

CHAPTER 4:

# PARAMETRIC METHODS

# Parametric Estimation

2

□  $X = \{x^t\}_t$  where  $x^t \sim p(x)$

“distributed as”

probability density function

□ Parametric estimation:

Assume a form for  $p(x|\theta)$  and estimate  $\theta$ , its sufficient statistics, using  $X$

e.g., Normal distribution  $N(\mu, \sigma^2)$  where  $\theta = \{\mu, \sigma^2\}$

Remember that a probability density function is defined as:

$$p(x_0) \equiv \lim_{\varepsilon \rightarrow 0} P(x_0 \leq x < x_0 + \varepsilon)$$

# Maximum Likelihood Estimation

3

- Likelihood of  $\theta$  given the sample  $\mathbf{X}$

$$l(\theta|\mathbf{X}) = p(\mathbf{X}|\theta) = \prod_t p(x^t|\theta)$$

- Log likelihood

$$L(\theta|\mathbf{X}) = \log l(\theta|\mathbf{X}) = \sum_t \log p(x^t|\theta)$$

- Maximum likelihood estimator (MLE)

$$\theta^* = \operatorname{argmax}_{\theta} L(\theta|\mathbf{X})$$

# Examples: Bernoulli/Multinomial

4

- Bernoulli: Two states, failure/success,  $x$  in  $\{0,1\}$  (e.g, coin)

$$P(x) = p_o^x (1 - p_o)^{(1-x)} \quad p_o \text{ is the parameter for probability of success}$$

$$\mathcal{L}(p_o | \mathbf{X}) = \log \prod_t p_o^{x^t} (1 - p_o)^{(1-x^t)}$$

$$\text{MLE: } p_o = \sum_t x^t / N$$

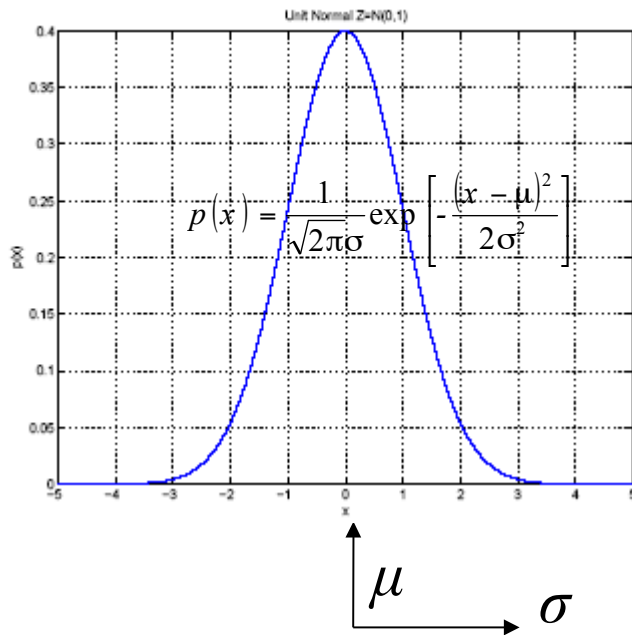
- Multinomial:  $K > 2$  states,  $x_i$  in  $\{0,1\}$  (e.g., die with 6 faces)

$$P(x_1, x_2, \dots, x_K) = \prod_i p_i^{x_i} \quad p_i \text{ is the parameter for probability of state } i$$

success

$$\mathcal{L}(p_1, p_2, \dots, p_K | \mathbf{X}) = \log \prod_t \prod_i p_i^{x_i^t}$$

# Gaussian (Normal) Distribution



□  $p(x) = N(\mu, \sigma^2)$

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

□ MLE for  $\mu$  and  $\sigma^2$ :

$$m = \frac{\sum_t x^t}{N}$$
$$s^2 = \frac{\sum_t (x^t - m)^2}{N}$$

# Bias and Variance

6

Unknown parameter  $\theta$

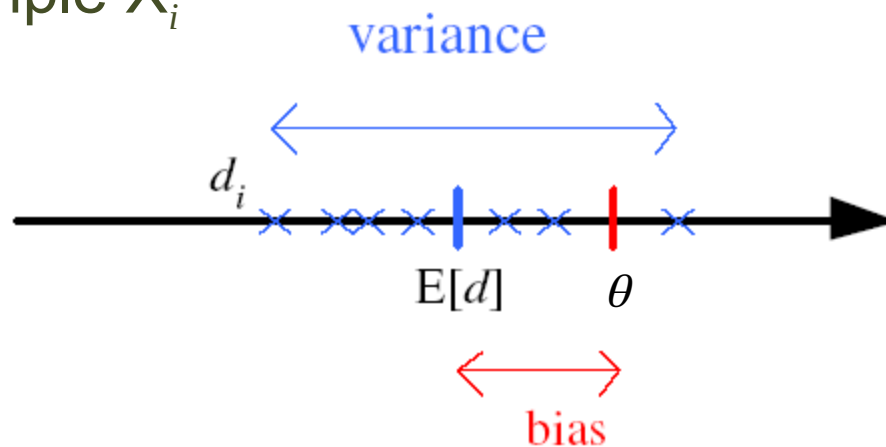
Estimator  $d_i = d(X_i)$  on sample  $X_i$

Bias:  $b_\theta(d) = E[d] - \theta$

Variance:  $E[(d - E[d])^2]$

Mean square error:

$$\begin{aligned} r(d, \theta) &= E[(d - \theta)^2] \\ &= (E[d] - \theta)^2 + E[(d - E[d])^2] \\ &= \text{Bias}^2 + \text{Variance} \end{aligned}$$



the smaller the bias and the variance, the better

# Bayes' Estimator

7

- Treat  $\theta$  as a random var with prior  $p(\theta)$  this prior may come from domain expertise, without looking at  $\mathbf{X}$
- Bayes' rule:  $p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta) p(\theta)}{p(\mathbf{X})}$  posterior density of  $\theta$       likelihood density
- Full:  $p(x|\mathbf{X}) = \int p(x|\theta) p(\theta|\mathbf{X}) d\theta$
- Maximum a Posteriori (MAP):  $\theta_{\text{MAP}} = \operatorname{argmax}_{\theta} p(\theta|\mathbf{X})$
- Maximum Likelihood (ML):  $\theta_{\text{ML}} = \operatorname{argmax}_{\theta} p(\mathbf{X}|\theta)$
- Bayes':  $\theta_{\text{Bayes'}} = E[\theta|\mathbf{X}] = \int \theta p(\theta|\mathbf{X}) d\theta$

# Comparing ML, MAP, and Bayes'

8

Let  $\Theta$  be the set of all possible solutions  $\theta$ 's

- **Maximum a Posteriori (MAP):** Selects the  $\theta$  that satisfies  $\theta_{\text{MAP}} = \operatorname{argmax}_{\theta} p(\theta|X)$
- **Maximum Likelihood (ML):** Selects the  $\theta$  that satisfies  $\theta_{\text{ML}} = \operatorname{argmax}_{\theta} p(X|\theta)$
- **Bayes':** Constructs the “weighted average” over all  $\theta$ 's in  $\Theta$ :  $\theta_{\text{Bayes'}} = E[\theta|X] = \int \theta p(\theta|X) d\theta$

