

## CH. 4: Parametric Methods

Estimate the parameters of a known model from a given sample.

Two kinds of models are considered:

(1) **Deterministic models**, e.g., **regression functions**

**Example:**

Linear model:  $g(x | \theta) = w_1 x + w_0$

Parameters:  $\theta = (w_0, w_1)^T$

Quadratic model:  $g(x | \theta) = w_2 x^2 + w_1 x + w_0$

Parameters:  $\theta = (w_0, w_1, w_2)^T$

(2) Nondeterministic models, e.g., probability  
distribution functions

**Example:**

Bernoulli probability:  $P(x | \theta) = p^x (1 - p)^{1-x}$

Parameters:  $\theta = p$

Gaussian distribution:  $p(x | \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Parameters:  $\theta = (\mu, \sigma)^T$

## 4.1 Probability Distribution Functions

- Two methods of parameter estimation:
  - (a) Maximum likelihood approach
  - (b) Bayesian approach

### (A) Maximum Likelihood Estimation (MLE)

Let  $X = \{\mathbf{x}^t\}_{t=1}^N$  be the sample drawn from some probability density family,  $p(\mathbf{x}|\boldsymbol{\theta})$ . Estimate  $\boldsymbol{\theta}^*$  that makes  $X$  from  $p(\mathbf{x}|\boldsymbol{\theta}^*)$  as likely as possible, i.e., maximizes the likelihood  $p(X|\boldsymbol{\theta})$  of  $\boldsymbol{\theta}$ .

- The **likelihood**  $p(X | \theta)$  of  $\theta$  given  $X = \{\mathbf{x}^t\}_{t=1}^N$  is the **probability** of obtaining  $X$ , which is the product  $\prod_{t=1}^N p(\mathbf{x}^t | \theta)$  of the probabilities of obtaining the individual examples, i.e.,  $p(X | \theta) = \prod_{t=1}^N p(\mathbf{x}^t | \theta)$

**Log-likelihood:**  $L(\theta | X) = \log p(X | \theta)$

$$= \log \prod_{t=1}^N p(\mathbf{x}^t | \theta) = \sum_{t=1}^N \log p(\mathbf{x}^t | \theta)$$

- **Maximum likelihood estimator:**  $\theta^* = \arg \max_{\theta} L(\theta | X)$

## A.1 Bernoulli Distribution

The random variable  $X$  takes values  $x \in \{0, 1\}$ .

Probability function:  $P(x | \theta) = p^x (1-p)^{1-x}$

$p$  is the only parameter to be estimated, i.e.,  $\theta = p$ .

Given a sample  $X = \{x^t\}_{t=1}^N$ , the log-likelihood of  $\theta$ :

$$\begin{aligned} L(\theta | X) &= \sum_{t=1}^N \log P(x^t | \theta) = \sum_{t=1}^N \log p^{x^t} (1-p)^{1-x^t} \\ &= \log p^{x^1} (1-p)^{1-x^1} + \dots + \log p^{x^N} (1-p)^{1-x^N} \\ &= \log \left( p^{x^1} (1-p)^{1-x^1} p^{x^2} (1-p)^{1-x^2} \dots p^{x^N} (1-p)^{1-x^N} \right) \\ &= \log \left( p^{\sum_t x^t} (1-p)^{N-\sum_t x^t} \right) = \log p^{\sum_t x^t} + \log (1-p)^{N-\sum_t x^t} \\ &= \sum_t x^t \log p + (N - \sum_t x^t) \log(1-p) \end{aligned}$$

The maximum likelihood estimate (MLE) of  $p$ :

Let  $dL / dp = 0$

$$\frac{dL}{dp} = \frac{d}{dp} \left( \sum_t x^t \log p + (N - \sum_t x^t) \log(1 - p) \right)$$

$$\begin{aligned} &= \frac{1}{p} \sum_t x^t - \frac{1}{1-p} (N - \sum_t x^t) = \frac{(1-p) \sum_t x^t}{p(1-p)} - \frac{p(N - \sum_t x^t)}{p(1-p)} \\ &= 0 \end{aligned}$$

$$\Leftrightarrow (1-p) \sum_t x^t - p(N - \sum_t x^t) = 0$$

$$\Leftrightarrow \sum_t x^t - p \sum_t x^t - pN + p \sum_t x^t = 0 \Leftrightarrow \sum_t x^t - pN = 0$$

$$\hat{p} = \sum_{t=1}^N x^t / N$$

## A.2 Binomial Distribution

In binomial trials,  $N$  identical independent Bernoulli trials (0/1) are conducted.

