

Parametric Models (Chpt 3-5)

Rui Kuang

Department of Computer Science and Engineering
University of Minnesota

Probabilistic Perspective

- We have seen classification models.
Classification decision is deterministic

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } h \text{ says } \mathbf{x} \text{ is positive} \\ 0 & \text{if } h \text{ says } \mathbf{x} \text{ is negative} \end{cases}$$

- What if we have cases with some uncertainty?
- Estimation of

$$p(C = 0 | \mathbf{x}) \text{ and } P(C = 1 | \mathbf{x})$$

Classification

- Credit scoring: Inputs are income and savings. Output is low-risk vs high-risk
- Input: $\mathbf{x} = [x_1, x_2]^T$, Output: C is in {0,1}
- Prediction:

$$\text{choose } \begin{cases} C = 1 & \text{if } P(C = 1 | x_1, x_2) > P(C = 0 | x_1, x_2) \\ C = 0 & \text{otherwise} \end{cases}$$

Coin Toss



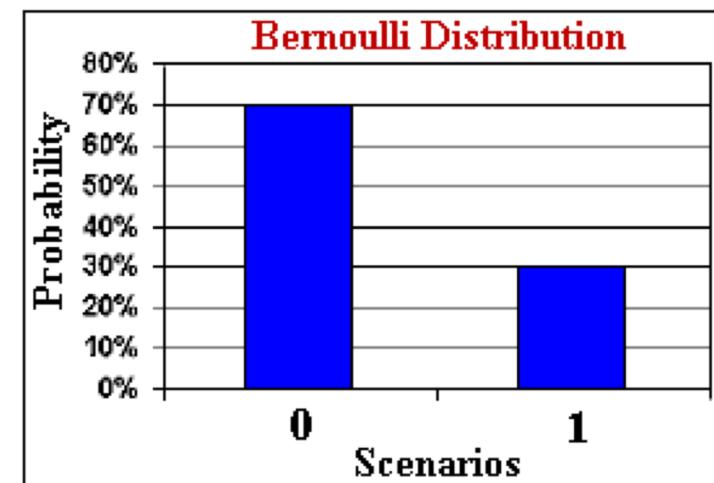
- Random var $X \in \{1, 0\}$

Bernoulli: $P\{X=1\} = p_o^X (1 - p_o)^{1-X}$

- Sample: $X = \{x^t\}_{t=1}^N = \{0, 1, \dots, 1, \dots\}$
- Prediction of next toss (no input):

Heads if $p_o > \frac{1}{2}$,

Tails otherwise



Parametric Estimation

- $\mathcal{X} = \{x^t\}$ where $x^t \sim p(x)$
- Parametric estimation:
 - Assume a form for $p(x|\theta)$ and estimate θ , its parameters, using X
 - e.g., $N(\mu, \sigma^2)$ where $\theta = \{\mu, \sigma^2\}$
 - $P(x) = p^x(1-p)^{1-x}$

Maximum Likelihood Estimation

- Likelihood of θ given the sample \mathcal{X}

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

- Log likelihood

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

- Maximum likelihood estimator (MLE)

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

Examples: Bernoulli

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

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

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

$$L(p|X) = \sum_t x^t \log p + (N - \sum_t x^t) \log(1-p)$$

$$\frac{\partial L(p|X)}{p} = \sum_t x^t / p - (N - \sum_t x^t) / (1-p) = 0$$

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

Examples: Multinomial (Categorical)

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

$$P(x_1, x_2, \dots, x_K) = \prod_{i=1 \dots K} p_i^{x_i}$$

$$\mathcal{L}(p_1, p_2, \dots, p_K | \chi) = \log \prod_{t=1 \dots N} \prod_{i=1 \dots K} p_i^{x_i^t}$$

$$L(p_1 \dots p_K | \chi) = \sum_t \sum_i x_i^t \log p_i \quad \text{with } \sum_i p_i = 1.$$

$$\frac{\partial (\sum_t \sum_i x_i^t \log p_i - \alpha (\sum_i p_i - 1))}{\partial p_i} = 0$$

$$\sum_t x_i^t / p_i - \alpha = 0$$

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

Gaussian (Normal) Distribution

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

$$L(\mu, \sigma | X) = \sum_t \log \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(x^t - \mu)^2}{2\sigma^2} \right]$$

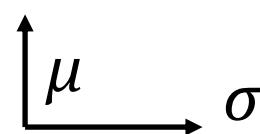
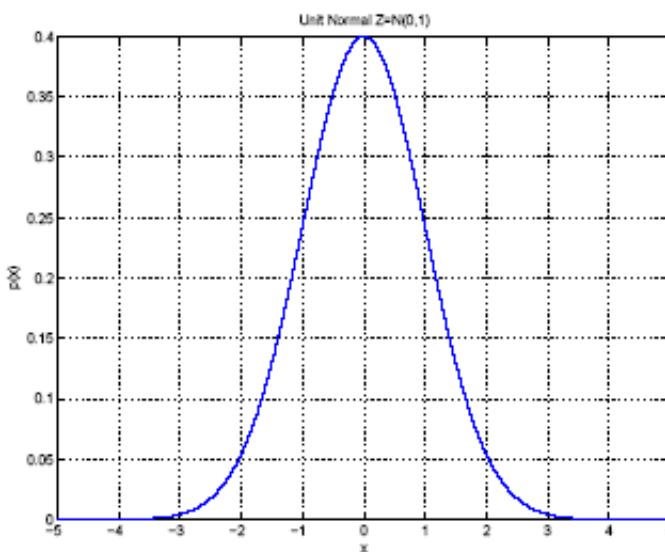
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(x - \mu)^2}{2\sigma^2} \right]$$

$$L(\mu, \sigma | X) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum_t (x^t - \mu)^2}{2\sigma^2}$$

$$\frac{\partial L(\mu, \sigma | X)}{\partial \mu} = \sum_t (x^t - \mu) = 0 \Rightarrow \mu = \frac{\sum_t x^t}{N}$$

$$\frac{\partial L(\mu, \sigma | X)}{\partial \sigma} = \frac{N}{\sigma} - \frac{\sum_t (x^t - \mu)^2}{\sigma^3} = 0$$

$$\sigma^2 = \frac{\sum_t (x^t - \mu)^2}{N}$$



Bayes' Rule

- How to get $P(C|x)$?

$$P(C|x) = \frac{prior}{evidence} \cdot likelihood$$
$$P(C|x) = \frac{P(C)p(x|C)}{p(x)}$$

prior *likelihood*
posterior ↘ ↘
 ↗

$$P(C=0) + P(C=1) = 1$$

$$p(x) = p(x|C=1)P(C=1) + p(x|C=0)P(C=0)$$

$$p(C=0|x) + P(C=1|x) = 1$$

Bayes' Rule: General Formula

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

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

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

Bayes' Rule Example

$$P(C | \mathbf{x}) = \frac{P(C) p(\mathbf{x} | C)}{p(\mathbf{x})}$$

prior *likelihood*
 ↓
↑ ← *evidence*
posterior

$$P(x=\text{high,med,low}|\text{acc}) = 0.6, 0.4, 0$$

$$P(x=\text{high,med,low}|\text{unacc}) = 0, 1/6, 5/6$$

What if we don't use Bayes' rule?

Safty (x)	Rating (C)
'high'	'acc'
'low'	'unacc'
'med'	'acc'
'high'	'acc'
'low'	'unacc'
'med'	'acc'
'high'	'acc'
'low'	'unacc'
'med'	'unacc'
'high'	'acc'
'low'	'unacc'
'med'	'acc'
'high'	'acc'
'low'	'unacc'
'med'	'acc'
'high'	'acc'

Parametric Classification

■ Discriminant function

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

or

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

■ Gaussians:

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

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

■ Given the sample

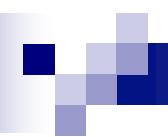
$$\mathcal{X} = \{x^t, r^t\}_{t=1}^N \quad x \in \Re \quad r_i^t = \begin{cases} 1 & \text{if } x^t \in C_i \\ 0 & \text{if } x^t \in C_j, j \neq i \end{cases}$$

