

# Introduction to ML

October 4<sup>th</sup>, 2021

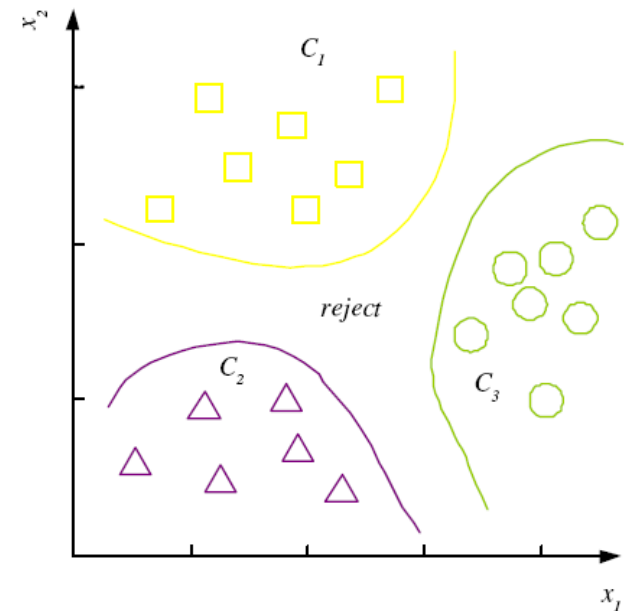
# Discriminant Functions

Classification rule: pick one such function that maximizes

choose  $C_i$  if  $g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$

$$g_i(\mathbf{x}) = \begin{cases} -R(\alpha_i | \mathbf{x}) \\ P(C_i | \mathbf{x}) \\ p(\mathbf{x} | C_i)P(C_i) \end{cases}$$

Three different discriminant functions



$K$  decision regions  $\mathcal{R}_1, \dots, \mathcal{R}_K$

$$\mathcal{R}_i = \{\mathbf{x} | g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$$

Discriminant functions:  
Divides feature space into  $K$  region  
For those inputs  $\mathbf{x}$ , find the function that give the largest value (use that function to carve out a region)

# Association Rules

- Association rule:  $X \rightarrow Y$
- *People who buy/click/visit/enjoy X are also likely to buy/click/visit/enjoy Y.*
  - Dependency, it's not cause and effect but good enough for making a business decision
- A rule implies association, not necessarily causation.

# Association measures

Joint probability,  
Make this large (increase the basis)  
**Significance** of the rule  
Useless is this number is low

- Support ( $X \rightarrow Y$ ):

$$P(X, Y) = \frac{\#\{\text{customers who bought } X \text{ and } Y\}}{\#\{\text{customers}\}}$$

- Confidence ( $X \rightarrow Y$ ):

$$P(Y | X) = \frac{P(X, Y)}{P(X)}$$

- conditional probability
- Should be larger than  $P(Y)$
- **Strength** of the association rule

- Lift ( $X \rightarrow Y$ ):

$$= \frac{\#\{\text{customers who bought } X \text{ and } Y\}}{\#\{\text{customers who bought } X\}}$$

$$= \frac{P(X, Y)}{P(X)P(Y)} = \frac{P(Y | X)}{P(Y)}$$

- If independent, this is 1
- If lift > 1 X makes Y more likely

# Example

| Transaction | Items in basket          |
|-------------|--------------------------|
| 1           | milk, bananas, chocolate |
| 2           | milk, chocolate          |
| 3           | milk, bananas            |
| 4           | chocolate                |
| 5           | chocolate                |
| 6           | milk, chocolate          |

SOLUTION:

milk  $\rightarrow$  bananas : Support = 2/6, Confidence = 2/4

bananas  $\rightarrow$  milk : Support = 2/6, Confidence = 2/2

milk  $\rightarrow$  chocolate : Support = 3/6, Confidence = 3/4

chocolate  $\rightarrow$  milk : Support = 3/6, Confidence = 3/5

In making a decision, take ‘high support + high confidence rule’

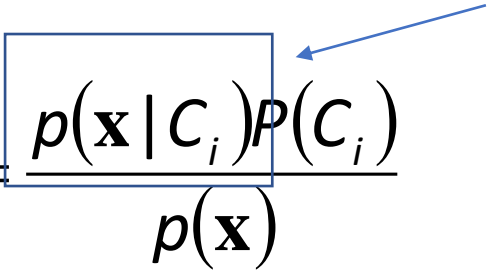
CHAPTER 4:

# Parametric Methods

Statistics based ML method

# Why do we need this?

## Using Bayes Decision Theory for Classification

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


The distribution function of  $P(\mathbf{x})$  when  $\mathbf{x}$  is coming from a class  $C_i$

Need this probability to be used as discriminant function

Collect data  
Use training set  
Learn the discriminant function (essentially estimate the parameters)


$$P(C_i) \geq 0 \text{ and } \sum_{i=1}^K P(C_i) = 1$$

choose  $C_i$  if  $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

# Parametric Estimation

- $\mathcal{X} = \{x^t\}_t$  where  $x^t \sim p(x)$
- Parametric estimation:  
Assume a form for  $p(x | \theta)$  and estimate  $\theta$ , its sufficient statistics, using  $X$   
e.g.,  $N(\mu, \sigma^2)$  where  $\theta = \{\mu, \sigma^2\}$

How about Bernoulli? Binomial?