Random variable  $X$  represents the number of 1s.

Probability function:  $P(x | \theta) = \binom{N}{x} p^x (1-p)^{N-x}$

$p$  is the only parameter to be estimated,

i.e.,  $\theta = p$ .

Given a sample  $X = \{x^t\}_{t=1}^n$ ,

derive the maximum likelihood estimate of  $\theta$ .

The log-likelihood of  $\theta$ :

$$\begin{aligned} L(\theta | X) &= \sum_{t=1}^n \log P(x^t | \theta) = \sum_{t=1}^n \log \binom{N}{x^t} p^{x^t} (1-p)^{N-x^t} \\ &= \log \binom{N}{x^1} p^{x^1} (1-p)^{N-x^1} + \log \binom{N}{x^2} p^{x^2} (1-p)^{N-x^2} \\ &\quad + \dots + \log \binom{N}{x^n} p^{x^n} (1-p)^{N-x^n} \\ &= \log \binom{N}{x^1} \dots \binom{N}{x^N} \left( p^{x^1} (1-p)^{N-x^1} \dots p^{x^n} (1-p)^{N-x^n} \right) \\ &\quad \begin{array}{c} \uparrow \text{-----} \uparrow \\ \downarrow \\ c \end{array} \end{aligned}$$



$$= \log c \left( p^{x^1} (1-p)^{N-x^1} \cdots p^{x^n} (1-p)^{N-x^n} \right)$$

$$\text{(where } c = \binom{N}{x^1} \cdots \binom{N}{x^n} \text{)}$$

$$= \log c \left( p^{\sum_t x^t} (1-p)^{Nn - \sum_t x^t} \right)$$

$$= \log c + \log p^{\sum_t x^t} + \log(1-p)^{Nn - \sum_t x^t}$$

$$= \log c + \sum_t x^t \log p + (Nn - \sum_t x^t) \log(1-p)$$

The maximum likelihood estimate (MLE) of  $p$ :

Let  $dL / dp = 0$

$$\begin{aligned}\frac{dL}{dp} &= \frac{d}{dp} (\log c + \sum_t x^t \log p + (Nn - \sum_t x^t) \log(1-p)) \\ &= \frac{\sum_t x^t}{p} - \frac{(Nn - \sum_t x^t)}{1-p} = \frac{(1-p)\sum_t x^t}{p(1-p)} - \frac{p(Nn - \sum_t x^t)}{p(1-p)} = 0\end{aligned}$$

$$\Leftrightarrow (1-p)\sum_t x^t - p(Nn - \sum_t x^t) = 0$$

$$\Leftrightarrow \sum_t x^t - p\sum_t x^t - pNn + p\sum_t x^t = 0 \Leftrightarrow \sum_t x^t - pNn = 0$$

$$\hat{p} = \sum_t x^t / Nn$$

$$\hat{p} = \sum_{t=1}^n x^t / Nn$$

## A.3 Gaussian (Normal) Distribution

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

Parameters to be estimated:  $\theta = (\mu, \sigma)^T$

Given a sample:  $X = \{x^t\}_{t=1}^N$

Log-likelihood:  $L(\theta | X) = \log \prod_{t=1}^N p(x^t | \theta)$

$$\begin{aligned} &= \log \left( \prod_{t=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x^t - \mu)^2}{2\sigma^2}\right] \right) \\ &= \log \left( \left[ \frac{1}{\sqrt{2\pi}\sigma} \right]^N \prod_{t=1}^N \exp\left[-\frac{(x^t - \mu)^2}{2\sigma^2}\right] \right) \end{aligned}$$