■ ML estimates are

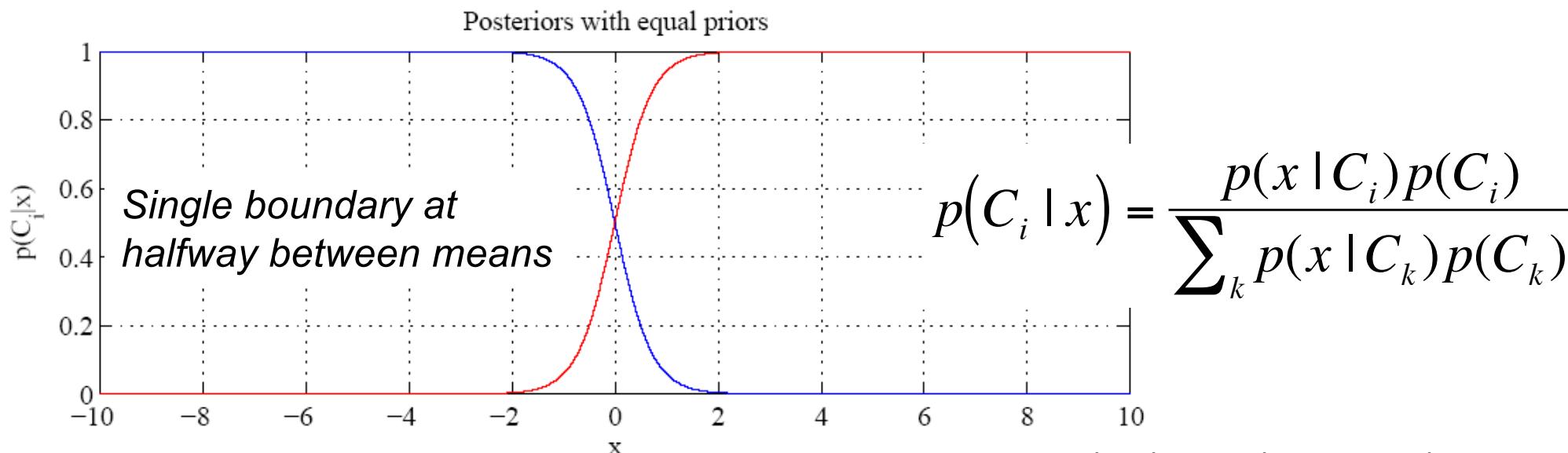
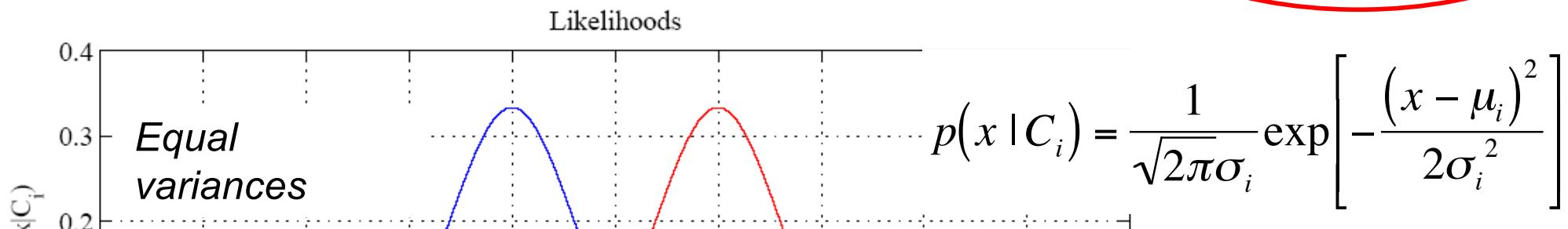
$$\hat{P}(C_i) = \frac{\sum r_i^t}{N} \quad m_i = \frac{\sum_t x^t r_i^t}{\sum_t r_i^t} \quad s_i^2 = \frac{\sum_t (x^t - m_i)^2 r_i^t}{\sum_t r_i^t}$$

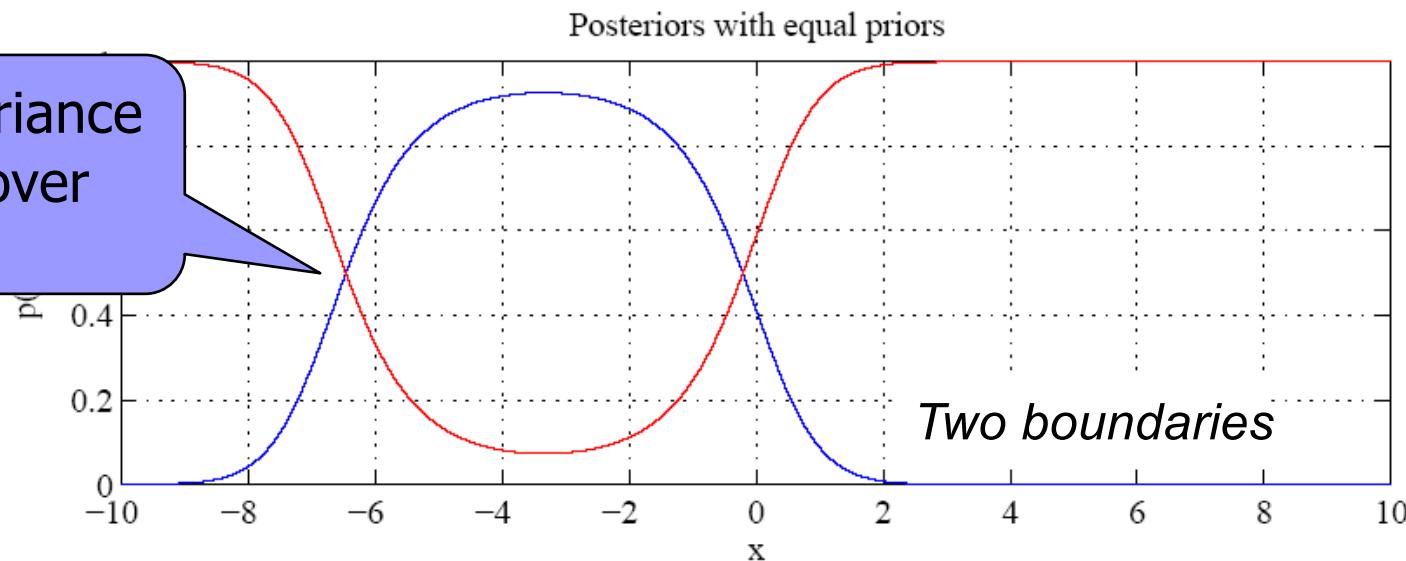
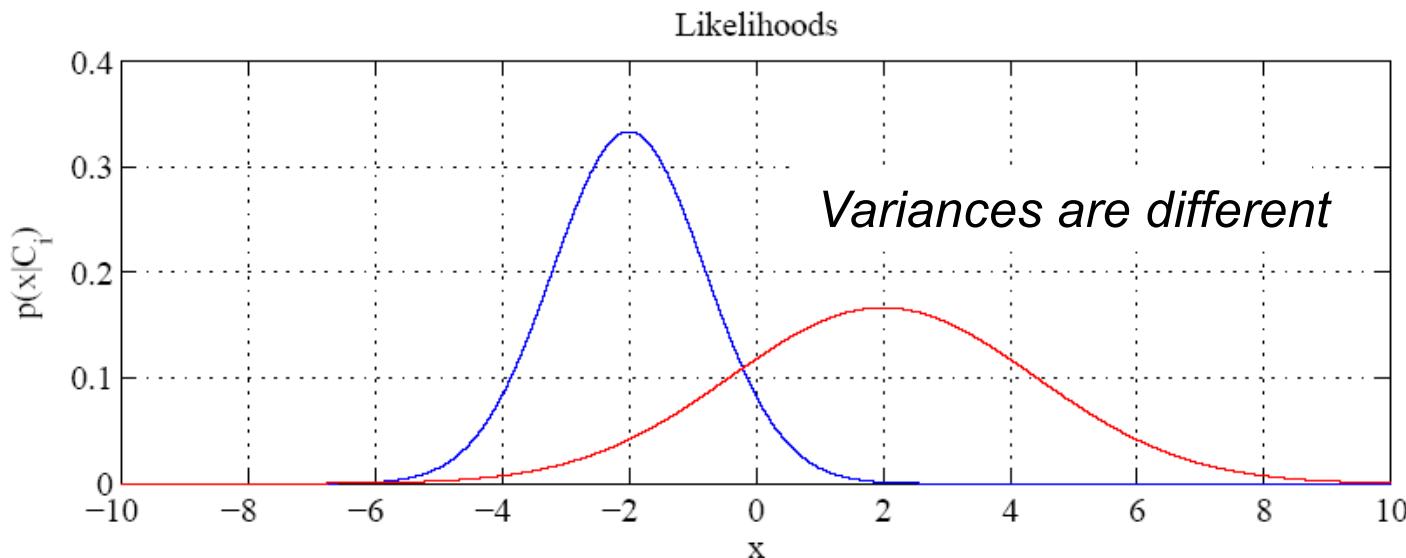
■ Discriminant becomes

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



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





Likelihood-based approach: estimate densities separately.