All that is required  
to define a  
probability  
distribution



# Maximum Likelihood Estimator

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

$$l(\vartheta | \mathcal{X}) = p(\mathcal{X} | \vartheta) = \prod_t p(x^t | \vartheta)$$

Estimate parameter  
of the model using  $\mathcal{X}$   
with this criterion

- Log likelihood

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

Joint factors to product  
(assume iid samples)

- Maximum likelihood estimator (MLE)

$$\vartheta^* = \operatorname{argmax}_{\vartheta} \mathcal{L}(\vartheta | \mathcal{X})$$

# Examples: Bernoulli/Multinomial

- **Bernoulli:** Two states, failure/success,  $x$  in  $\{0,1\}$

$$P(x) = p_o^x (1 - p_o)^{(1-x)}$$

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

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

- **Multinomial:**  $K > 2$  states,  $x_i$  in  $\{0,1\}$

$$P(x_1, x_2, \dots, x_K) = \prod_i p_i^{x_i}$$

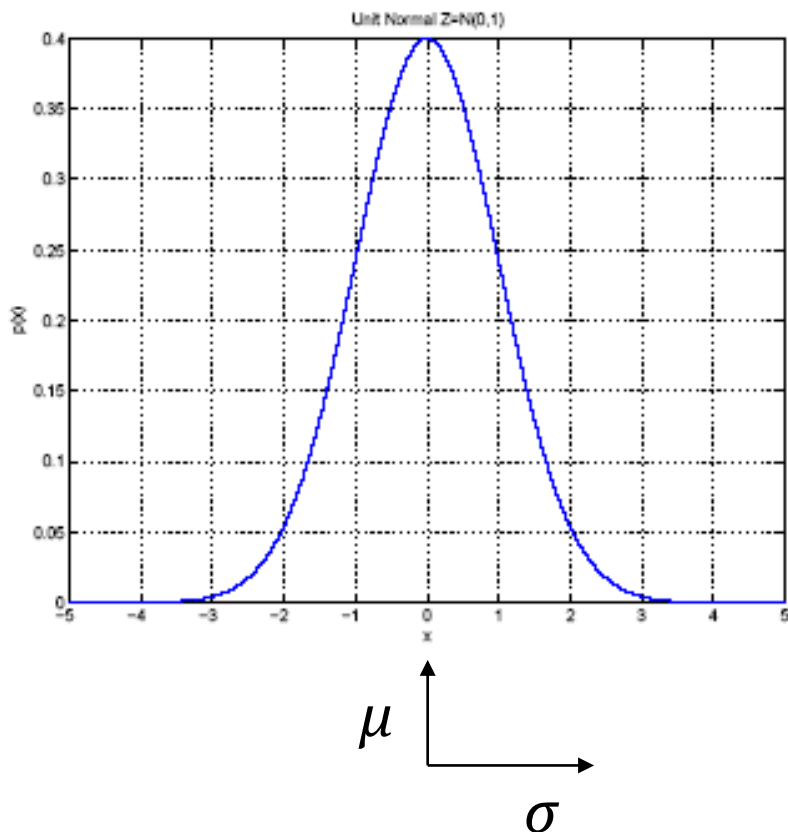
$$\mathcal{L}(p_1, p_2, \dots, p_K | \mathcal{X}) = \log \prod_t \prod_i p_i^{x_i^t}$$

$$\text{MLE: } p_i = \sum_t x_i^t / N$$

Like  
average

# Gaussian (Normal) Distribution

- Used for continuous-value data



- $p(x) = \mathcal{N}(\mu, \sigma^2)$

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

- 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

We derive a function to estimate the parameter values of distribution

They are a function of rvs, also a rv by itself, understand the property of such a learner is important

Unknown parameter  $\theta$

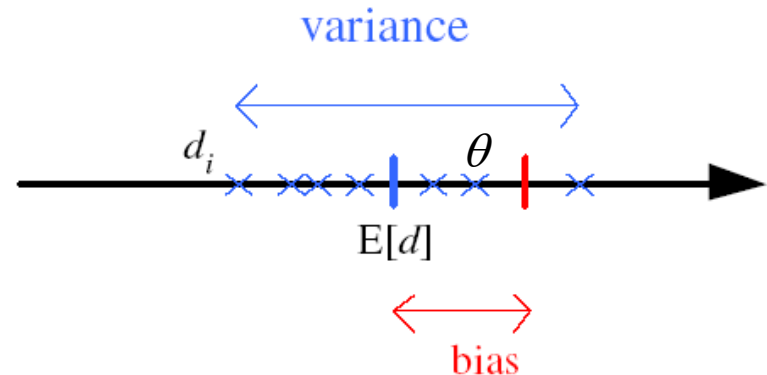
Estimator  $d_i = d(X_i)$  on sample  $X_i$

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

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

Mean square error:

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



Errors in our estimator consists of two things:

1. Bias<sup>2</sup>
2. Variance

Take home thinking: what would minimize this?