$$\begin{aligned}
&= \log \left[ \frac{1}{\sqrt{2\pi\sigma}} \right]^N + \log \left( \prod_{t=1}^N \exp \left[ -\frac{(x^t - \mu)^2}{2\sigma^2} \right] \right) \\
&= N \log \left[ \frac{1}{\sqrt{2\pi\sigma}} \right] + \sum_{t=1}^N \log \exp \left[ -\frac{(x^t - \mu)^2}{2\sigma^2} \right] \\
&= -N \log \left[ \sqrt{2\pi\sigma} \right] + \sum_{t=1}^N \left[ -\frac{(x^t - \mu)^2}{2\sigma^2} \right] \\
&= -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum_{t=1}^N (x^t - \mu)^2}{2\sigma^2}
\end{aligned}$$

The maximum likelihood estimate (MLE) of  $(\mu, \sigma)$ :

Let  $\partial L / \partial \mu = 0$ ,  $\partial L / \partial \sigma = 0$

$$\frac{\partial L}{\partial \mu} = \frac{\partial}{\partial \mu} \left( -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum_{t=1}^N (x^t - \mu)^2}{2\sigma^2} \right)$$

$$= \frac{-1}{2\sigma^2} \frac{\partial}{\partial \mu} \sum_{t=1}^N (x^t - \mu)^2 = \frac{-1}{2\sigma^2} \sum_{t=1}^N \frac{\partial}{\partial \mu} (x^t - \mu)^2$$

$$= \frac{-1}{2\sigma^2} \sum_{t=1}^N (-2(x^t - \mu)) = \frac{1}{\sigma^2} \sum_{t=1}^N (x^t - \mu) = 0$$

$$\Rightarrow \sum_{t=1}^N (x^t - \mu) = 0, \sum_{t=1}^N x^t - N\mu = 0, \hat{\mu} = \frac{1}{N} \sum_{t=1}^N x^t$$

$$\begin{aligned}
\frac{\partial L}{\partial \sigma} &= \frac{\partial}{\partial \sigma} \left( -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum_{t=1}^N (x^t - \mu)^2}{2\sigma^2} \right) \\
&= \frac{\partial}{\partial \sigma} \left( -N \log \sigma - \frac{\sum_{t=1}^N (x^t - \mu)^2}{2\sigma^2} \right) \\
&= \frac{-N}{\sigma} - \frac{1}{2} \sum_{t=1}^N (x^t - \mu)^2 \frac{\partial}{\partial \sigma} \left( \frac{1}{\sigma^2} \right) \\
&= \frac{-N}{\sigma} + \frac{1}{\sigma^3} \sum_{t=1}^N (x^t - \mu)^2 = 0 \\
\Rightarrow -N + \frac{1}{\sigma^2} \sum_{t=1}^N (x^t - \mu)^2 &= 0, \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{t=1}^N (x^t - \mu)^2
\end{aligned}$$

## (B) Bayes' Estimation

-- Treat  $\theta$  as a random variable with prior density  $p(\theta)$ .

Use Bayes' rule to combine the likelihood  $p(X | \theta)$

to obtain the posterior  $p(\theta | X) = p(X | \theta)p(\theta) / p(X)$

**Bayes' estimator:**  $\theta_{Bayes} = E[\theta | X] = \int \theta p(\theta | X) d\theta$

i.e., the posterior expected value of  $\theta$ .

**ML estimator:**  $\theta_{ML} = \arg \max_{\theta} p(X | \theta)$

**MAP estimator:**  $\theta_{MAP} = \arg \max_{\theta} p(\theta | X)$

**Example:** Given a sample  $X = \{x^t\}_{t=1}^N$ ,

Suppose  $x^t \sim N(\theta, \sigma^2)$ ,  $p(x^t) = \frac{1}{(2\pi)^{1/2} \sigma} \exp\left[-\frac{(x^t - \theta)^2}{2\sigma^2}\right]$

$\theta = \mu$  unknown,  $\sigma$  known

$$\begin{aligned} \text{Likelihood of } \theta : p(X | \theta) &= \prod_{t=1}^N \frac{1}{(2\pi)^{1/2} \sigma} \exp\left[-\frac{(x^t - \theta)^2}{2\sigma^2}\right] \\ &= \frac{1}{(2\pi)^{N/2} \sigma^N} \exp\left[-\frac{\sum_t (x^t - \theta)^2}{2\sigma^2}\right] \end{aligned}$$