Note that if the  $\theta$ 's in  $\Theta$  are uniformly distributed then  $\theta_{\text{MAP}} = \theta_{\text{ML}}$

since  $\theta_{\text{MAP}} = \operatorname{argmax}_{\theta} p(\theta|X)$

$$= \operatorname{argmax}_{\theta} p(X|\theta) p(\theta)/p(X) = \operatorname{argmax}_{\theta} p(X|\theta) = \theta_{\text{ML}}$$



# Bayes' Estimator: Example

- If  $x^t \sim N(\theta, \sigma_o^2)$  and  $\theta \sim N(\mu, \sigma^2)$   
where  $\sigma_o^2$ ,  $\mu$ , and  $\sigma^2$  are known

- then:

- $\theta_{ML} = m$  *sample mean*

- $\theta_{MAP} = \theta_{\text{Bayes}'} =$

$$E[\theta | \mathbf{X}] = \frac{N/\sigma_o^2}{N/\sigma_o^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma_o^2 + 1/\sigma^2} \mu$$

Note: if the  $\theta$ 's in  $\Theta$  are normally distributed and  $p(\mathbf{X}|\theta)$  is normal, then  $p(\theta|\mathbf{X})$  is normal and  $\theta_{MAP} = \theta_{\text{Bayes}'}$

# Parametric Classification

10

Remember from Chapter 3, where  $g_i$  is a discriminant function

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})}$$
$$= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)}$$

$$g_i(x) = p(x | C_i)P(C_i)$$

or

$$g_i(x) = \log p(x | C_i) + \log P(C_i)$$

If  $p(x | C_i)$  are Gaussian, then: (here “log” is natural log)

$$p(x | C_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right]$$

$$g_i(x) = -\frac{1}{2}\log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$

- Given the sample  $\mathbf{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$

$$x \in \mathfrak{R} \quad r_i^t = \begin{cases} 1 & \text{if } x^t \in \mathcal{C}_i \\ 0 & \text{if } x^t \in \mathcal{C}_j, j \neq i \end{cases}$$

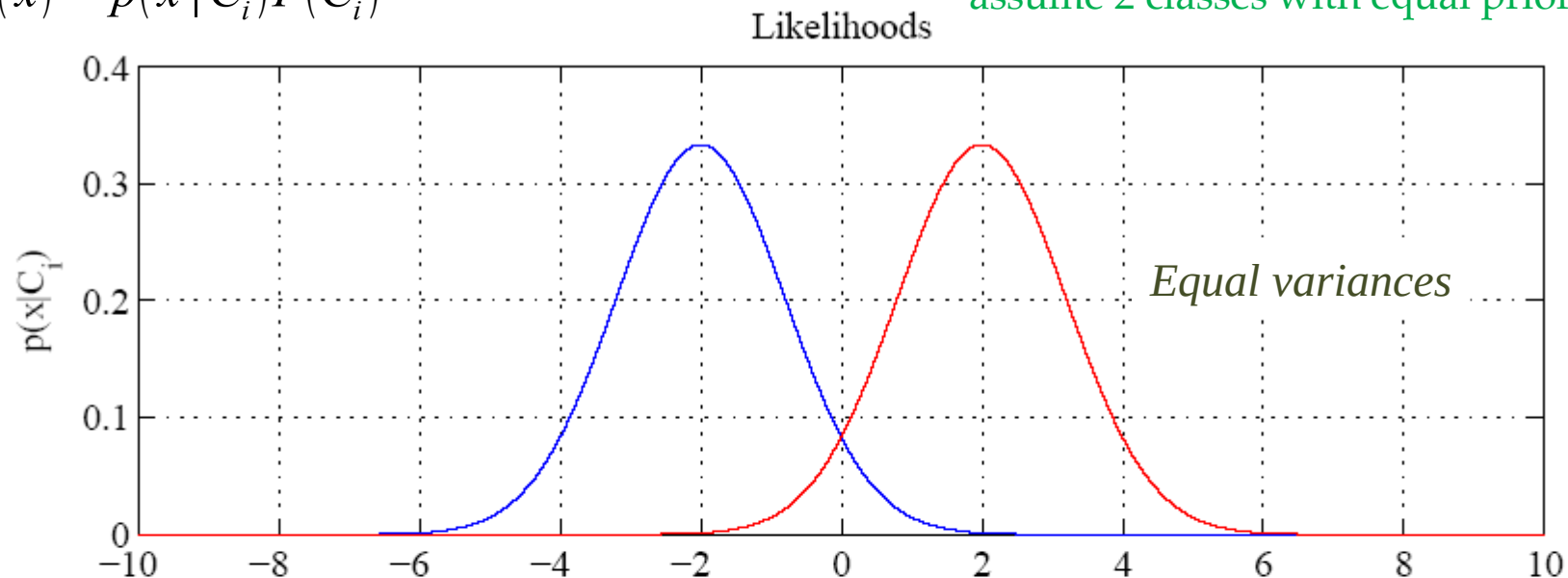
- ML estimates are

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad m_i = \frac{\sum_t x^t r_i^t}{\sum_t r_i^t} \quad s_i^2 = \frac{\sum_t (x^t - m_i)^2 r_i^t}{\sum_t r_i^t}$$

- Discriminant  $g_i(\mathbf{x}) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(\mathbf{x} - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$