# Bayes' Estimator

-different from MLE

Parameter of the function  
is a random variable with a  
distribution

- Treat  $\vartheta$  as a random var with prior  $p(\vartheta)$
- Bayes' rule:  $p(\vartheta|\mathcal{X}) = p(\mathcal{X}|\vartheta) p(\vartheta) / p(\mathcal{X})$
- Full:  $p(x|\mathcal{X}) = \int p(x|\vartheta) p(\vartheta|\mathcal{X}) d\vartheta$
- Maximum a Posteriori (MAP) estimation:

$$\vartheta_{\text{MAP}} = \operatorname{argmax}_{\vartheta} p(\vartheta|\mathcal{X})$$

- Maximum Likelihood (ML) estimation:  $\vartheta_{\text{ML}} = \operatorname{argmax}_{\vartheta} p(\mathcal{X}|\vartheta)$

- Bayes':  $\vartheta_{\text{Bayes'}} = E[\vartheta|\mathcal{X}] = \int \vartheta p(\vartheta|\mathcal{X}) d\vartheta$



Find the **conditional expected** value of the parameter  
random variable given collected data

# Bayes' Estimator: Example

- $x^t \sim \mathcal{N}(\vartheta, \sigma_o^2)$  and  $\vartheta \sim \mathcal{N}(\mu, \sigma^2)$
- $\vartheta_{\text{ML}} = m$
- $\vartheta_{\text{MAP}} = \vartheta_{\text{Bayes'}} =$

$$E[\theta | \mathcal{X}] = \frac{N/\sigma_0^2}{N/\sigma_0^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma_0^2 + 1/\sigma^2} \mu$$

Weighted average between mean of the prior and the mean of the data distribution

What happens when N is very large?  
(you learn everything from data)

# Parametric Classification


-once done (learning), you also obtain the actual 'discriminant' function that can be used for classification

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

or

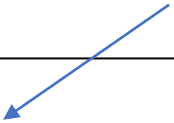
$$g_i(x) = \log p(x | C_i) + \log P(C_i)$$

This is your discriminant function for classification



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

The actual implementation when using Gaussian distribution



An example of learning  
for two class problem

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

$$x \in \mathfrak{R}$$

$$r_i^t = \begin{cases} 1 & \text{if } x^t \in C_i \\ 0 & \text{if } x^t \in C_j, j \neq i \end{cases}$$

- 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

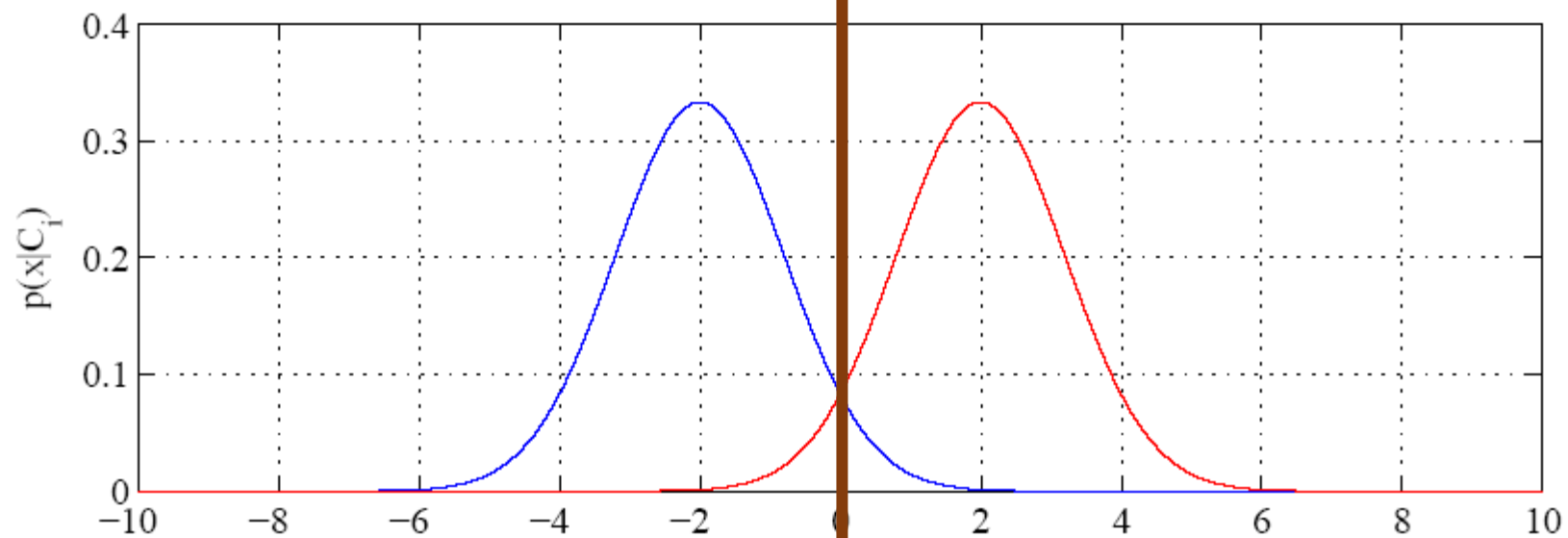
Once done, use the  
following function on  
test set