Assume prior:  $\theta \sim N(\mu_0, \sigma_0^2)$ ,  $\mu_0, \sigma_0$ : known

$$p(\theta) = \frac{1}{(2\pi)^{1/2} \sigma_0} \exp\left[-\frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right]$$

Posterior:

$$p(\theta | X) = p(X | \theta) p(\theta) / p(X), \quad p(X) = \sum_{\theta} p(X | \theta) p(\theta)$$



$$p(\theta | X) = \frac{1}{(2\pi)^{(N+1)/2} \sigma^N \sigma_0} \exp \left[ - \left( \frac{\sum_t (x^t - \theta)^2}{2\sigma^2} + \frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right) \right] / p(X)$$

$$\theta_{Bay} = E[\theta | X] = \int \theta p(\theta | X) d\theta = \frac{1}{(2\pi)^{(N+1)/2} \sigma^N \sigma_0} \cdot$$

$$\int \theta \exp \left[ - \left( \frac{\sum_t (x^t - \theta)^2}{2\sigma^2} + \frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right) \right] / p(X) d\theta$$

$$= \frac{N / \sigma^2}{N / \sigma^2 + 1 / \sigma_0^2} m + \frac{1 / \sigma_0^2}{N / \sigma^2 + 1 / \sigma_0^2} \mu_0$$

$$N \text{ large} \rightarrow \theta_{Bay} \text{ close to } m = \sum_{t=1}^N x^t;$$

$$\sigma_0^2 \text{ small or } N \text{ small} \rightarrow \theta_{Bay} \text{ close to } \mu_0.$$

## 4.2 Regression Functions

-- Given a sample  $(x^t, r^t)_{t=1}^N$ , determine the parameters  $\theta$  of the function  $f$  for  $r = f(x)$ .

Example:

**Linear function**  $f(x|\theta = (w_0, w_1)^T) = w_1 x + w_0$

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

$$\frac{\partial}{\partial w_0} E(w_0, w_1 | X) = 0, \quad \sum_t r^t = N w_0 + w_1 \sum_t x^t,$$

$$\frac{\partial}{\partial w_1} E(w_0, w_1 | X) = 0, \quad \sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2$$

In vector-matrix form  $A\mathbf{w} = \mathbf{r}$ , where

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

The solution:  $\mathbf{w} = A^{-1}\mathbf{r}$

## Polynomial function

$$f(x|\boldsymbol{\theta} = (w_k, \dots, w_0)) = w_k x^k + \dots + w_1 x + w_0$$

$$\begin{aligned} E(\boldsymbol{\theta} | X) &= \frac{1}{2} \sum_{t=1}^N \left[ r^t - f(x^t | \boldsymbol{\theta}) \right]^2 \\ &= \frac{1}{2} \sum_{t=1}^N \left[ r^t - (w_k (x^t)^k + \dots + w_1 (x^t) + w_0) \right]^2 \end{aligned}$$

$$\frac{\partial E(\boldsymbol{\theta} | X)}{\partial w_0} = 0, \dots, \frac{\partial E(\boldsymbol{\theta} | X)}{\partial w_k} = 0$$

$$\sum_t r^t = w_0 \sum_t (x^t)^0 + \dots + w_k \sum_t (x^t)^k,$$

$$\dots,$$

$$\sum_t r^t (x^t)^k = w_0 \sum_t (x^t)^k + \dots + w_k \sum_t (x^t)^{2k}$$

In vector-matrix form  $A\mathbf{w} = \mathbf{r}$ , where

$$A = \begin{bmatrix} \sum_t (x^t)^0 & \sum_t (x^t)^1 & \dots & \sum_t (x^t)^k \\ \sum_t (x^t)^1 & \sum_t (x^t)^2 & \dots & \sum_t (x^t)^{k+1} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_t (x^t)^k & \sum_t (x^t)^{k+1} & \dots & \sum_t (x^t)^{2k} \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_k \end{bmatrix}, \quad \mathbf{r} = \begin{bmatrix} \sum_t r^t (x^t)^0 \\ \sum_t r^t (x^t)^1 \\ \vdots \\ \sum_t r^t (x^t)^k \end{bmatrix}$$