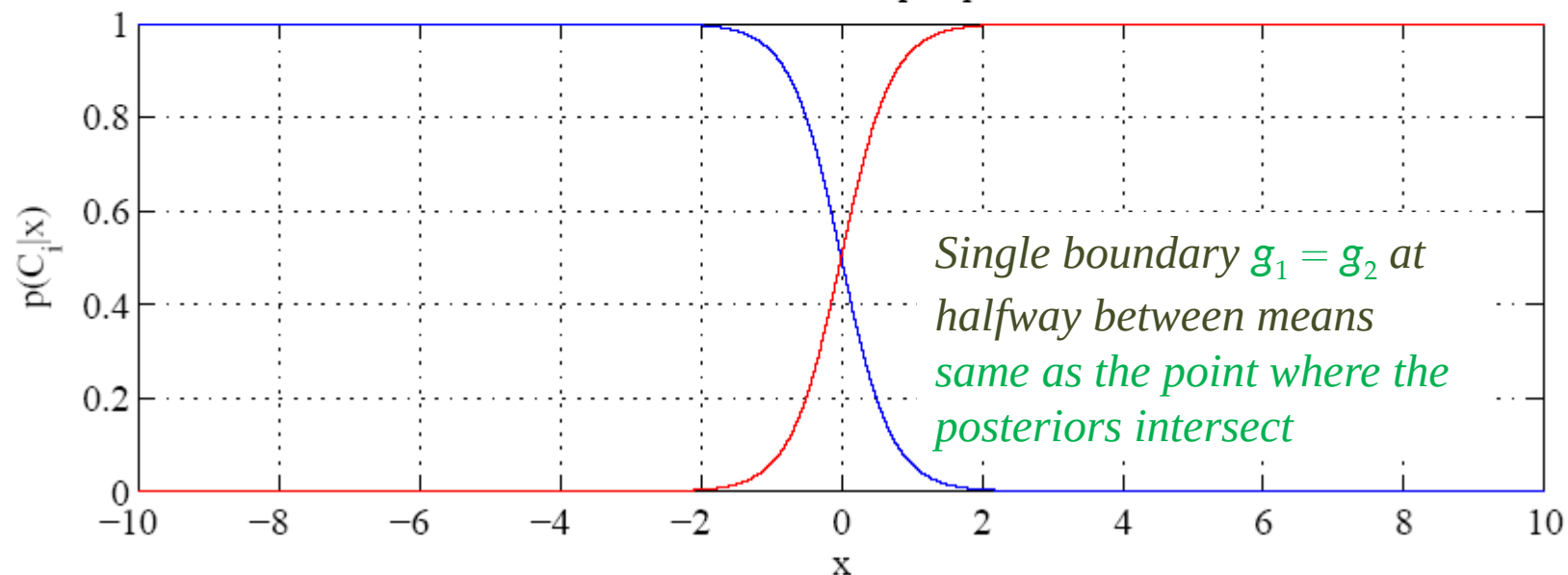
$$g_i(x) = p(x | C_i)P(C_i)$$

assume 2 classes with equal priors  $P(C_i)$

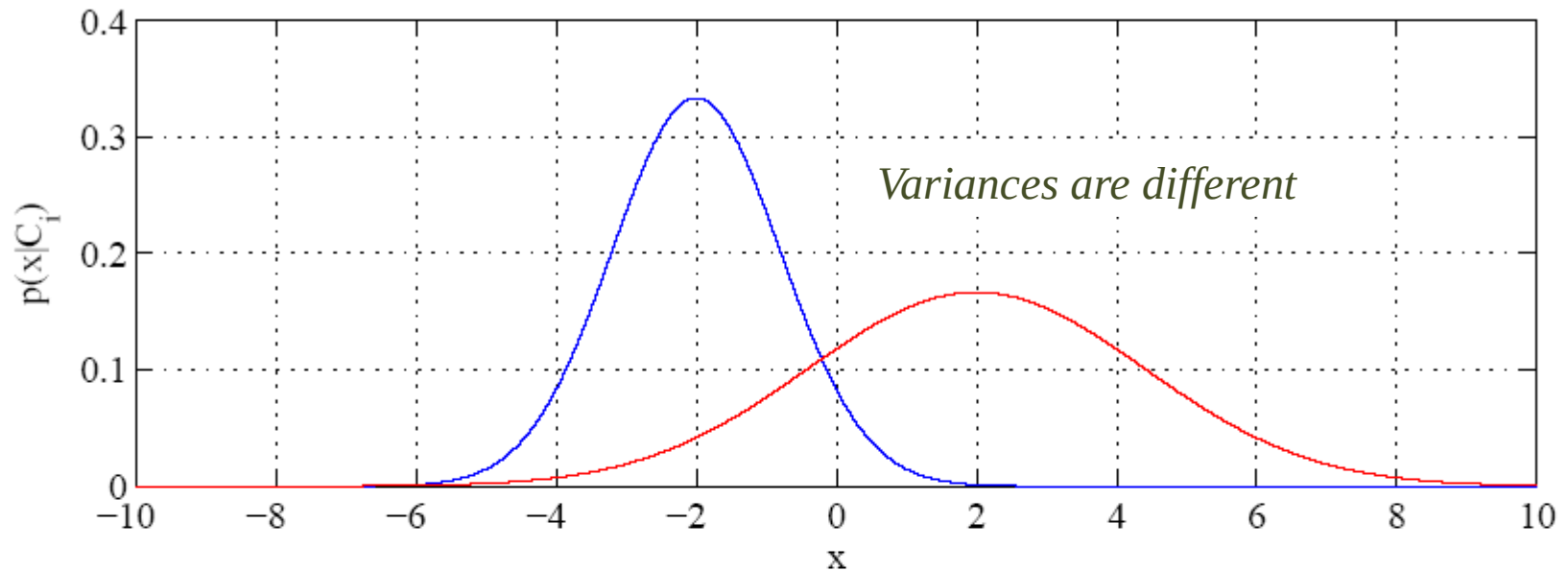


In this case,  $g_i = - (x - m_i)^2$ , thus we assign  $x$  to class with nearest mean