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

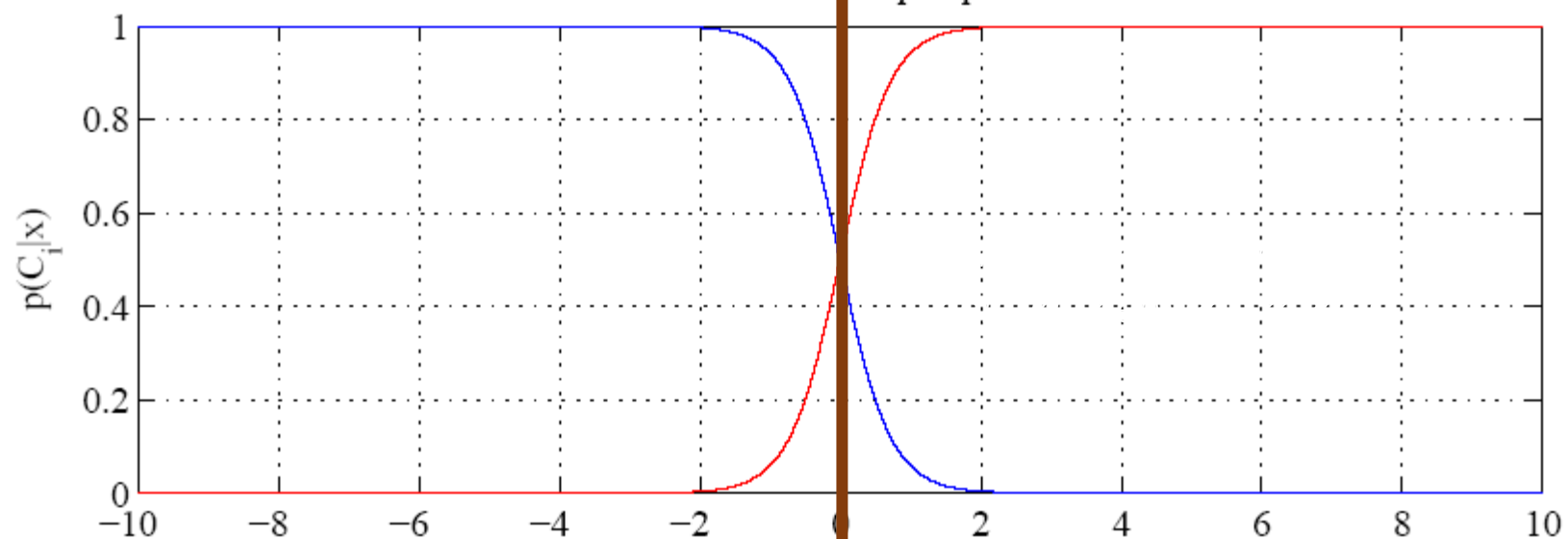
Input  $x$  (test sample) for each  $i$ , look at the class label  $i$ , for the one that is max, pick that class  $i$