Find  $\mathbf{w}$  that minimizes  $e = \|\mathbf{A}\mathbf{w} - \mathbf{r}\|^2$ .

$$\text{Let } \frac{de}{d\mathbf{w}} = \frac{d}{d\mathbf{w}} \|\mathbf{A}\mathbf{w} - \mathbf{r}\|^2 = 2\mathbf{A}^T (\mathbf{A}\mathbf{w} - \mathbf{r}) = \mathbf{0}$$

$$\mathbf{A}^T (\mathbf{A}\mathbf{w} - \mathbf{r}) = \mathbf{0}, \quad \mathbf{A}^T \mathbf{A}\mathbf{w} - \mathbf{A}^T \mathbf{r} = \mathbf{0}, \quad \mathbf{A}^T \mathbf{A}\mathbf{w} = \mathbf{A}^T \mathbf{r}$$

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

$$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T : \text{ the } \textcolor{red}{\text{pseudoinverse}} \text{ of } \mathbf{A}.$$

## Relate parameter estimations of deterministic and nondeterministic models

Consider equation  $r = f(x)$ . Introduce uncertainty  $\mathcal{E}$  into the equation  $r = f(x) + \mathcal{E}$  and view  $\mathcal{E}$  and  $r$  as random variables.

Assume  $\mathcal{E} \sim p(e)$ , from the probability theorem,

$$r \sim p(r) = p(e) |de / dr|$$

$$r = f(x) + e, \quad e = r - f(x), \quad \left| \frac{de}{dr} \right| = 1$$

$$p(r) = p(e) |de / dr| = p(e)$$

**Example:** Given  $r = f(x) + \varepsilon$ , let  $\varepsilon \sim p(e) = N(0, \sigma^2)$ ,

$$\text{i.e., } p(e) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{e^2}{2\sigma^2}\right].$$

From theorem,  $r \sim p(r) = p(e) |de/dr|$ .

$$\begin{aligned} p(r) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{e^2}{2\sigma^2}\right] \left| \frac{de}{dr} \right| \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{[r - \hat{f}(x|\theta)]^2}{2\sigma^2}\right] \cdot 1 \end{aligned} \quad \left( \begin{array}{l} r = f(x) + e \\ e = r - f(x) \\ |de/dr| = 1 \end{array} \right)$$

where  $\hat{f}(x|\theta)$ : the estimator of  $f(x)$  up to  $\theta$ .

The log-likelihood of  $\theta$  :  $L(\theta | X) = \log \prod_{t=1}^N p(r^t | x^t)$

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



$$L(\boldsymbol{\theta} | X) = -N \log \sqrt{2\pi} \sigma - \frac{1}{2\sigma^2} E(\boldsymbol{\theta} | X)$$

$\therefore$  Maximizing  $L(\boldsymbol{\theta} | X) \Leftrightarrow$  Minimizing  $E(\boldsymbol{\theta} | X)$

(Probabilistic functions)    (Regression functions)

## 4.3 Model Complexity: Bias and Variance

Given  $M$  samples  $X_i = \{x_i^t, r_i^t\}_{t=1}^M$ ,  $t = 1, \dots, N$  to fit  $f(x)$ . Let  $\hat{f}_i(x)$ ,  $i = 1, \dots, M$ , be the estimates.