Multivariate Data

- Multiple measurements (sensors)
- d inputs/features/attributes: d -variate
- N instances/observations/examples

$$\mathbf{X} = \begin{bmatrix} X_1^1 & X_2^1 & \dots & X_d^1 \\ X_1^2 & X_2^2 & \dots & X_d^2 \\ \vdots \\ X_1^N & X_2^N & \dots & X_d^N \end{bmatrix}$$

Discrete Features

- **Binary** features: $p_{ij} \equiv p(x_j = 1 | C_i)$
if x_j are independent (Naive Bayes')

$$p(x | C_i) = \prod_{j=1}^d p_{ij}^{x_j} (1 - p_{ij})^{(1-x_j)}$$

Example: Text classification. Each dimension is a word, set to 1 if word in the text

	x^1	x^2	x^3	x^4
Dim1: "the" =	1	0	1	1
Dim2: "hello" =	0	1	0	1
Dim3: "and" =	1	1	0	1
Dim4: "happy" =	1	0	0	1

Discrete Features

- Likelihood of Naive Bayes': the discriminant is linear

$$\begin{aligned}g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\&= \sum_j \left[x_j \log p_{ij} + (1 - x_j) \log (1 - p_{ij}) \right] + \log P(C_i)\end{aligned}$$

Estimated parameters

$$\hat{p}_{ij} = \frac{\sum_t x_j^t r_i^t}{\sum_t r_i^t}$$

Discrete Features

- Multinomial (1-of- n_j) features: $x_j = \{v_1, v_2, \dots, v_{n_j}\}$

$$p_{ijk} \equiv p(z_{jk}=1 | C_i) = p(x_j=v_k | C_i)$$

if x_j are independent

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d \prod_{k=1}^{n_j} p_{ijk}^{z_{jk}}$$

$$g_i(\mathbf{x}) = \sum_j \sum_k z_{jk} \log p_{ijk} + \log P(C_i)$$

$$\hat{p}_{ijk} = \frac{\sum_t z_{jk}^t r_i^t}{\sum_t r_i^t}$$

In class C_i , variable x_j is category v_k

Multivariate Parameters

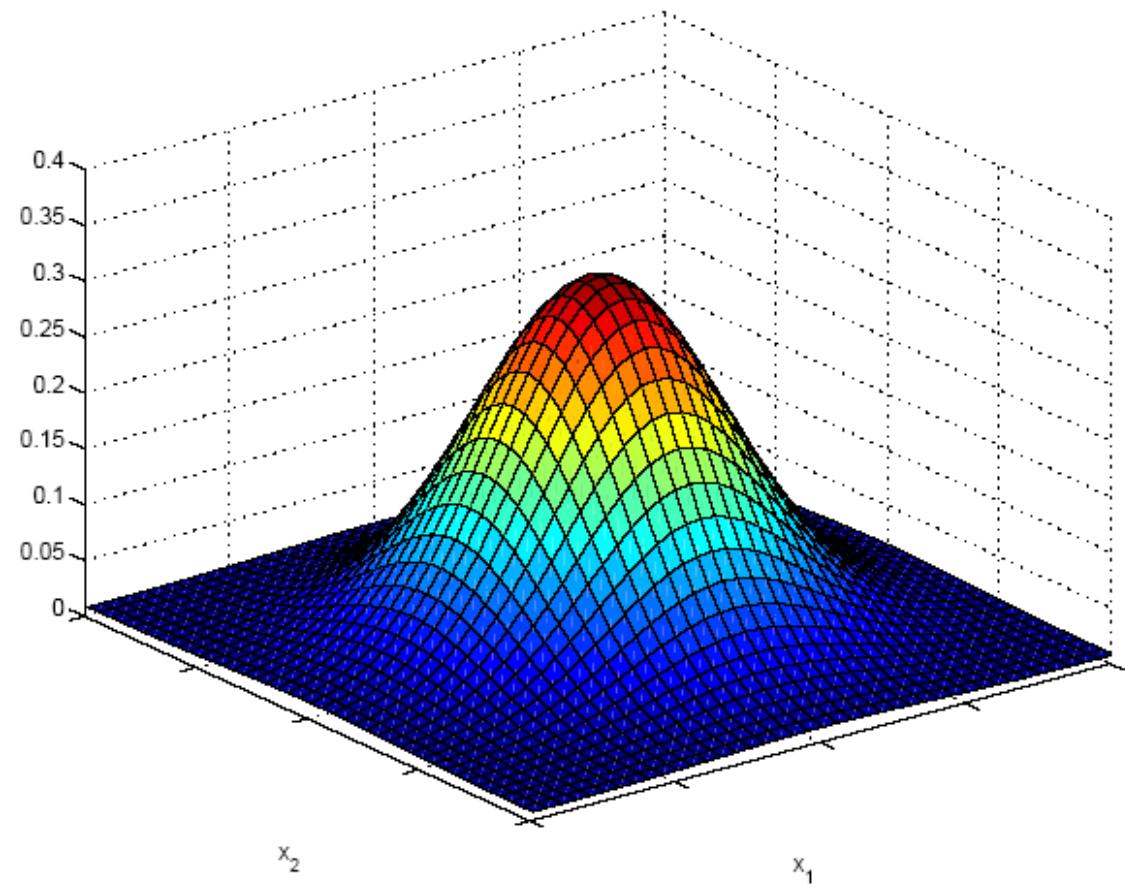
$$\text{Mean: } E[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^T$$

$$\text{Covariance: } \sigma_{ij} \equiv \text{Cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)]$$

$$\text{Correlation: } \text{Corr}(X_i, X_j) \equiv \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$$

$$\Sigma \equiv \text{Cov}(\mathbf{X}) = E\left[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\right] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & & & \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{bmatrix}$$

Multivariate Normal Distribution



$$\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$$

Independent Inputs: Naive Bayes

- If x_i are independent, offdiagonals of Σ are 0, Mahalanobis distance reduces to weighted (by $1/\sigma_i$) Euclidean distance:

$$p(\mathbf{x}) = \prod_{i=1}^d p_i(x_i) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp \left[-\frac{1}{2} \sum_{i=1}^d \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2 \right]$$

- If variances are also equal, reduces to Euclidean distance

Multivariate Normal Distribution

- Mahalanobis distance: $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ measures distance from \mathbf{x} to $\boldsymbol{\mu}$ by rotation and normalization with $\boldsymbol{\Sigma}$
 - normalizes for difference in variances and correlations
 - Variable with larger variance receive less weight
 - Two highly correlated variables contribute less.

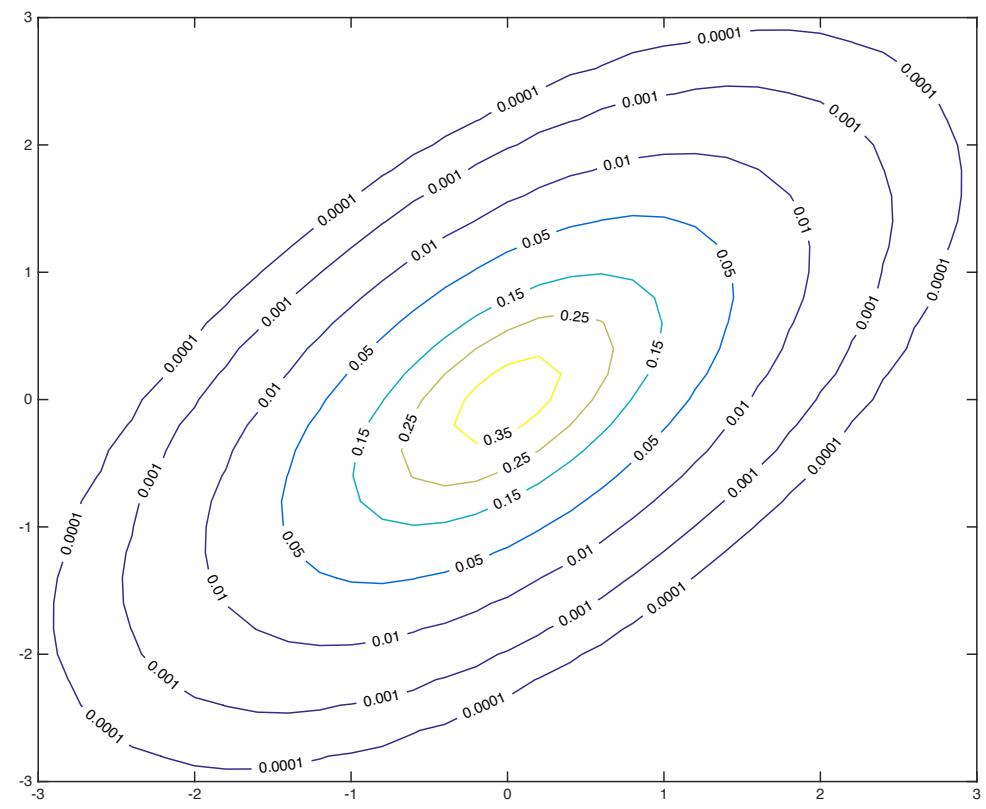
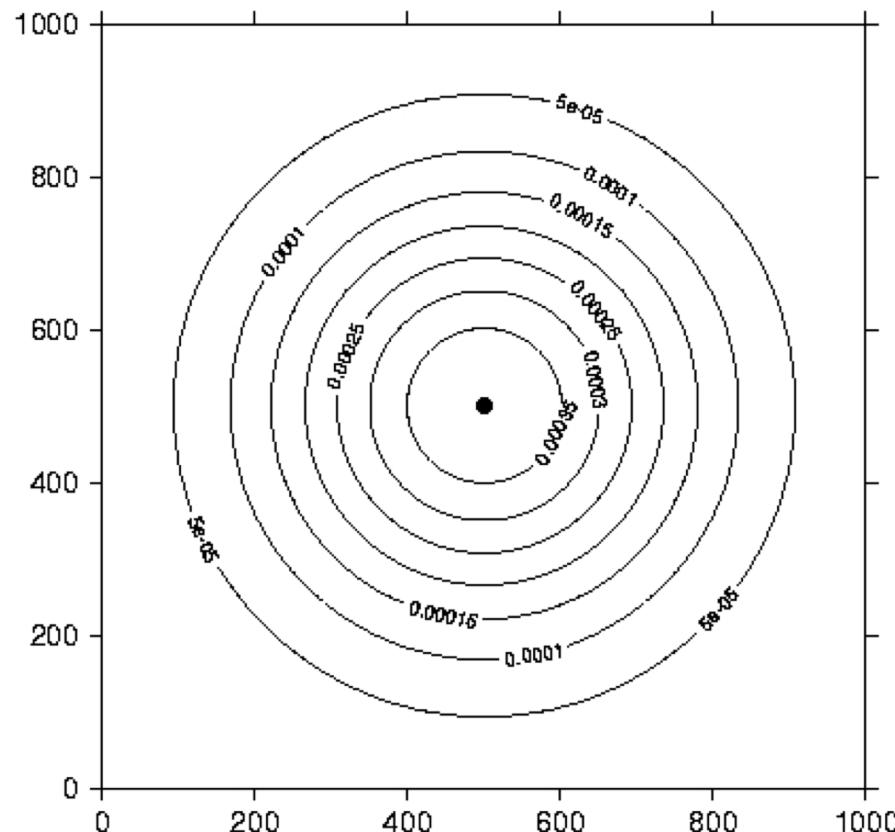
Multivariate Normal Distribution