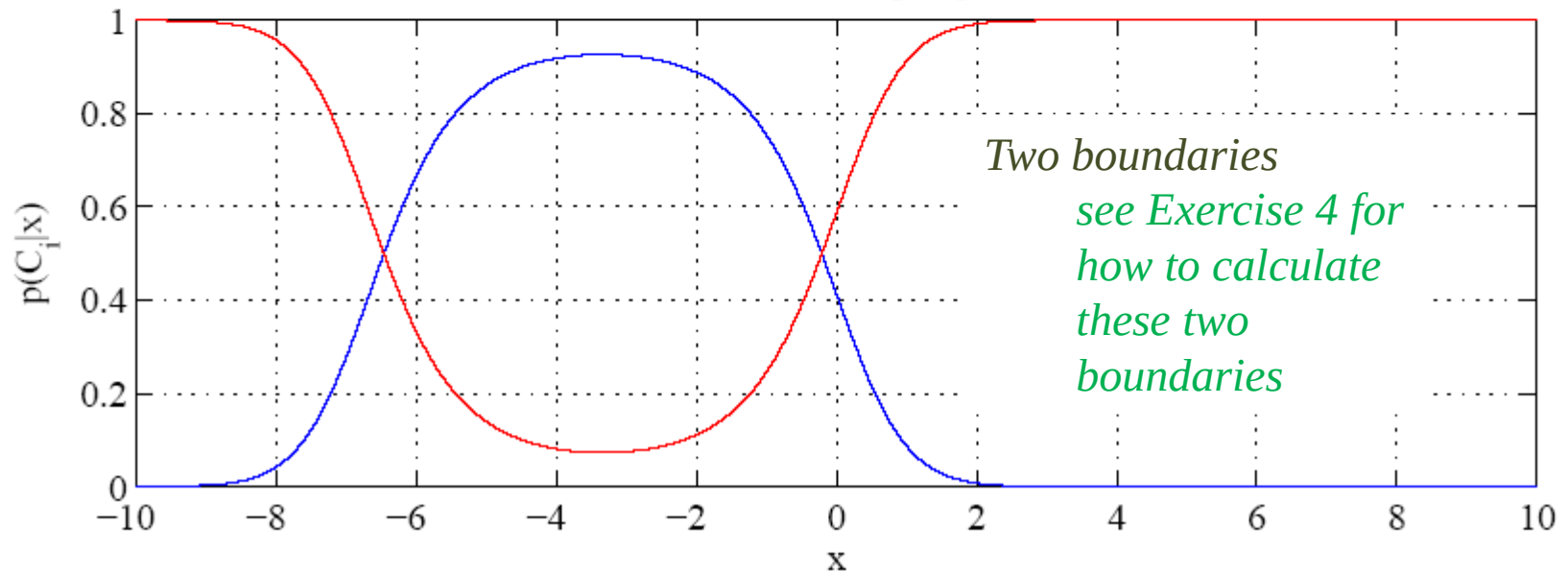
Posterior with equal priors



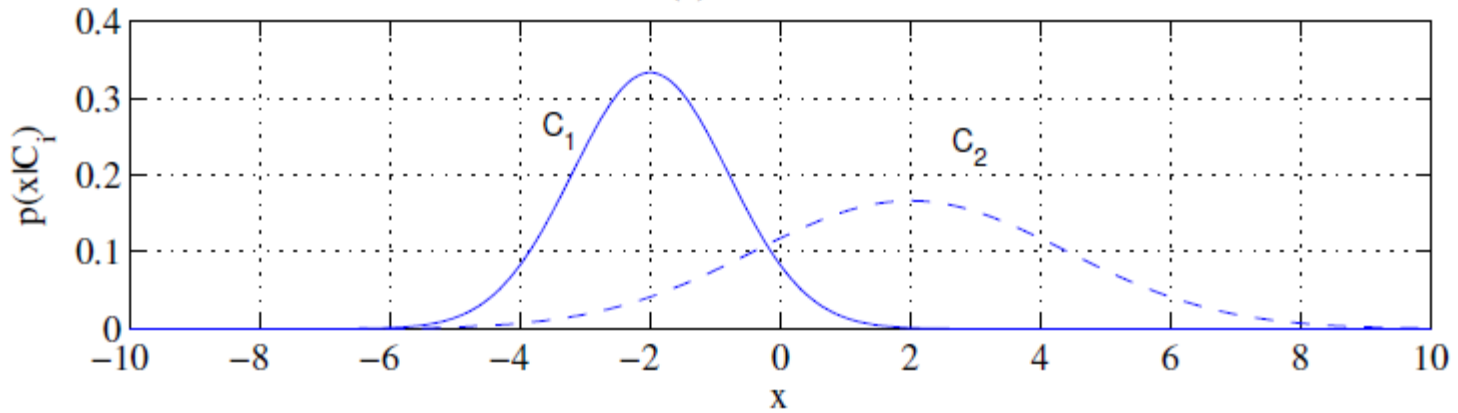
Likelihoods



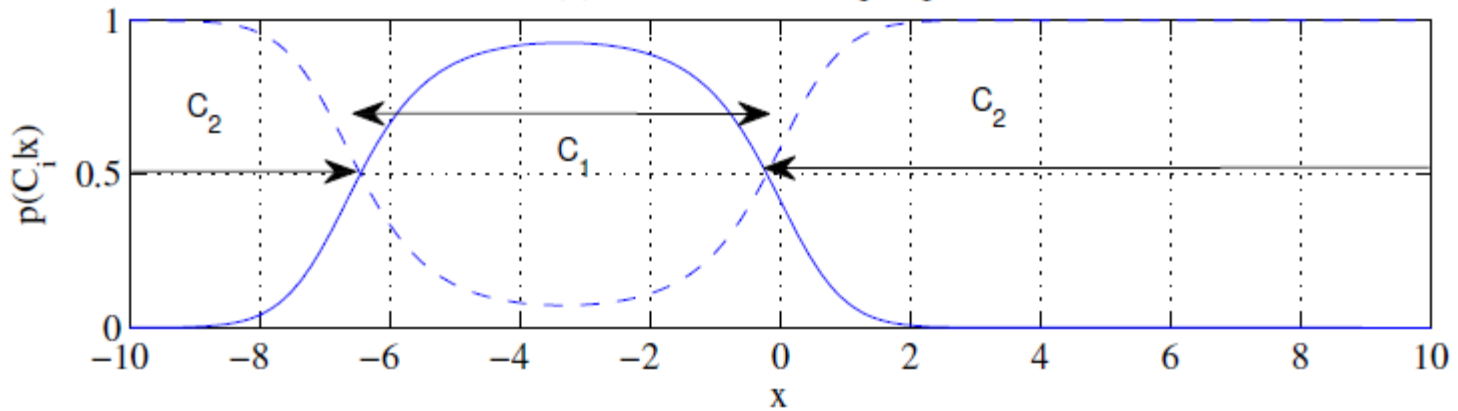
Posteriors with equal priors



(a) Likelihoods

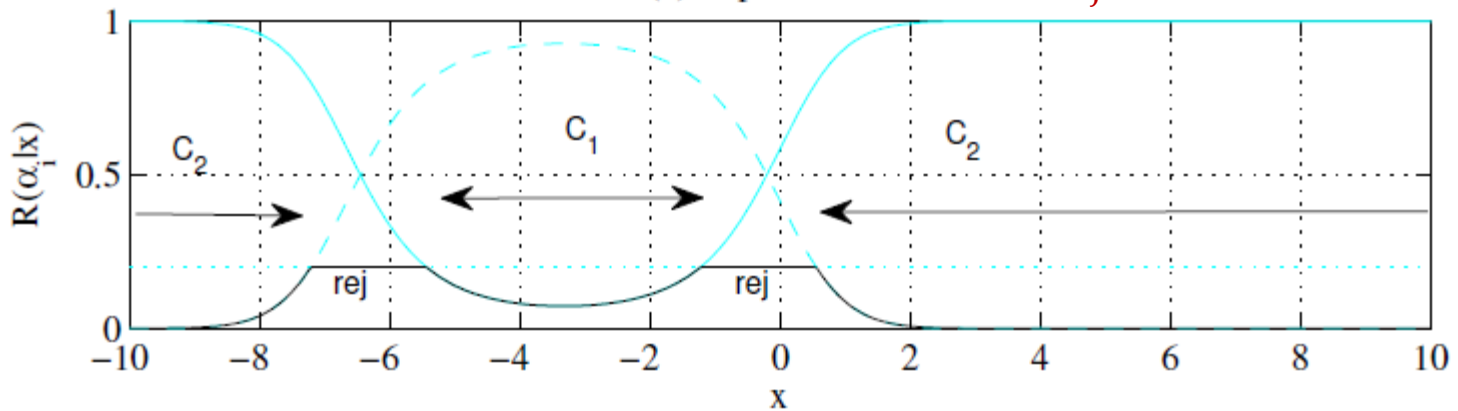


(b) Posteriors with equal priors



(c) Expected risks

Reject with  $\lambda = 0.2$



# Parametric Classification – Notes

15

- Note that the equations for parametric classification derived in the last few slides assume the data follows a Gaussian distribution
- Hence, you need to determine whether or not your univariate data  $X$  follows a Gaussian distribution

*using a statistical test for normality (e.g., Shapiro–Wilk test, Kolmogorov–Smirnov test, Lilliefors test, ...)*

before you consider applying these equations

# Regression

16

$$r = f(x) + \varepsilon$$

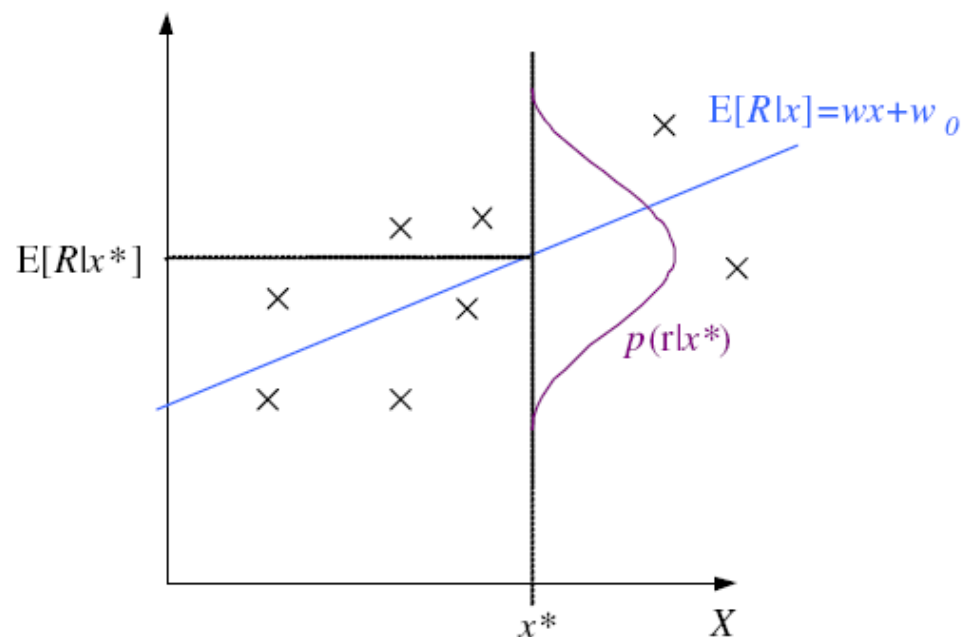
$$\text{estimator} : g(x | \theta)$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$p(r | x) \sim \mathcal{N}(g(x | \theta), \sigma^2)$$