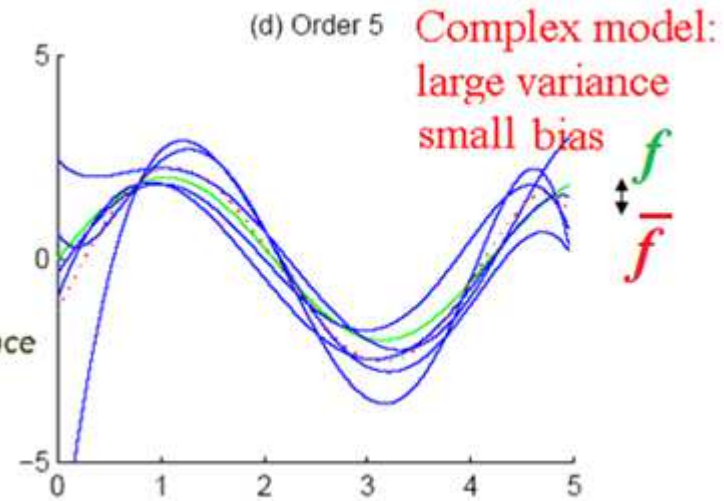
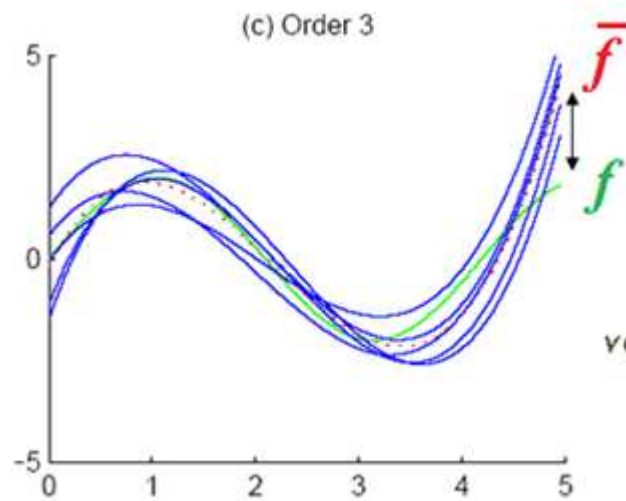
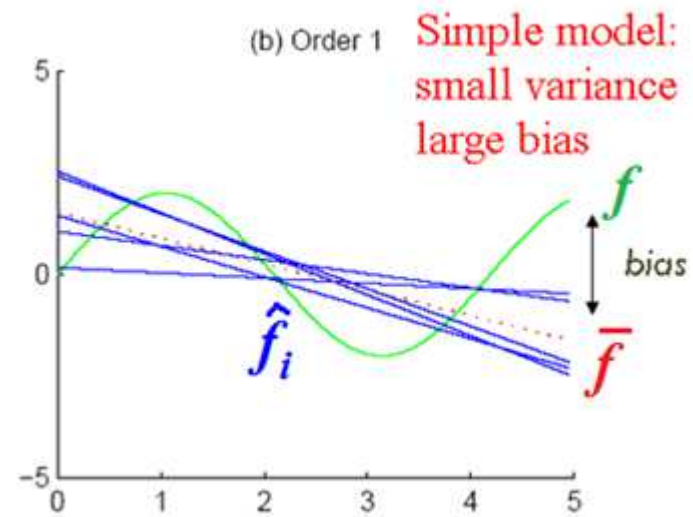
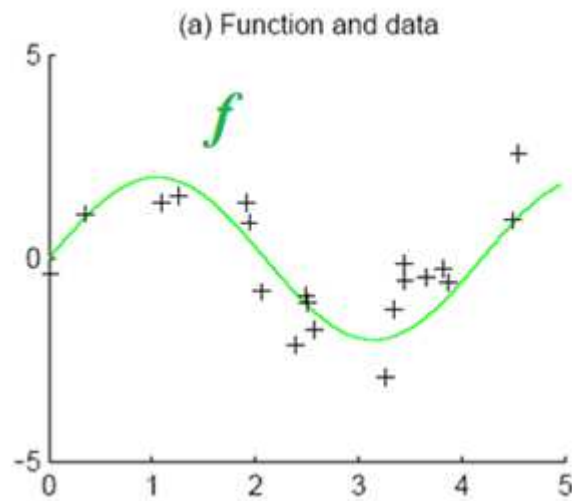
$$\text{Let } \bar{f}(x) = \frac{1}{M} \sum_i \hat{f}_i(x)$$

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

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

**Bias/Variance Dilemma** -- Increase model complexity,

➡ bias decreases (a better fit to data) and  
variance increases (fit varies more with data)



## Appendix: Derivatives

- $\frac{df(x)}{dx}$ , e.g.,  $f(x) = 3x^2 - 4x + 15$ ,  $\frac{df(x)}{dx} = 6x - 4$