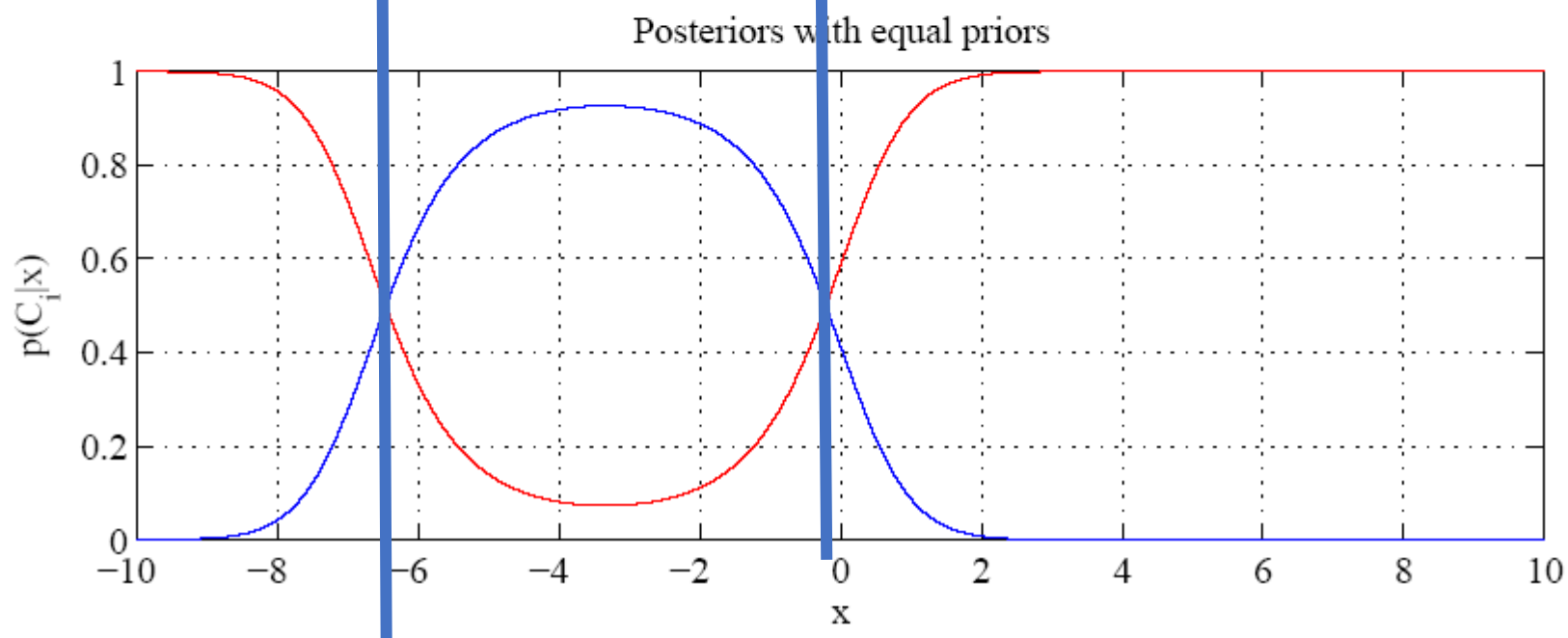
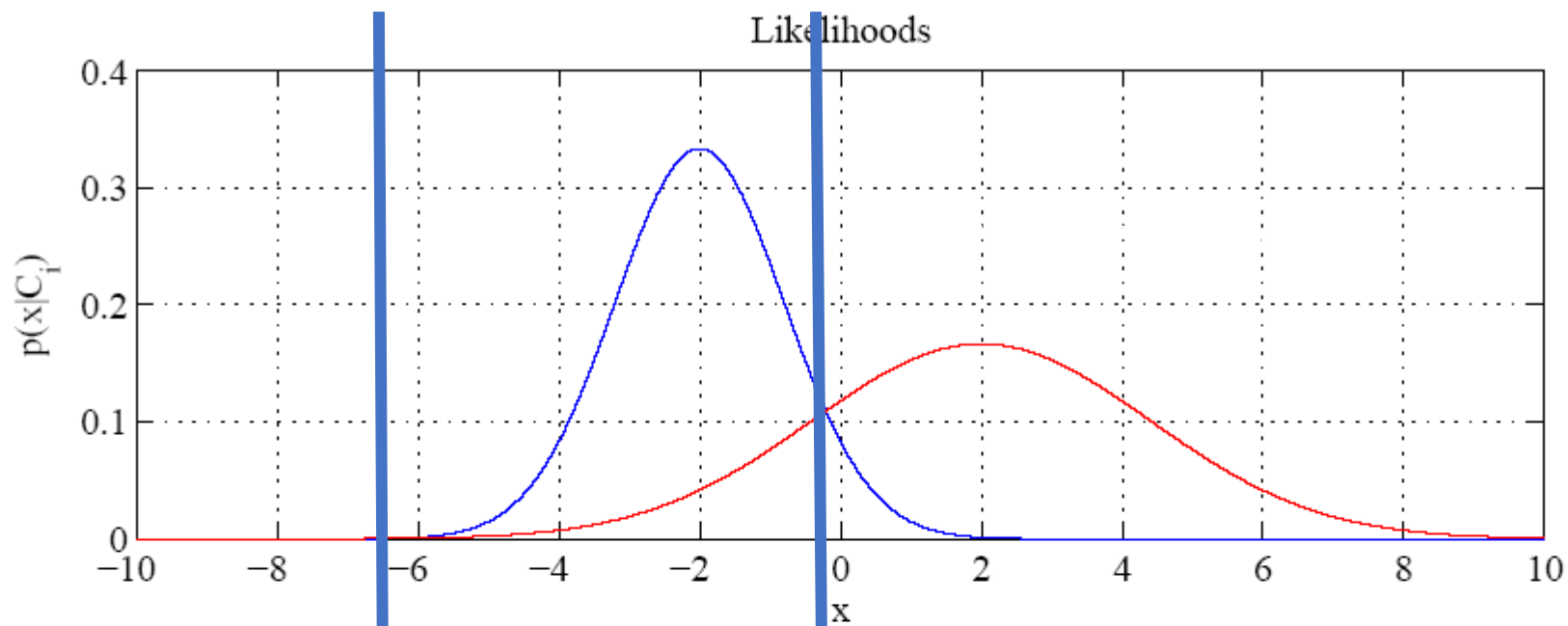