$$L(\theta | \mathbf{X}) = \log \prod_{t=1}^N p(x^t, r^t)$$

$$= \log \prod_{t=1}^N p(r^t | x^t) + \log \prod_{t=1}^N p(x^t)$$



second term can be ignored



# Regression: From LogL to Error

17

$$\begin{aligned} L(\theta | \mathbf{X}) &= \log \prod_{t=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{[r^t - g(x^t | \theta)]^2}{2\sigma^2} \right] \\ &= -N \log \sqrt{2\pi}\sigma - \frac{1}{2\sigma^2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2 \\ E(\theta | \mathbf{X}) &= \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2 \end{aligned}$$

Maximizing the log likelihood above is the same as minimizing the Error.

$\theta$  that minimizes Error is called “least squares estimate”

Trick: When maximizing a likelihood  $l$  that contains exponents, instead minimize error  $E = -\log l$

# Linear Regression

$$g(x^t | w_1, w_0) = w_1 x^t + w_0$$

$$\sum_t r^t = Nw_0 + w_1 \sum_t x^t$$

$$\sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2$$

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix}$$

$$\mathbf{A} \mathbf{w} = \mathbf{y} \quad \text{and so} \quad \mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$$

# Polynomial Regression

$$g(x^t | w_k, \dots, w_2, w_1, w_0) = w_k (x^t)^k + \dots + w_2 (x^t)^2 + w_1 x^t + w_0$$

$$\mathbf{D} = \begin{bmatrix} 1 & x^1 & (x^1)^2 & \dots & (x^1)^k \\ 1 & x^2 & (x^2)^2 & \dots & (x^2)^k \\ \vdots & & & & \\ 1 & x^N & (x^N)^2 & \dots & (x^N)^k \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

$$\mathbf{A} \mathbf{w} = \mathbf{y} \quad \text{where} \quad \mathbf{A} = (\mathbf{D}^T \mathbf{D}) \quad \text{and} \quad \mathbf{y} = \mathbf{D}^T \mathbf{r}$$

$$\text{then } \mathbf{w} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{r}$$

# Maximum Likelihood and Least Squares

20

Taken from Tom Mitchell's Machine Learning textbook:

*“... under certain assumptions (\*) any learning algorithm that minimizes the squared error between the output hypothesis predictions and the training data will output a maximum likelihood hypothesis.”*

(\*): Assumption:

“the observed training target values are generated by adding random noise to the true target value, where the random noise is drawn independently for each example from a Normal distribution with zero mean.”

# Other Error Measures

21

- Square Error:

$$E(\theta | \mathbf{X}) = \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2$$

- Relative Square Error ( $E_{\text{RSE}}$ ):

$$E(\theta | \mathbf{X}) = \frac{\sum_{t=1}^N [r^t - g(x^t | \theta)]^2}{\sum_{t=1}^N [r^t - \bar{r}]^2}$$

- Absolute Error:

$$E(\theta | \mathbf{X}) = \sum_t |r^t - g(x^t | \theta)|$$

- $\epsilon$ -sensitive Error:

$$E(\theta | \mathbf{X}) = \sum_t 1(|r^t - g(x^t | \theta)| > \epsilon) (|r^t - g(x^t | \theta)| - \epsilon)$$

- $R^2$ : Coefficient of Determination:  $R^2 = 1 - E_{\text{RSE}}$  (see next slide)

# Coefficient of Determination: $R^2$

22

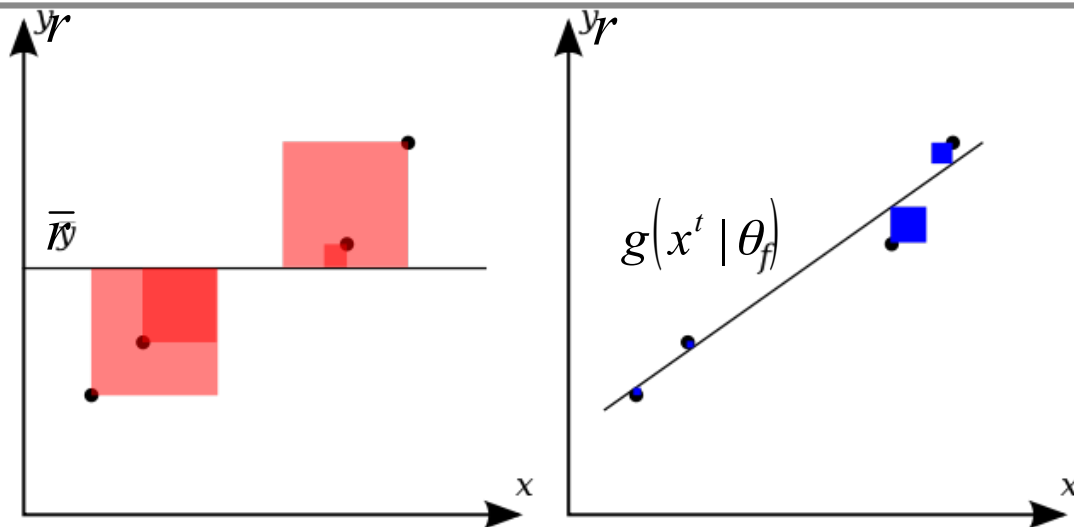


Figure adapted from Wikipedia (Sept. 2015)

"Coefficient of Determination" by Orzetto - Own work. Licensed under CC BY-SA 3.0 via Commons –

[https://commons.wikimedia.org/wiki/File:Coefficient\\_of\\_Determination.svg#/media/File:Coefficient\\_of\\_Determination.svg](https://commons.wikimedia.org/wiki/File:Coefficient_of_Determination.svg#/media/File:Coefficient_of_Determination.svg)

$$R^2 = 1 - \frac{\sum_{t=1}^N [r^t - g(x^t | \theta)]^2}{\sum_{t=1}^N [r^t - \bar{r}]^2}$$

*The closer  $R^2$  is to 1 the better*  
as this means that the  $g(.)$  estimates  
(on the right graph) fit the data well  
in  
comparison to the simple estimate  
given by the average value (on the