- Covariances are obtained by rotation matrix P

$$(x - \mu)^T \Sigma^{-1} (x - \mu)$$

$$= (x - \mu)^T P^T \begin{pmatrix} 1/\sigma_1^2 & & & \\ & 1/\sigma_2^2 & & \\ & & \ddots & \\ & & \cdots & \\ & & & 1/\sigma_d^2 \end{pmatrix} P (x - \mu)$$

Interpreting Probability Density Contour



All the points on a contour curve has the same Mahalanobis distance to the mean.

Multivariate Normal Distribution

- Bivariate: $d = 2$

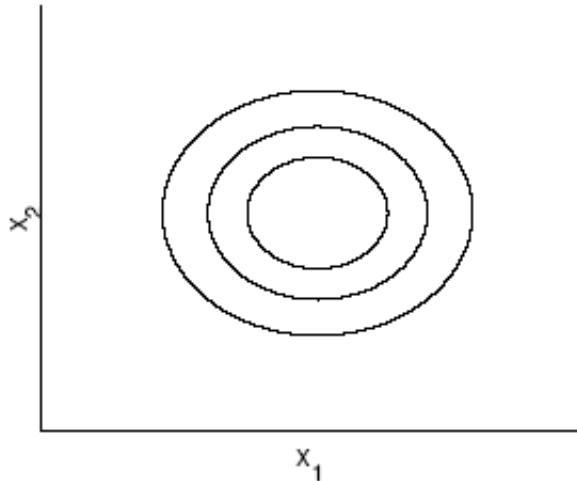
$$\Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}(z_1^2 - 2\rho z_1 z_2 + z_2^2)\right]$$

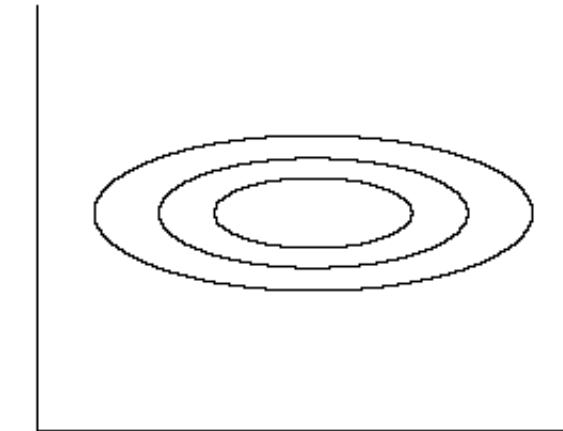
$$z_i = (x_i - \mu_i)/\sigma_i$$

Bivariate Normal

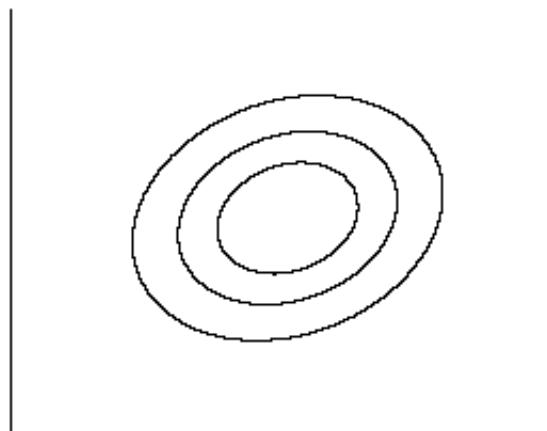
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$$



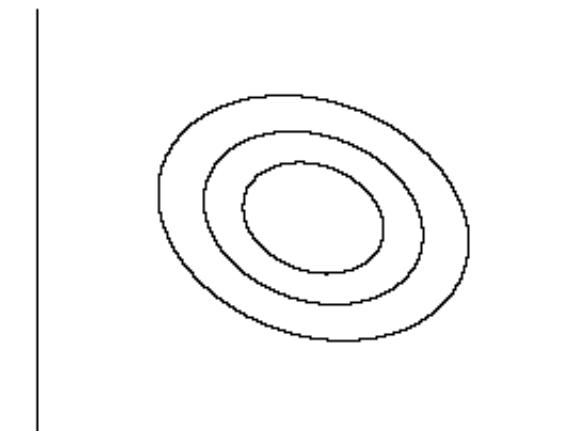
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$$