Likelihoods

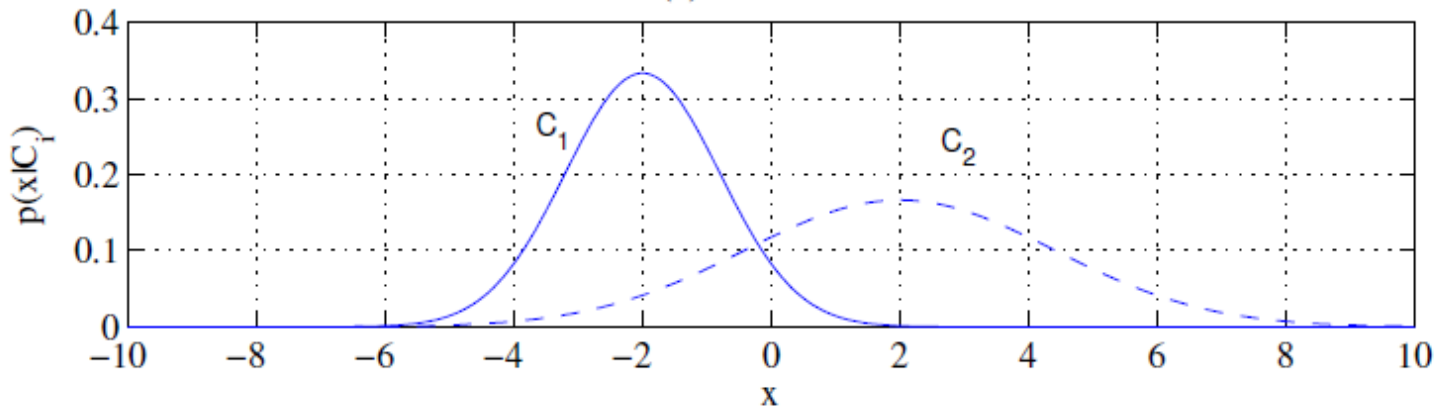


Posteriors with equal priors

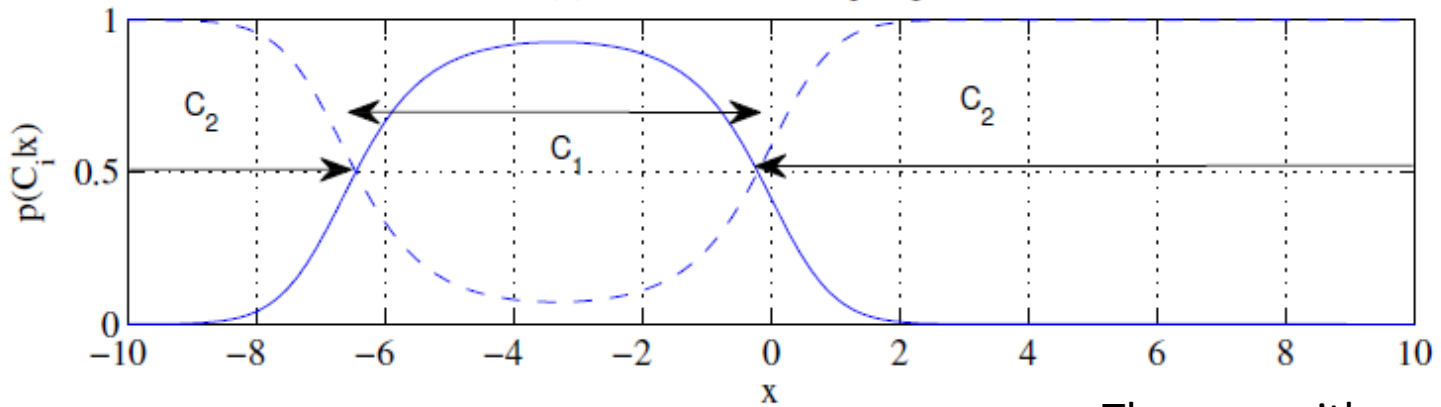




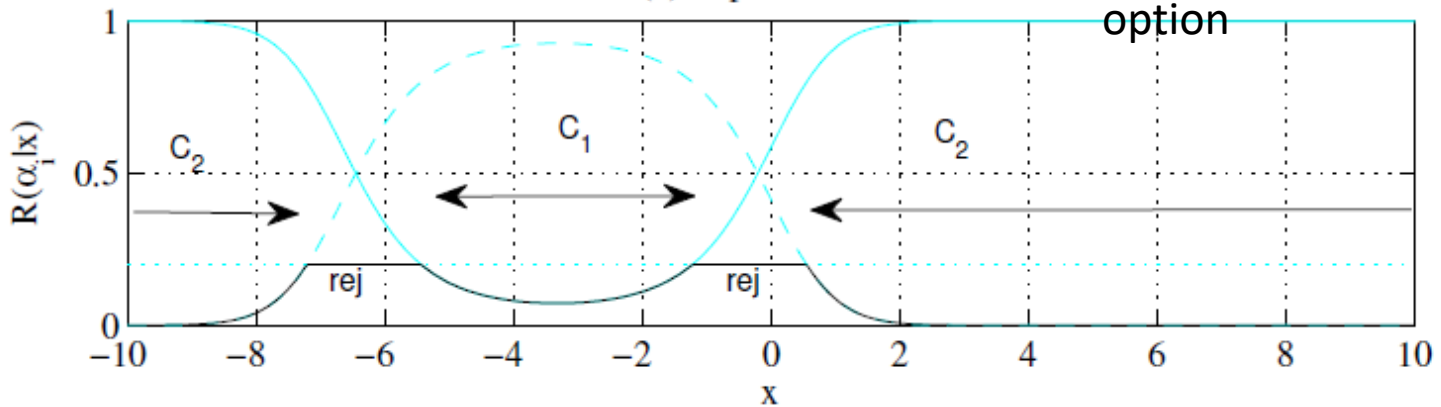
(a) Likelihoods



(b) Posteriors with equal priors



(c) Expected risks



# Regression

measurement      Input  $x$ ,  $f$  is what we want

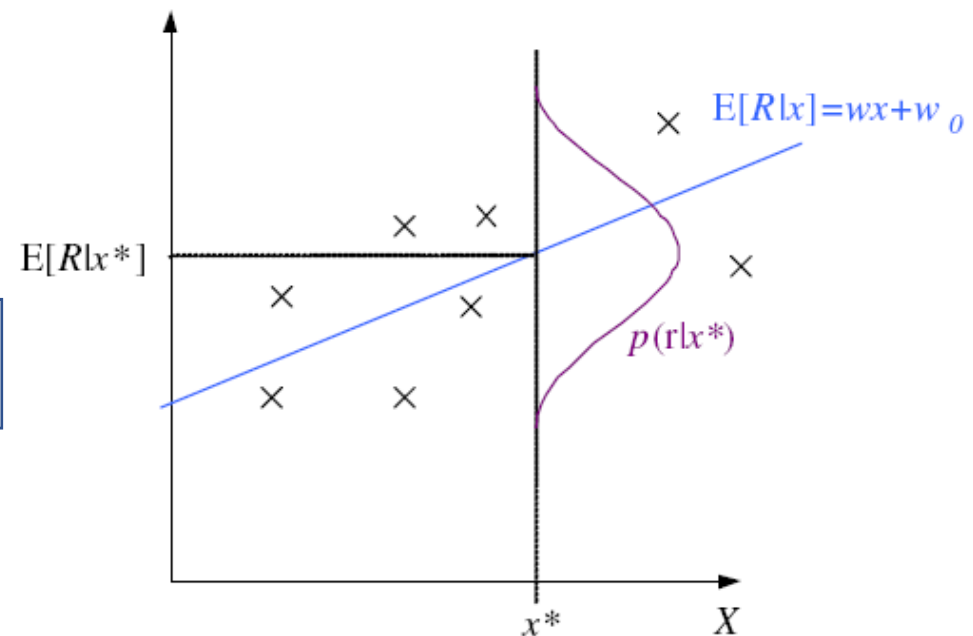
$$r = f(x) + \varepsilon$$

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

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

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

A value out of  $g()$  add to a Gaussian distribution  $\rightarrow$  a Gaussian distribution



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

Likelihood function

Joint = condition  $x$  unconditioned

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

# Regression: From LogL to Error

Expected  
value

$$\mathcal{L}(\theta | \mathcal{X}) = \log \prod_{t=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{[r^t - g(x^t | \theta)]^2}{2\sigma^2} \right]$$

$$= -N \log \sqrt{2\pi}\sigma - \frac{1}{2\sigma^2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2$$

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

Least square estimation

# Linear Regression

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

This is our loss function

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

$$\sum_t r^t = N w_0 + w_1 \sum_t x^t$$

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

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

$$\mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$$

# Polynomial Regression

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

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

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

# Other Error Measures

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

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

If it is close to 1, it is almost like guessing the average of the data at all time

- Absolute Error:  $E(\vartheta | \mathcal{X}) = \sum_t |r^t - g(x^t | \vartheta)|$   
 $E(\vartheta | \mathcal{X}) = \sum_t 1(|r^t - g(x^t | \vartheta)| > \varepsilon) (|r^t - g(x^t | \vartheta)| - \varepsilon)$

Error always decreases with more complex model  
 So how do we choose and select model?