- $\frac{df(\mathbf{x})}{d\mathbf{x}} = \nabla_{\mathbf{x}} f(\mathbf{x}) = \left( \frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_d} \right)^T$ ,

where  $\mathbf{x} = (x_1 \ x_2 \ \dots \ x_d)$

e.g.,  $f(\mathbf{x}) = f(x, y) = 3x - 4xy + 15y - 6$ ,

$$\frac{df(\mathbf{x})}{d\mathbf{x}} = \left( \frac{\partial f(x, y)}{\partial x} \quad \frac{\partial f(x, y)}{\partial y} \right)^T = \begin{pmatrix} 3 - 4y \\ -4x + 15 \end{pmatrix}$$

- $$\frac{d\mathbf{f}(x)}{dx} = \frac{d(f_1(x) \cdots f_n(x))^T}{dx} = \left( \frac{df_1(x)}{dx} \cdots \frac{df_n(x)}{dx} \right)^T$$

e.g.,  $\mathbf{f}(x) = (f_1(x) \ f_2(x))^T = ((3x-4) \ (x^2+5x-6))^T$ ,

$$\begin{aligned} \frac{d\mathbf{f}(x)}{dx} &= \frac{d(f_1(x) \ f_2(x))^T}{dx} = \frac{d((3x-4) \ (x^2+5x-6))^T}{dx} \\ &= \left( \frac{d(3x-4)}{dx} \quad \frac{d(x^2+5x-6)}{dx} \right)^T \\ &= (3 \ 2x+5)^T \end{aligned}$$

$$\begin{aligned}
\bullet \quad \frac{df(\mathbf{x})}{d\mathbf{x}} &= \frac{d(f_1(\mathbf{x}) \cdots f_n(\mathbf{x}))^T}{d\mathbf{x}} = \left( \frac{df_1(\mathbf{x})}{d\mathbf{x}} \cdots \frac{df_n(\mathbf{x})}{d\mathbf{x}} \right)^T \\
&= (\nabla_{\mathbf{x}} f_1(\mathbf{x}) \quad \nabla_{\mathbf{x}} f_2(\mathbf{x}) \quad \cdots \quad \nabla_{\mathbf{x}} f_n(\mathbf{x}))^T \\
&= \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_d} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_2(\mathbf{x})}{\partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n(\mathbf{x})}{\partial x_1} & \frac{\partial f_n(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_n(\mathbf{x})}{\partial x_d} \end{pmatrix} = J(\mathbf{f})
\end{aligned}$$

Jacobian matrix of  $\mathbf{f}$

$$\begin{aligned}
\bullet \quad \frac{dF(\mathbf{x})}{d\mathbf{x}} &= \frac{d}{d\mathbf{x}} \begin{pmatrix} f_{11} & f_{12} & \cdots & f_{1m} \\ f_{21} & f_{22} & \cdots & f_{2m} \\ \cdots & \cdots & \cdots & \cdots \\ f_{n1} & f_{n2} & \cdots & f_{nm} \end{pmatrix} = \begin{pmatrix} \frac{df_{11}}{d\mathbf{x}} & \cdots & \frac{df_{1m}}{d\mathbf{x}} \\ \cdots & \cdots & \cdots \\ \frac{df_{n1}}{d\mathbf{x}} & \cdots & \frac{df_{nm}}{d\mathbf{x}} \end{pmatrix} \\
&= \begin{pmatrix} J_{11} & J_{12} & \cdots & J_{1m} \\ J_{21} & J_{22} & \cdots & J_{2m} \\ \cdots & \cdots & \cdots & \cdots \\ J_{n1} & J_{n2} & \cdots & J_{nm} \end{pmatrix}, \text{ where } J_{ij} = J(f_{ij}(\mathbf{x}))
\end{aligned}$$

Hessian matrix of  $F$

- $\frac{df(X)}{dX}$ , where  $X = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \cdots & \cdots & \cdots \\ x_{m1} & \cdots & x_{mn} \end{pmatrix}$

- $\frac{df(X)}{dX}$

- $\frac{dF(X)}{dX}$