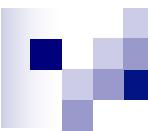


$$\text{Cov}(x_1, x_2) > 0$$

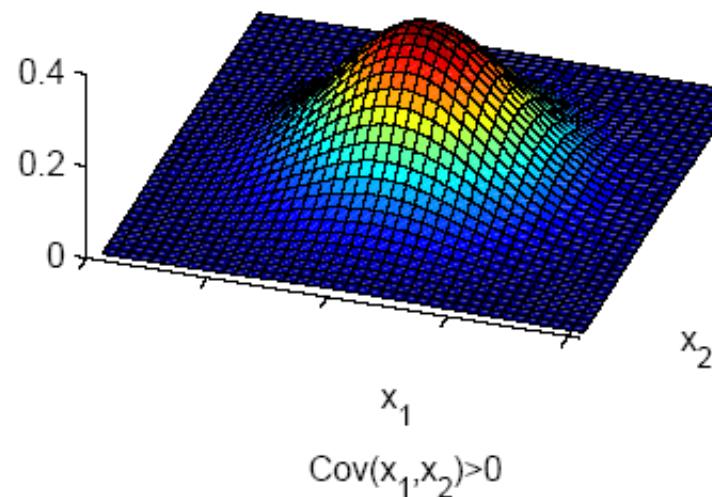


$$\text{Cov}(x_1, x_2) < 0$$

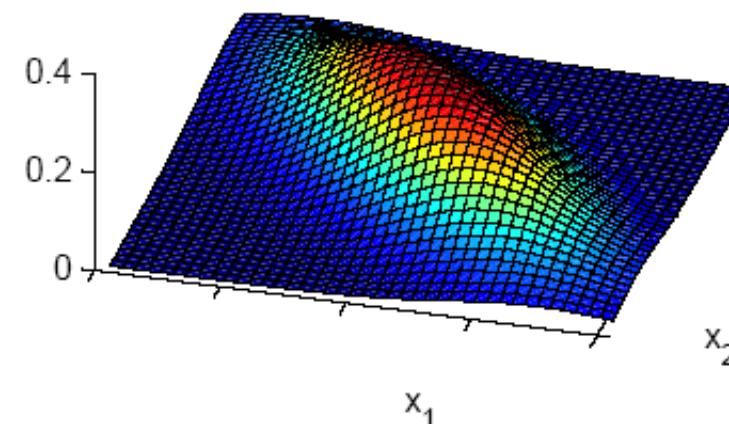
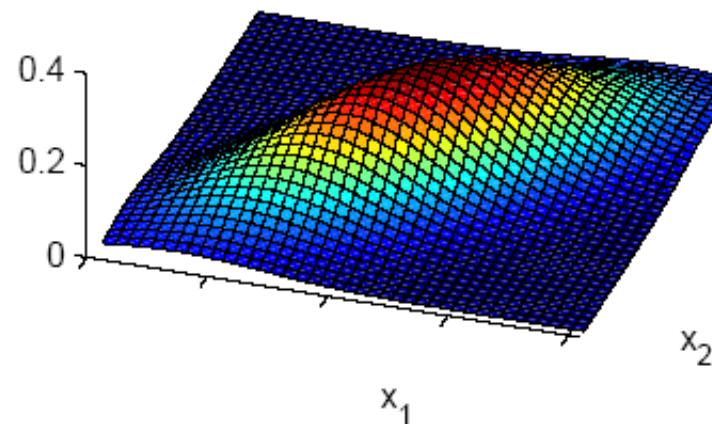
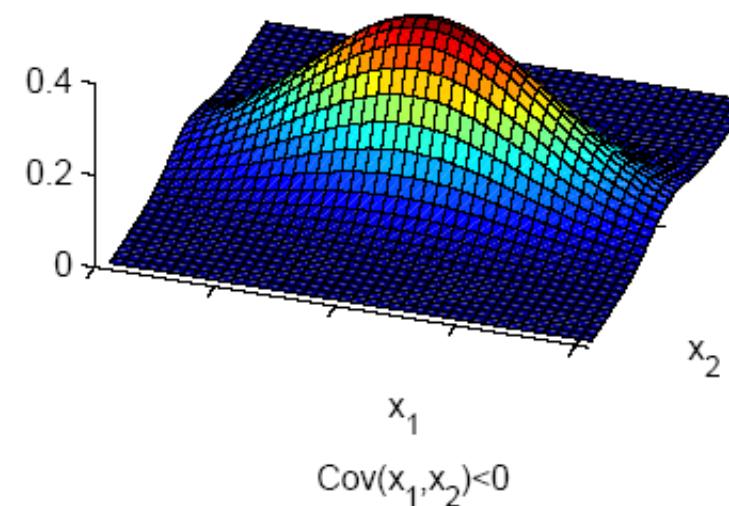




$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$



$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$



Max Likelihood for Multivariate Gaussian

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$L(\boldsymbol{\mu}, \Sigma | \chi) = \sum_{t=1}^N -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}^t - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}^t - \boldsymbol{\mu})$$

$$\text{max over } \boldsymbol{\mu}: \frac{\partial L(\boldsymbol{\mu}, \Sigma | \chi)}{\partial \boldsymbol{\mu}} = \frac{\partial \sum_{t=1}^N -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}^t - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}^t - \boldsymbol{\mu})}{\partial \boldsymbol{\mu}}$$

$$\sum_{t=1}^N (\mathbf{x}^t - \boldsymbol{\mu})^T \Sigma^{-1} = 0$$

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{t=1}^N \mathbf{x}^t$$

Max Likelihood for Multivariate Gaussian

$$\max \text{ over } \Sigma : \frac{\partial L(\mu, \Sigma | \chi)}{\partial \Sigma^{-1}} = \frac{\partial \sum_{t=1}^N -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}^t - \mu)^T \Sigma^{-1} (\mathbf{x}^t - \mu)}{\partial \Sigma^{-1}}$$
$$= \frac{\partial \sum_{t=1}^N -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} \text{trace}[(\mathbf{x}^t - \mu)(\mathbf{x}^t - \mu)^T \Sigma^{-1}]}{\partial \Sigma^{-1}}$$
$$\frac{\partial \ln |A^{-1}|}{\partial A} = -(A^{-1})^T \quad \frac{\partial \text{trace}[BA]}{\partial A} = B^T$$
$$= \frac{N}{2} \Sigma - \frac{1}{2} \sum_{t=1}^N (\mathbf{x}^t - \mu)(\mathbf{x}^t - \mu)^T$$
$$\Sigma = \frac{\sum_{t=1}^N (\mathbf{x}^t - \mu)(\mathbf{x}^t - \mu)^T}{N}$$

Parametric Classification

- If $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \Sigma_i)$

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right]$$

- Discriminant functions

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

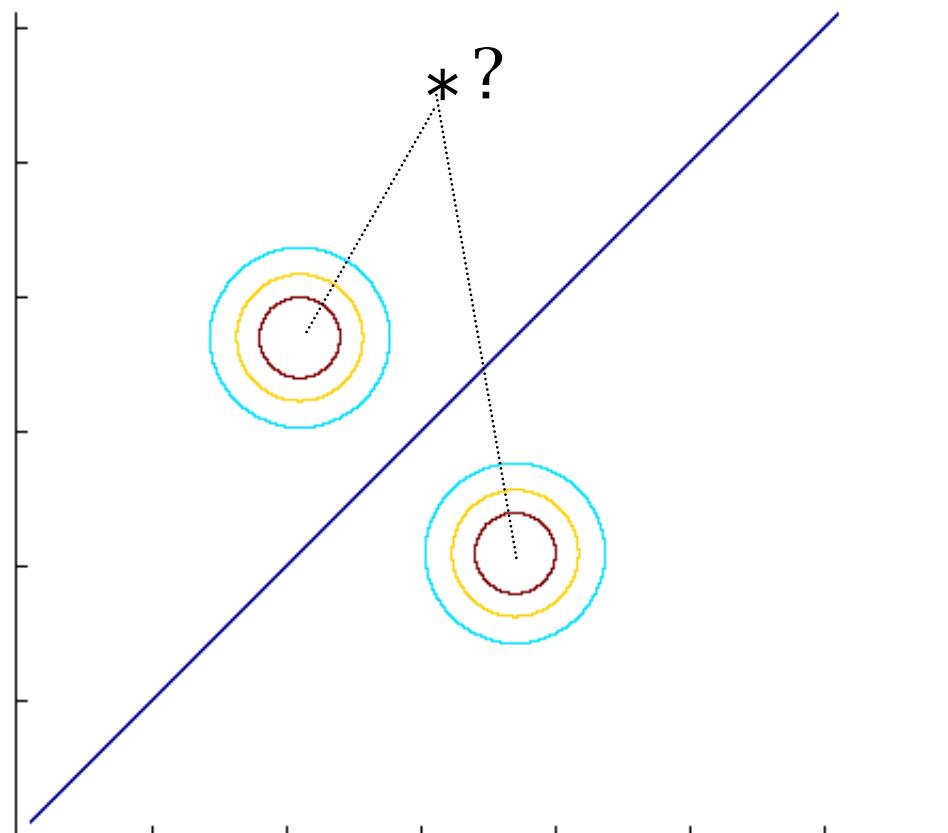
$$= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log P(C_i)$$

Model Selection

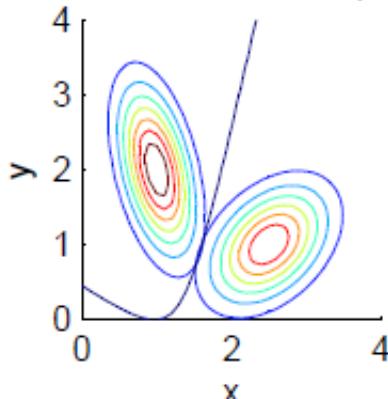
<i>Assumption</i>	<i>Covariance matrix</i>	<i>No of parameters</i>
Shared, Hyperspheric	$\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$	1
Shared, Axis-aligned	$\mathbf{S}_i = \mathbf{S}$, with $s_{ij} = 0$	d
Shared, Hyperellipsoidal	$\mathbf{S}_i = \mathbf{S}$	$d(d+1)/2$
Different, Hyperellipsoidal	\mathbf{S}_i	$K d(d+1)/2$

- As we increase complexity (less restricted \mathbf{S}), bias decreases and variance increases
- Assume simple models (allow some bias) to control variance (regularization)

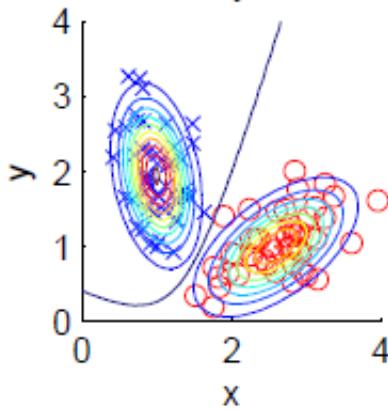
Diagonal \mathbf{S} , equal variances



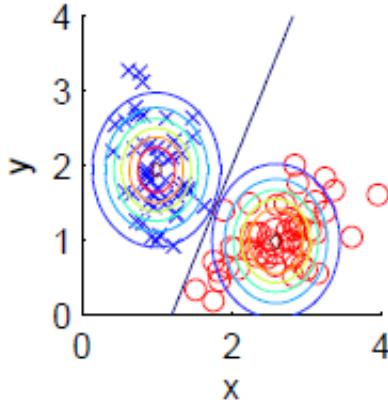
Population likelihoods and posteriors



Arbitrary covar.

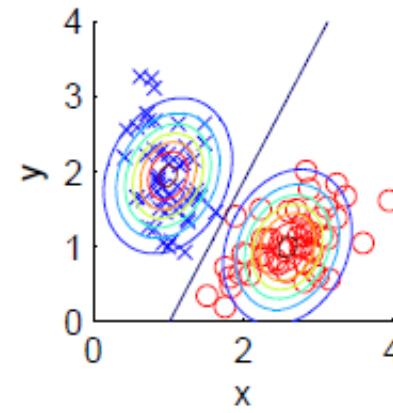


Diag. covar.