# Bias and Variance

23

- Given  $X = \{x^t, r^t\}$  drawn from unknown pdf  $p(x, r)$   
Expected Square Error of the  $g(\cdot)$  estimate at  $x$ :

$$E[(r - g(x))^2 | x] = E[(r - E[r | x])^2 | x] + (E[r | x] - g(x))^2$$

*noise*  
*variance of noise added*  
*does not depend on  $g(\cdot)$  or  $X$*

*squared error*  
 *$g(x)$  deviation from  $E[r|x]$*   
*depend on  $g(\cdot)$  and  $X$*

it may be that for a sample  $X$ ,  $g(\cdot)$  is a very good fit, and for another sample it is not

- Given samples  $X$ 's, all of size  $N$  drawn from the same joint density  $p(x, r)$ . Expected value,

$$E_X[E[r | x] - g(x)] = (E[r | x] - E_X[g(x)])^2 + E_X[(g(x) - E_X[g(x)])^2]$$

*bias*  
*how much  $g(\cdot)$  is wrong*  
*disregarding effect of varying samples*

*variance*  
*how much  $g(\cdot)$  fluctuates*  
*as samples varies*

# Estimating Bias and Variance

24

- $M$  samples  $X_i = \{x_i^t, r_i^t\}$ ,  $i=1, \dots, M$   
are used to fit  $g_i(x)$ ,  $i=1, \dots, M$

$$\text{Bias}^2(g) = \frac{1}{N} \sum_t \left[ \bar{g}(x^t) - f(x^t) \right]^2$$

$$\text{Variance}(g) = \frac{1}{NM} \sum_t \sum_i \left[ g_i(x^t) - \bar{g}(x^t) \right]^2$$

$$\text{where } \bar{g}(x) = \frac{1}{M} \sum_i g_i(x)$$



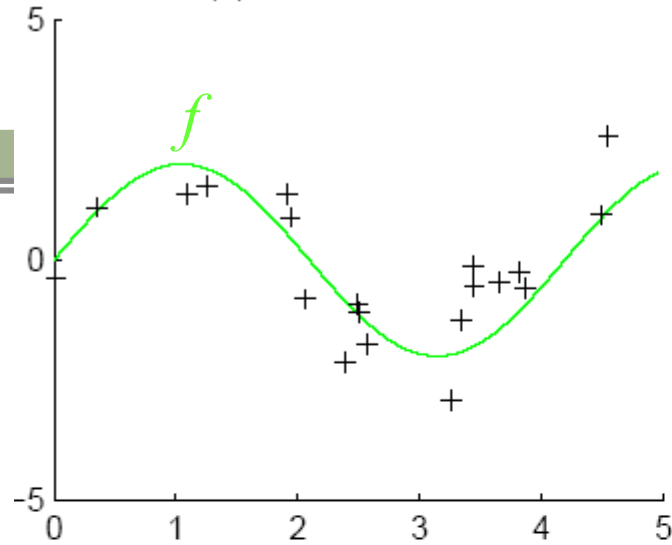
# Bias/Variance Dilemma

25

- Examples:
  - ▣  $g_i(x)=2$  has no variance and high bias
  - ▣  $g_i(x)=\sum_t r_i^t/N$  has lower bias but higher variance
- As we increase complexity of  $g(\cdot)$ ,  
bias decreases (a better fit to data) and  
variance increases (fit varies more with data)
- Bias/Variance dilemma: (Geman et al., 1992)

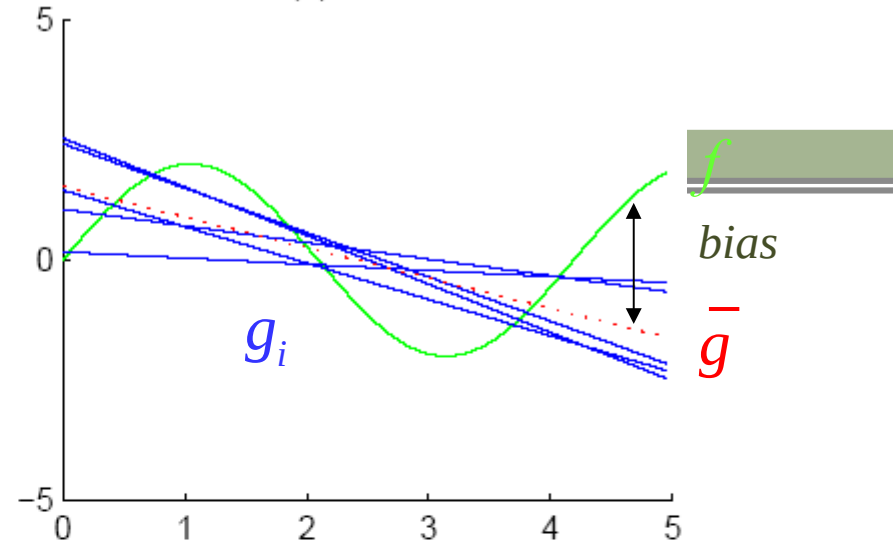
added noise  $\sim N(0,1)$

(a) Function and data



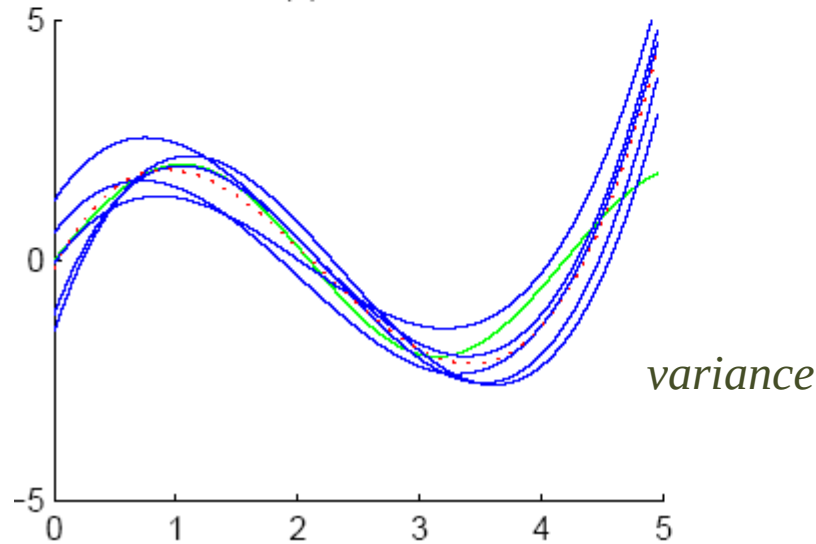
A linear regression for each of 5 samples

(b) Order 1

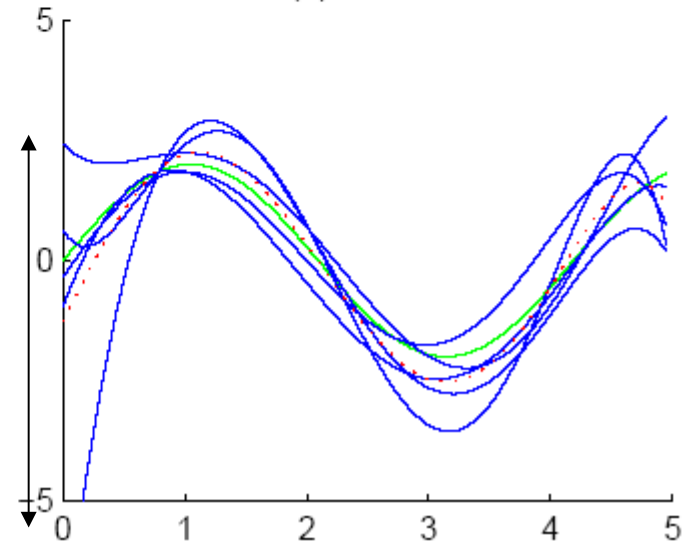


Polynomial regression of order 3 (left) and order 5 (right) for each of 5 samples

