagging: Bootstrap Aggregating

■ For  $b = 1, \dots, B$  do

$S_b$  = bootstrap replicate of  $S$

Apply learning algorithm to  $S_b$  to learn  $h_b$

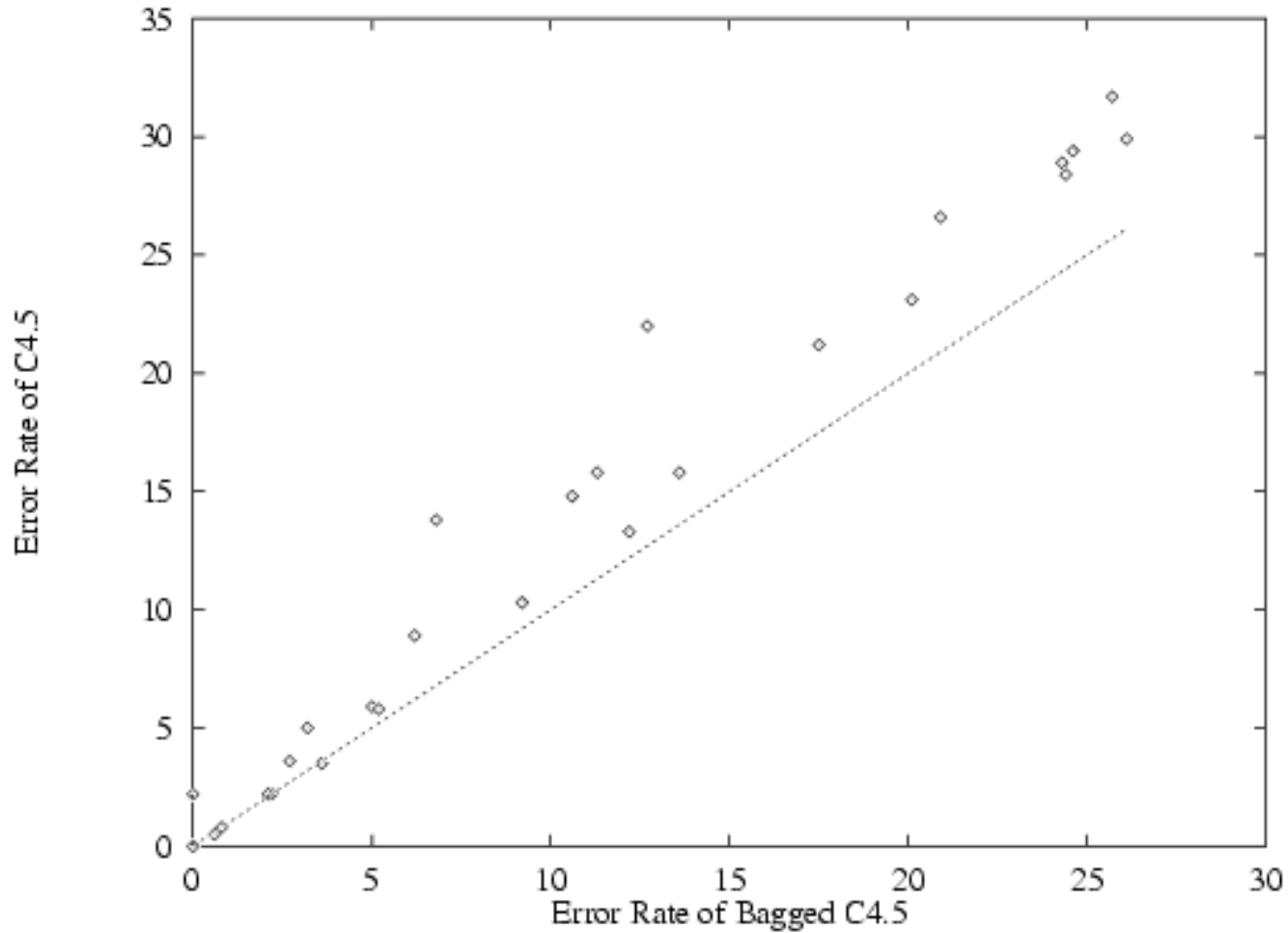
■ Classify new points by unweighted vote:

$$\frac{1}{B} \sum_b h_b(x) > 0$$

# Estimated Bias and Variance of Bagging

- If we estimate bias and variance using the same  $B$  bootstrap samples, we will have:
  - Bias =  $\bar{y} - y$  [same as before]
  - Variance =  $\frac{1}{K-1} \sum_k (\bar{y} - \bar{y})^2 = 0$
- Hence, according to this approximate way of estimating variance, bagging removes the variance while leaving bias unchanged.
- In reality, bagging only *reduces* variance and tends to slightly increase bias

# Bagging Decision Trees (Freund & Schapire)





# Bayesian Ensembles: Bayesian Model Averaging

- $P(y^*|x^*, S) = \int_{\mathbf{w}} P(\mathbf{w})P(S|\mathbf{w})P(y^*|x^*, \mathbf{w})d\mathbf{w}$

- where  $P(S|\mathbf{w}) = \prod_i P(y_i|\mathbf{w})P(x_i|y_i, \mathbf{w})$

- This is rarely practical to evaluate, but suppose we could sample some good values for  $\mathbf{w}$ :  $(\mathbf{w}_1, \dots, \mathbf{w}_C)$

- We could approximate the integral by a sum:

$$P(y^*|x^*, S) = \sum_c P(\mathbf{w}_c|S)P(y^*|x^*, \mathbf{w}_c)$$

- This is called Bayesian Model Averaging

# Review of Part 2

- Goal: making accurate predictions on *new* data points
  - the problem of Overfitting
  - occurs when the model becomes too complex for the amount of data
- Complexity can be controlled
  - regularization penalty
  - early stopping
  - choosing  $\lambda$  by holdout or cross-validation
- Bias-variance error decomposition
  - Squared loss can be decomposed into  $\text{bias}^2 + \text{variance} + \text{noise}$
  - bias and variance can be (approximately) measured using bootstrapping
  - they provide a diagnostic tool for machine learning
- Bagging is an ensemble method that applies bootstrapping to reduce variance
  - Bagging a low-bias, high-variance classifier can produce excellent results
- Bayesian analysis provides an alternative view of
  - regularization penalty = log prior
  - ensemble methods = integrating out the prior

# Questions?

# Questions to think about

- Most machine learning methods involve complex, flexible models
  - decision trees
  - support vector machines
  - neural networks (esp. deep ones)

Hence, complexity control (variance management) is a central challenge

“drop out” is a cool new technique in this area
- The bias-variance analysis was done for regression. Can it be extended to classification?
  - Yes, see James (2003) “Variance and bias for general loss functions” *Machine Learning*.