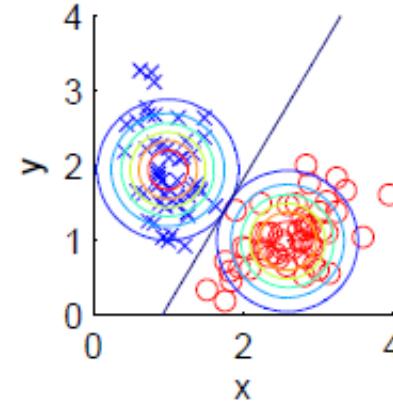


x

Shared covar.



Equal var.



x

Regression

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

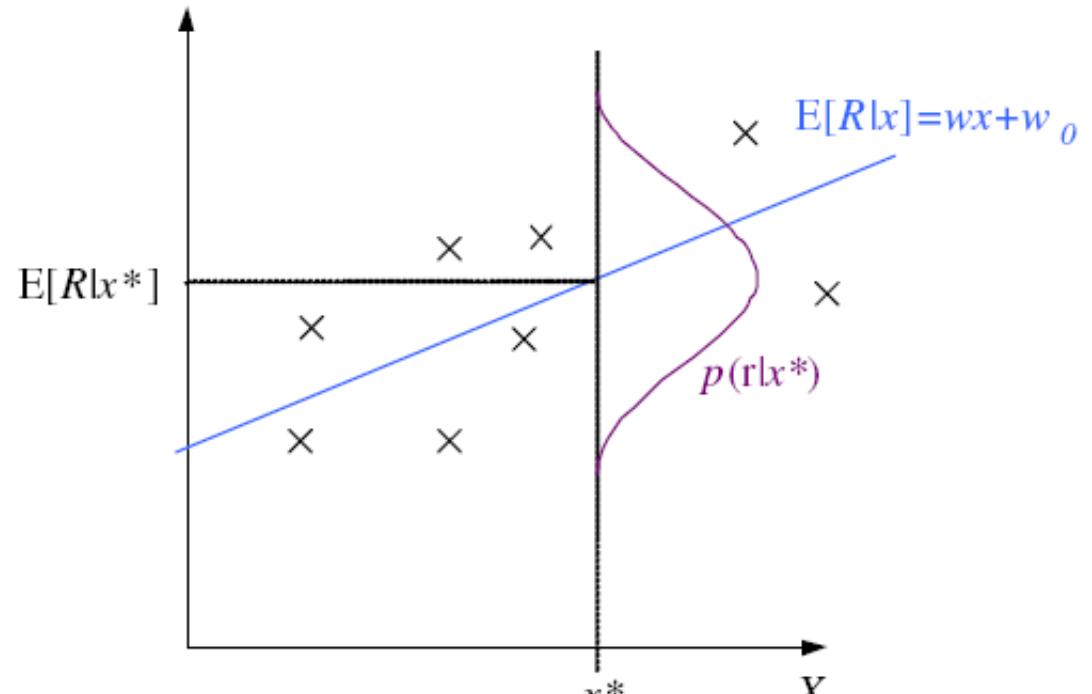
estimator : $g(x | \theta)$

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

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

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

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



Regression: From LogL to Error

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

- Maximize the log likelihood is the same as minimize the error function.

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

$$E(w_1, w_0 | \mathcal{X}) = \frac{1}{2} \sum_{t=1}^N [r^t - (w_1 x^t + w_0)]^2$$

$$\frac{\partial E(w_1, w_0 | \mathcal{X})}{\partial w_0} = \sum_{t=1}^N [(r^t - w_1 x^t - w_0)(-1)] = 0 \Leftrightarrow \sum_t r^t = N w_0 + w_1 \sum_t x^t$$

$$\frac{\partial E(w_1, w_0 | \mathcal{X})}{\partial w_1} = \sum_{t=1}^N [(r^t - w_1 x^t - w_0)(-x^t)] = 0 \Leftrightarrow \sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2$$

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix} \quad \mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$$

Multivariate Regression

$$r^t = g(x^t | w_0, w_1, \dots, w_d) + \varepsilon$$

■ Multivariate linear model

$$w_0 + w_1 x_1^t + w_2 x_2^t + \cdots + w_d x_d^t$$

$$E(w_0, w_1, \dots, w_d | \mathcal{X}) = \frac{1}{2} \sum_t [r^t - w_0 - w_1 x_1^t - \cdots - w_d x_d^t]^2$$

Multivariate Regression

$$E(w_0, w_1, \dots, w_d | \mathcal{X}) = \frac{1}{2} \sum_t [r^t - w_0 - w_1 x_1^t - \dots - w_d x_d^t]^2$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}, X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \cdots & x_d^1 \\ 1 & x_1^2 & x_2^2 & \cdots & x_d^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^N & x_2^N & \cdots & x_d^N \end{bmatrix}$$

$$\min_w \|\mathbf{r} - Xw\|^2$$

$$\frac{\partial \|\mathbf{r} - Xw\|^2}{\partial w} = -2X^T(\mathbf{r} - Xw) = 0 \Rightarrow w = (X^T X)^{-1} X^T \mathbf{r}$$

Multivariate Regression

$$E(w_0, w_1, \dots, w_d | \mathcal{X}) = \frac{1}{2} \sum_t [r^t - w_0 - w_1 x_1^t - \dots - w_d x_d^t]^2$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}, X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \cdots & x_d^1 \\ 1 & x_1^2 & x_2^2 & \cdots & x_d^2 \\ \vdots & & & & \\ 1 & x_1^N & x_2^N & \cdots & x_d^N \end{bmatrix}$$

$$\min_w \|\mathbf{r} - Xw\|^2$$

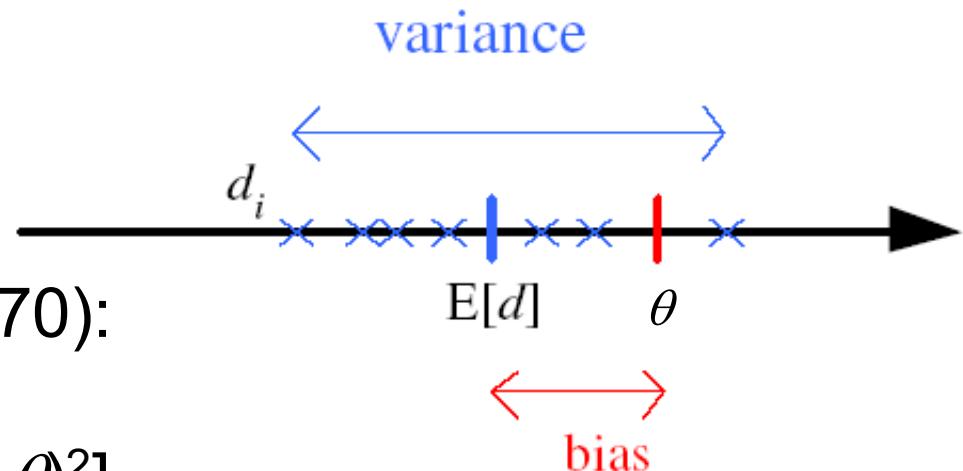
$$\frac{\partial \|\mathbf{r} - Xw\|^2}{\partial w} = -2X^T(\mathbf{r} - Xw) = 0 \Rightarrow w = (X^T X)^{-1} X^T \mathbf{r}$$

Evaluating an Estimator: Bias and Variance

Unknown parameter θ , 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 (page 70):