(c) Order 3



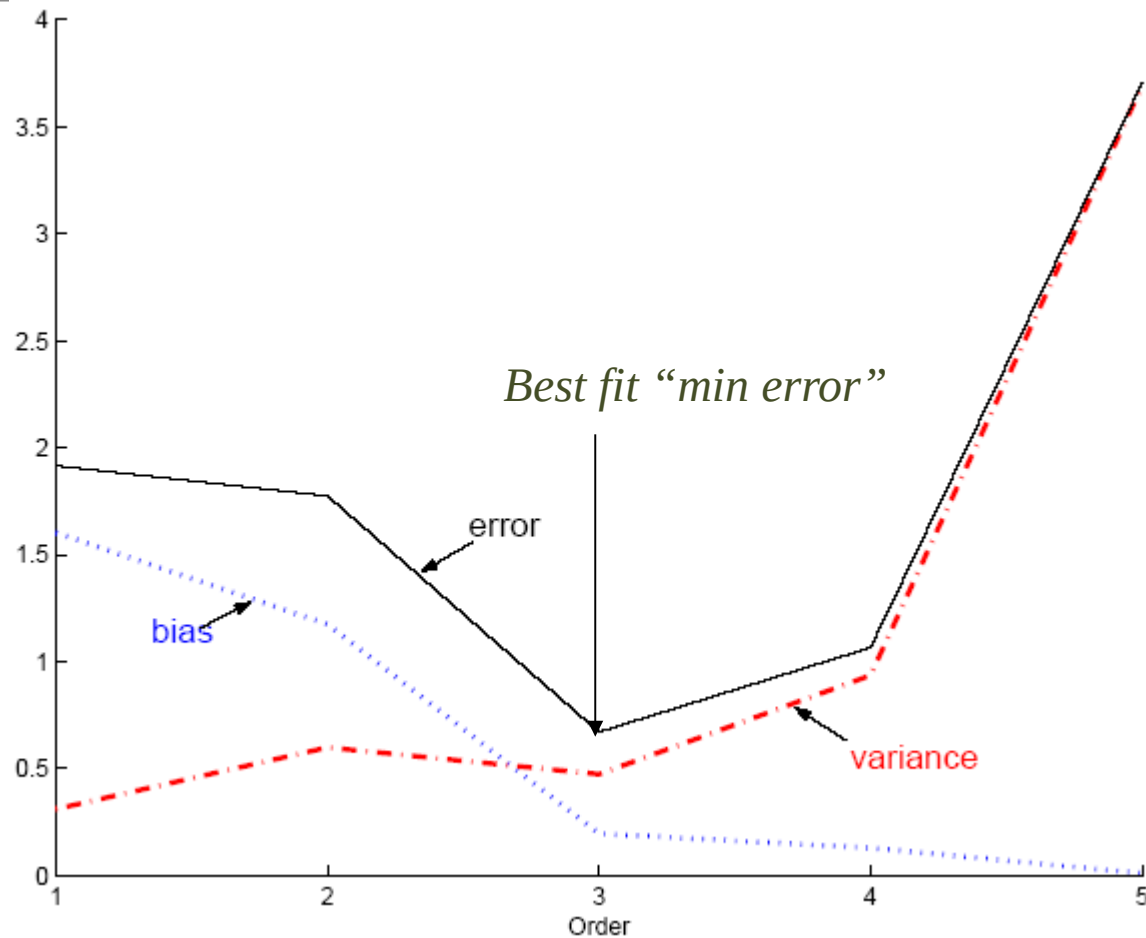
(d) Order 5



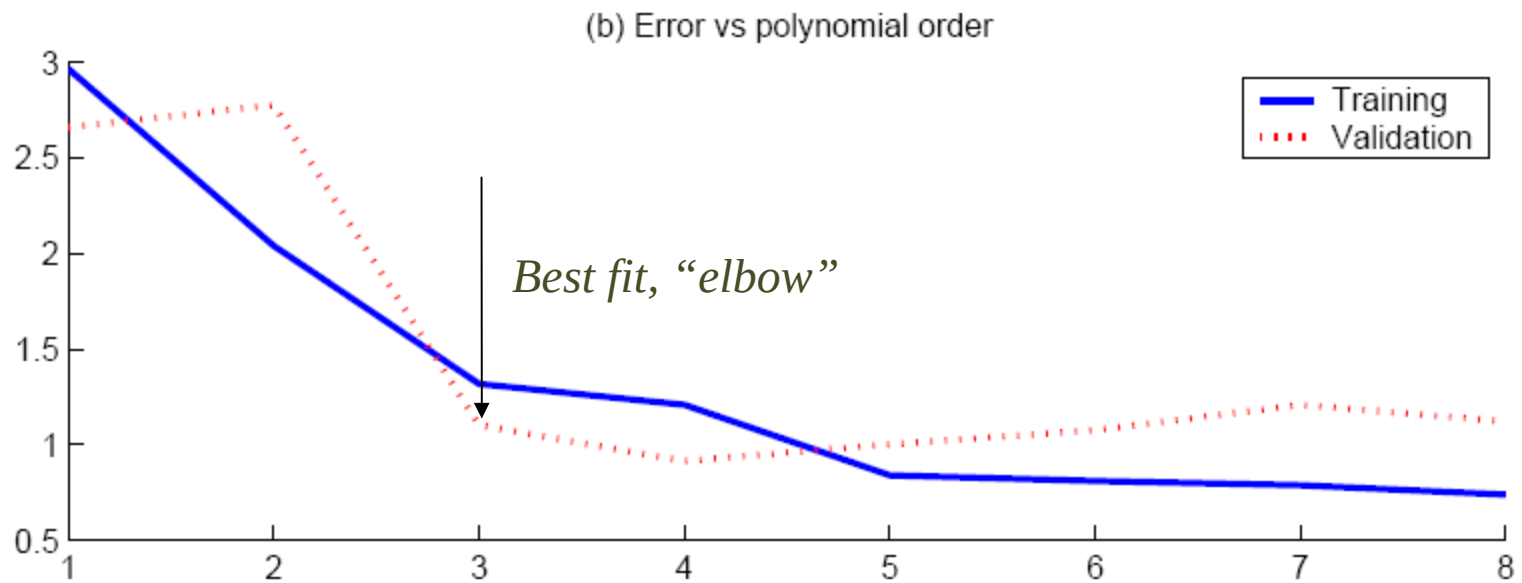
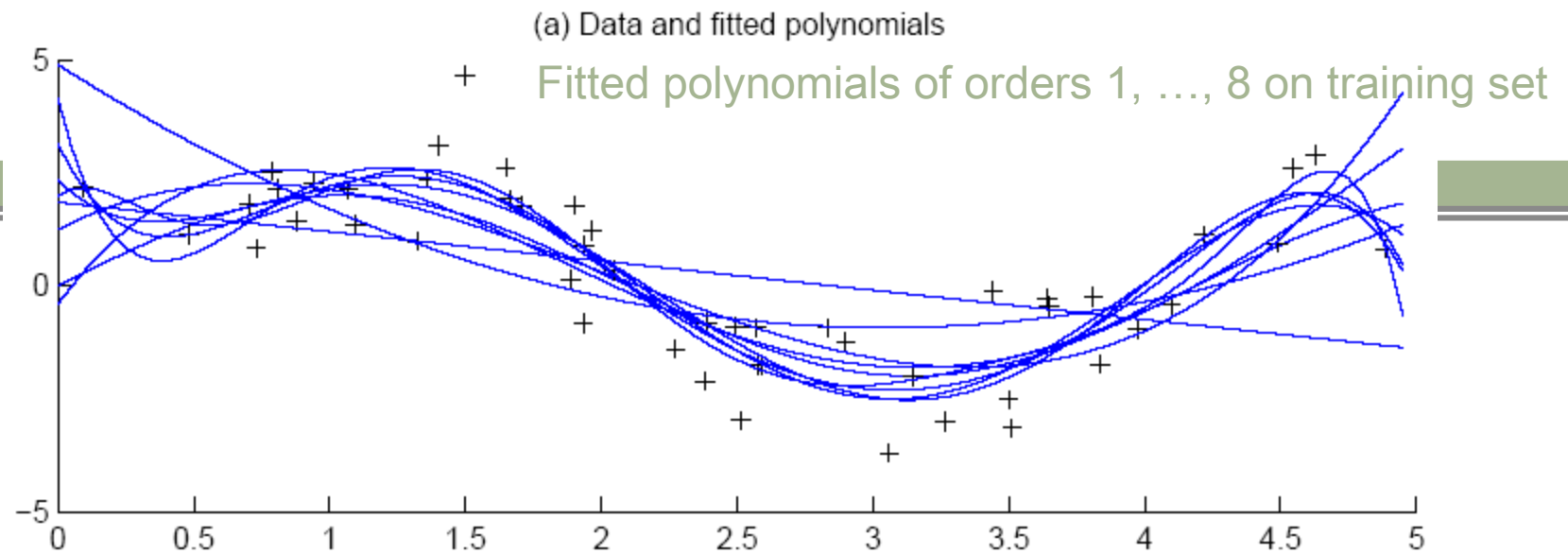
dotted red line in each plot is the average of the five  $\hat{f}_i(x)$