# Bias and Variance for regression

No g, just pure noise

Difference between real output and  
our estimate output  
Quantify how well on average our  $g(x)$   
is on the training set

$$E[(r - g(x))^2 | x] = \underbrace{E[(r - E[r | x])^2 | x]}_{\text{noise}} + \underbrace{(E[r | x] - g(x))^2}_{\text{squared error}}$$

$$E_x[(E[r | x] - g(x))^2] = \underbrace{(E[r | x] - E_x[g(x)])^2}_{\text{bias}} + \underbrace{E_x[(g(x) - E_x[g(x)])^2]}_{\text{variance}}$$

To quantify how well on average our  $g(x)$  is on different dataset, we take the average over  $X$

Side note:

$$\text{Var}(x) = E(X^2) - E(X)^2$$

$$\text{Var}(x) = E((x - E(x))^2)$$

# Estimating Bias and Variance

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

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

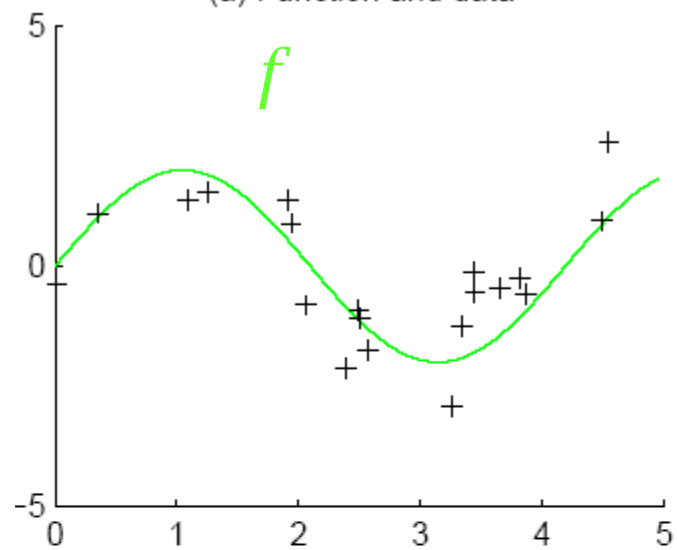
$$\text{Variance}(g) = \frac{1}{NM} \sum_t \sum_i [g_i(x^t) - \bar{g}(x^t)]^2$$

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

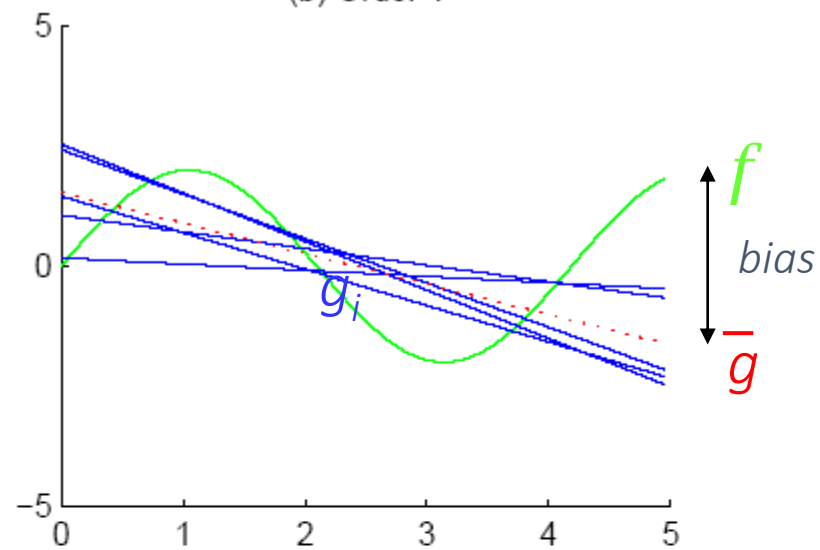
# Bias/Variance Dilemma

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

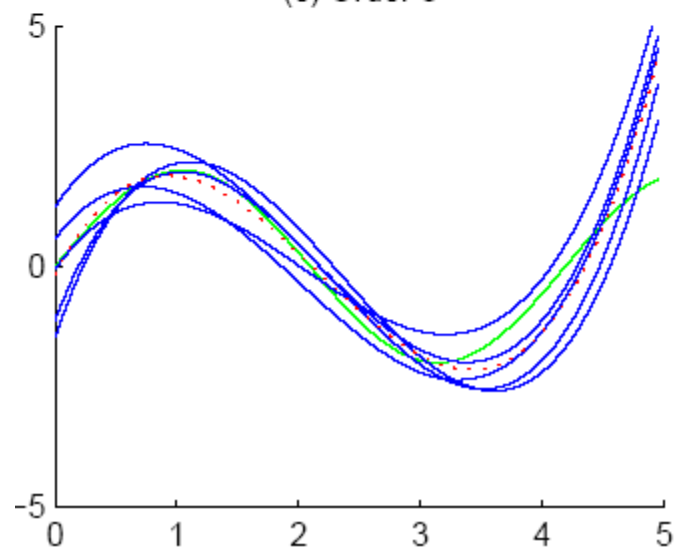
(a) Function and data



(b) Order 1



(c) Order 3



(d) Order 5