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

Bias/Variance Dilemma

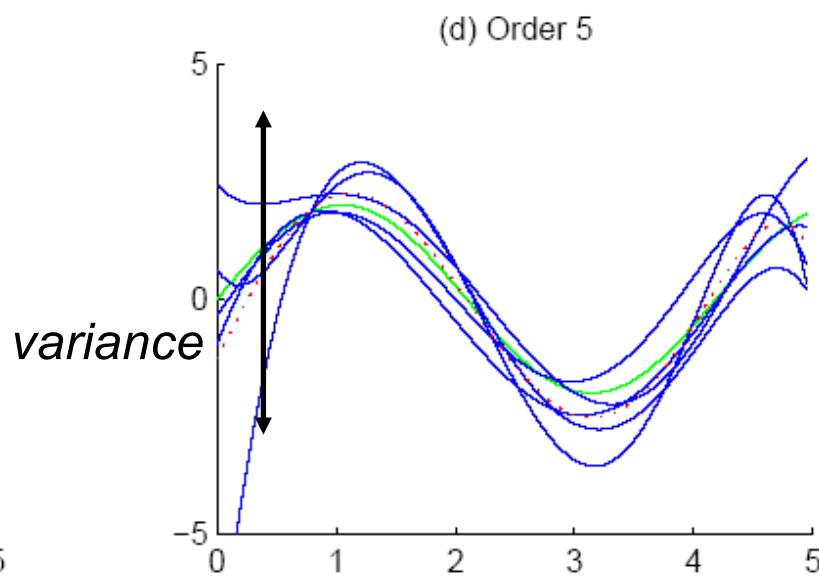
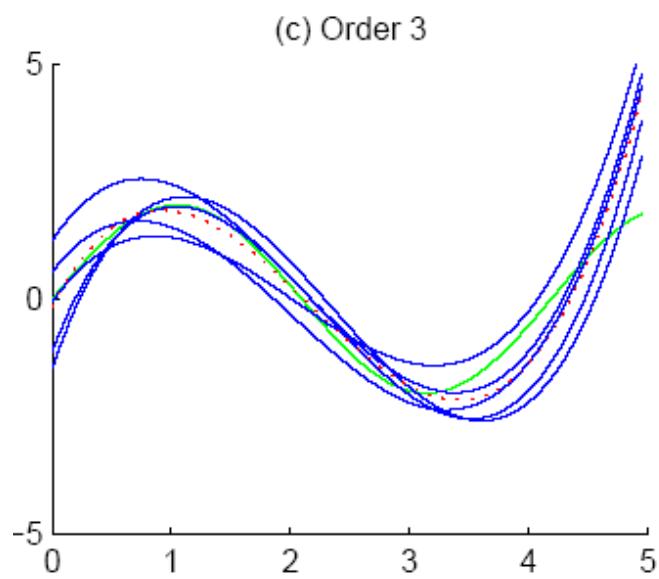
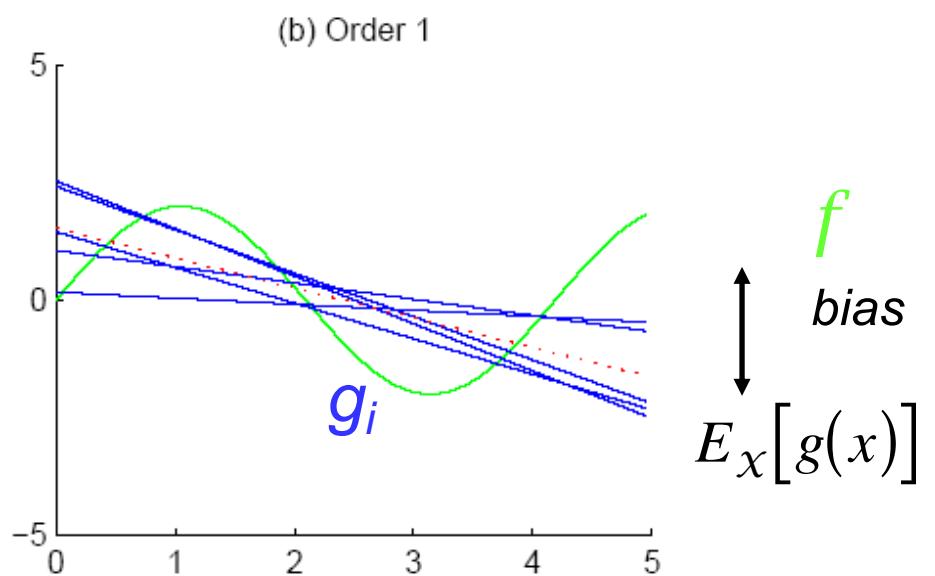
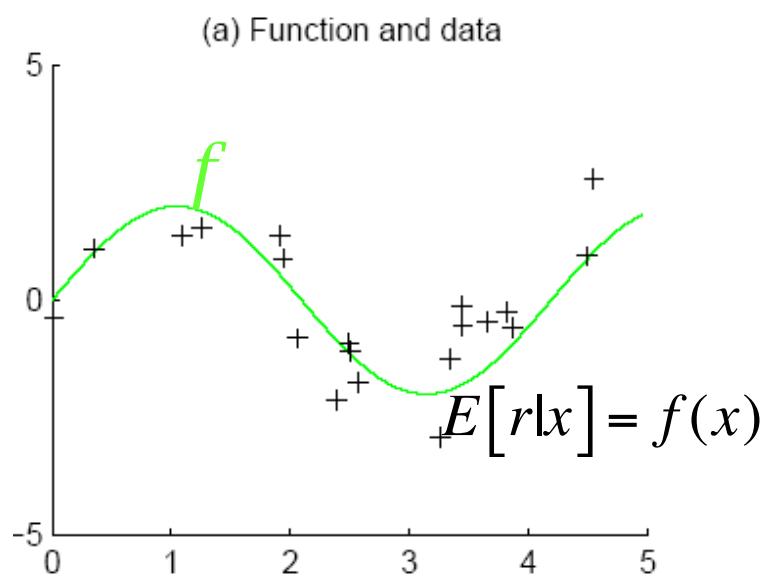
- Example: $g_i(x)=2$