# Polynomial Regression

27



Same settings as plots on the previous slide, but using 100 models instead of 5



Settings as plots on previous slides, but using training and validation sets (50 instances each)

# Model Selection (1 of 2 slides)

## Methods to fine-tune model complexity

29

- Cross-validation: Measures generalization accuracy by testing on data unused during training
- Regularization: Penalizes complex models  
 $E' = \text{error on data} + \lambda \text{ model complexity}$  (where  $\lambda$  is a penalty weight)

*the lower the better*

- Other measures of “goodness of fit” with complexity penalty:

- ▣ Akaike’s information criterion (AIC)

$$\text{AIC} \equiv \log p(X|\theta_{ML}, M) - k(M)$$

- ▣ Bayesian information criterion (BIC)

$$\text{BIC} \equiv \log p(X|\theta_{ML}, M) - k(M) \log(N)/2$$

where:  $M$  is a model

$\log p(X|\theta_{ML}, M)$ : is the *log likelihood* of  $M$  where  $M$ ’s parameters  $\theta_{ML}$  have been estimated using maximum likelihood

$k(M)$  = number of adjustable parameters in  $\theta_{ML}$  of the model  $M$

$N$  = size of sample  $X$

*For both AIC and BIC, the higher the better*

# Model Selection (2 of 2 slides)

## Methods to fine-tune model complexity

30

- Minimum description length (MDL): Kolmogorov complexity, shortest description of data

Given a dataset  $X$ ,

$\text{MDL}(M) = \text{Description length of model } M +$

Description length of data in  $X$  not correctly described by  $M$

*the lower the better*

- Structural risk minimization (SRM)

- ▣ Uses a set of models and their complexities (measured usually using their number of free parameters or their VC-dimension)
- ▣ Selects the simplest model in terms of order and best in terms of empirical error on the data

# Bayesian Model Selection

31

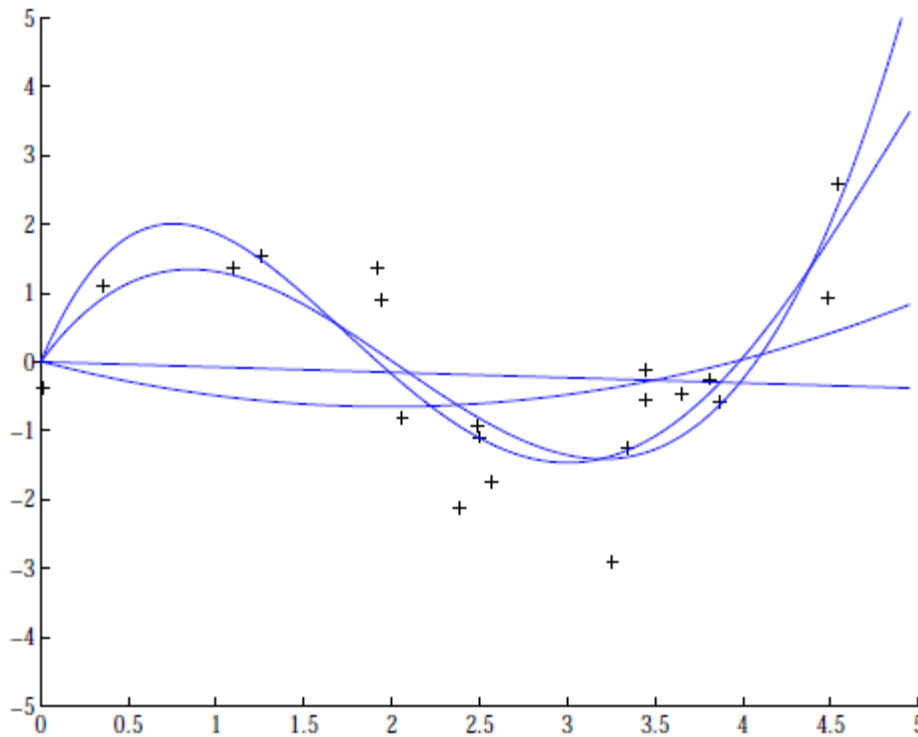
- Used when we have Prior on models,  $p(\text{model})$

$$p(\text{model} | \text{data}) = \frac{p(\text{data} | \text{model}) p(\text{model})}{p(\text{data})}$$

- Regularization, when prior favors simpler models
- Bayes, MAP of the posterior,  $p(\text{model} | \text{data})$
- Average over a number of models with high posterior (voting, ensembles: Chapter 17)

# Regression example

32



Coefficients increase in magnitude as order increases:

1: [-0.0769, 0.0016]

2: [0.1682, -0.6657, 0.0080]

3: [0.4238, -2.5778, 3.4675, -0.0002]

4: [-0.1093, 1.4356, -5.5007, 6.0454, -0.0019]

Regularization (L2): 
$$E(\mathbf{w} | \mathbf{X}) = \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t | \mathbf{w})]^2 + \lambda \sum_i w_i^2$$