has no variance and high bias

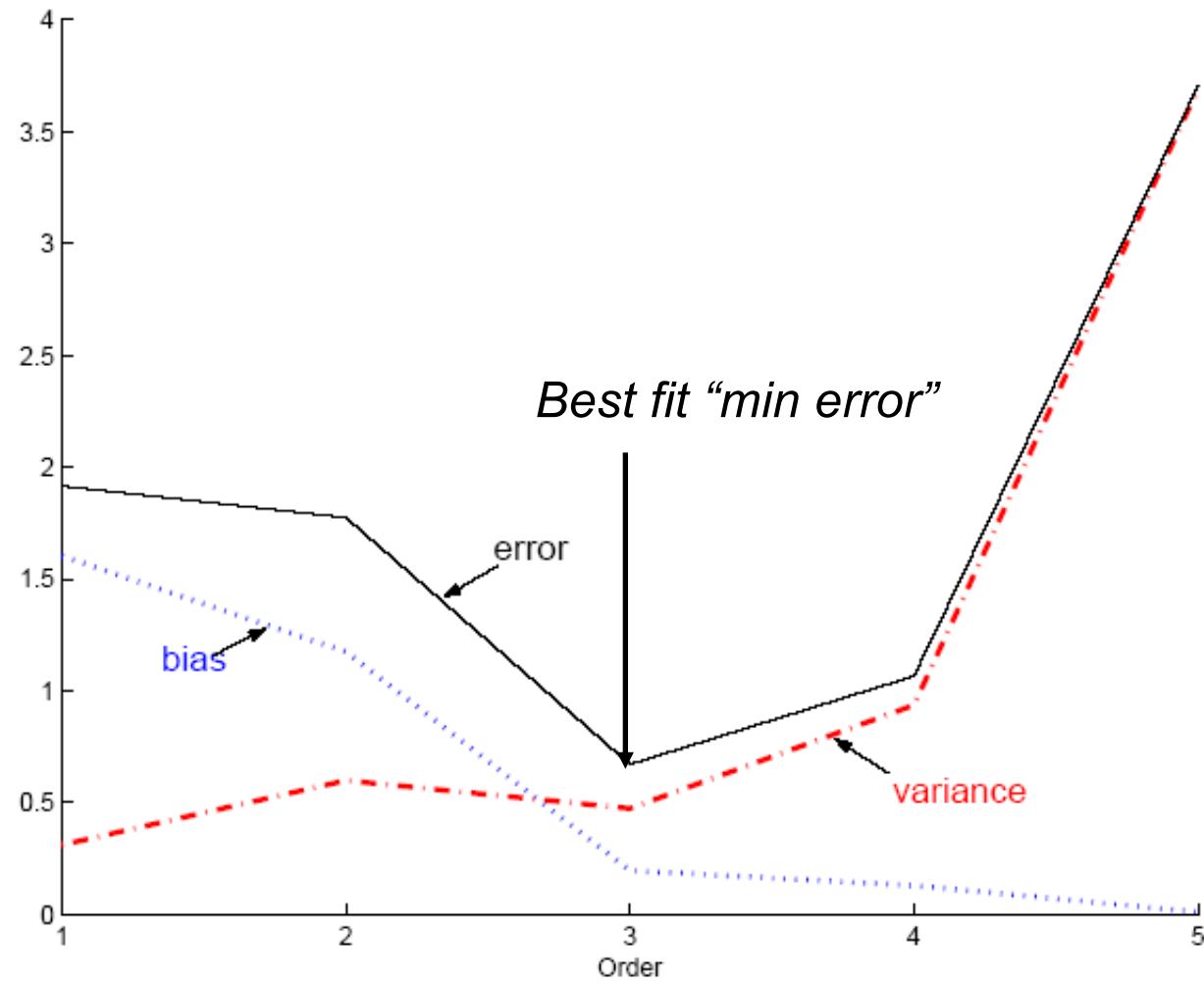
$$g_i(x) = \sum_t r_i^t / N$$

has lower bias with variance

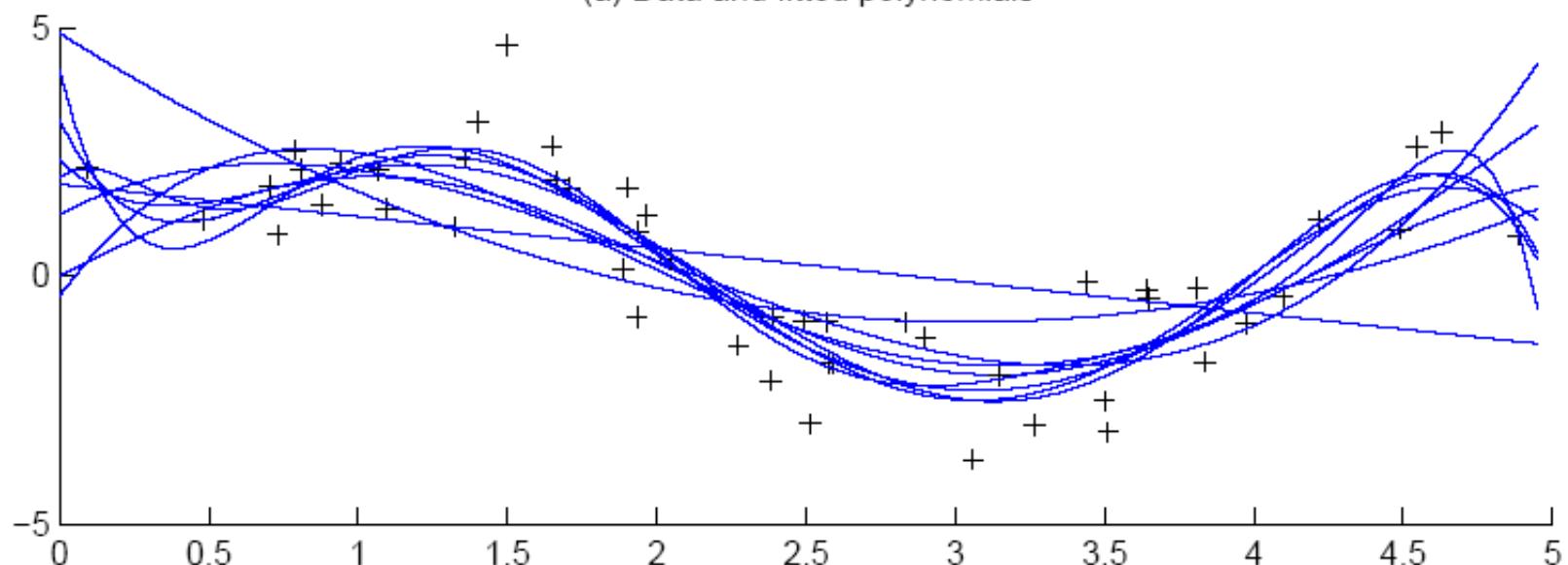
- As we increase complexity, bias decreases (a better fit to data) and variance increases (fit varies more with data)
- Bias/Variance dilemma



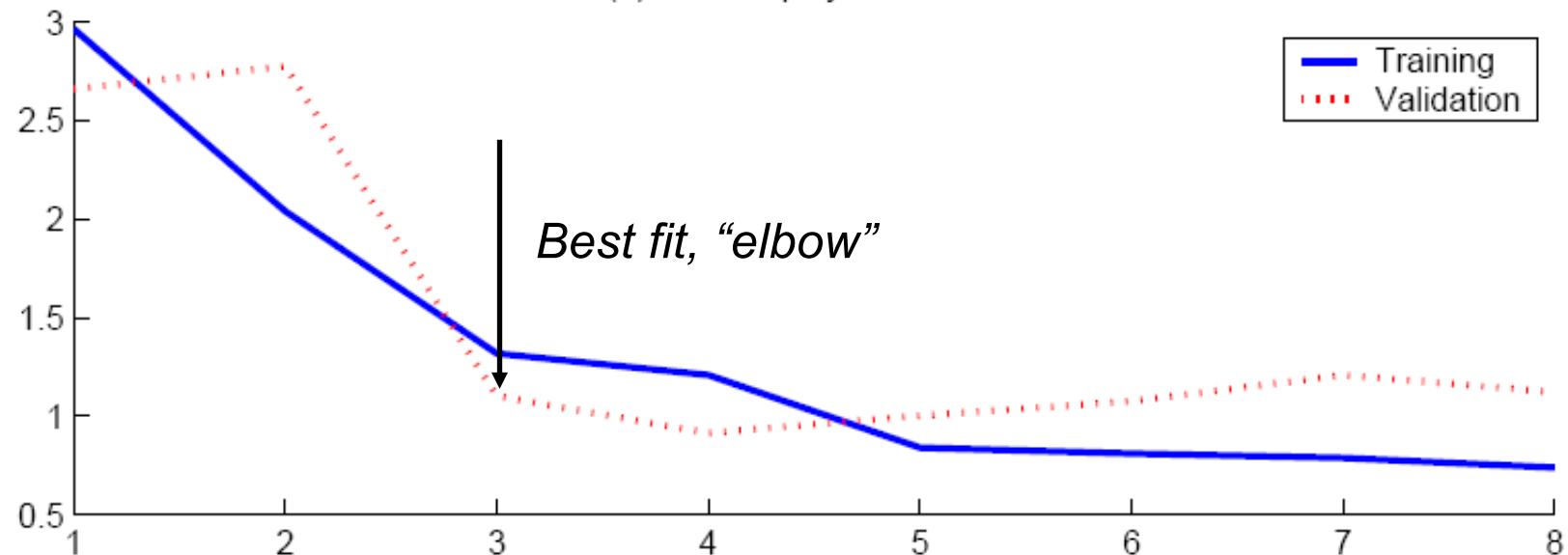
Polynomial Regression



(a) Data and fitted polynomials



(b) Error vs polynomial order

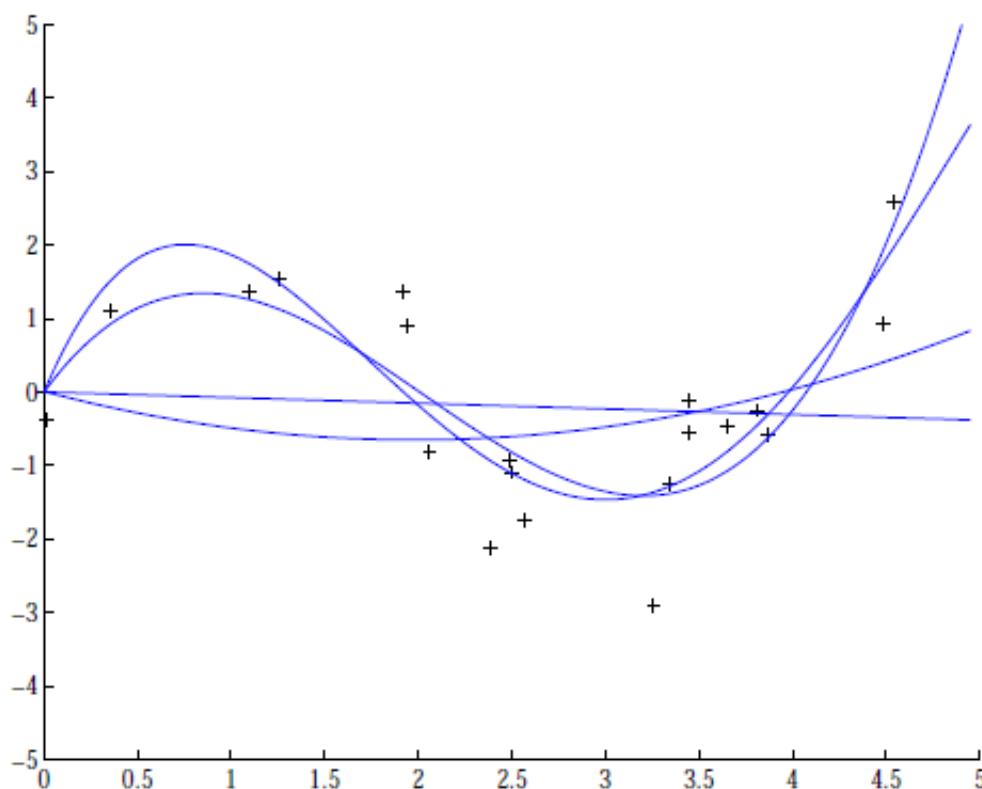


Model Selection

- **Cross-validation:** Measure generalization accuracy by testing on data unused during training
- **Regularization:** Penalize complex models
 $E' = \text{error on data} + \lambda \text{ model complexity}$
- **Regression with penalty on w:**

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

Regression example



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]

$$\text{regularization : } E(\mathbf{w} \mid \mathcal{X}) = \frac{1}{2} \sum_{t=1}^N \left[r^t - g(x^t \mid \mathbf{w}) \right]^2 + \lambda \sum_i w_i^2$$

Regression example

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

When λ is small ($\rightarrow 0$), the model is unregularized and achieve 0 training error with complex model.

When λ is larger ($+\infty$), $g(x) = 0$ large training error with simplest model.