Lecture Slides for

INTRODUCTION TO

# Machine Learning
## 2nd Edition

ETHEM ALPAYDIN
© The MIT Press, 2010

*alpaydin@boun.edu.tr*
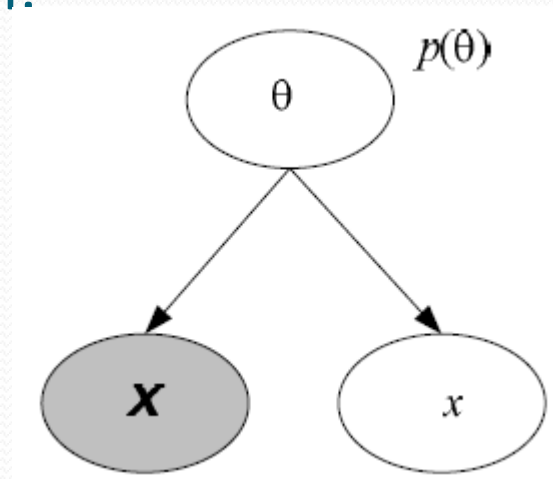*http://www.cmpe.boun.edu.tr/~ethem/i2ml2e*

# Bayesian Estimation

# Rationale

- Bayes' Rule:

$$p(\theta \mid \mathrm{X}) = \frac{p(\theta)p(\mathrm{X} \mid \theta)}{p(\mathrm{X})}$$

- Generative model:

# Estimating the Parameters of a Distribution: Discrete case

- $x_i^t = 1$ if in instance $t$ is in state $i$, probability of state $i$ is $q_i$
- Dirichlet prior, $\alpha_i$ are hyperparameters

$$Dirichlet(\mathbf{q} \mid \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_K)} \prod_{i=1}^{K} q_i^{\alpha_i - 1}$$

- Sample likelihood

$$p(X \mid \mathbf{q}) = \prod_{t=1}^{N} \prod_{i=1}^{K} q_i^{x_i^t}$$

- Posterior

$$p(\mathbf{q} \mid \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + N_1)\cdots\Gamma(\alpha_K + N_K)} \prod_{i=1}^{K} q_i^{\alpha_i + N_i - 1}$$

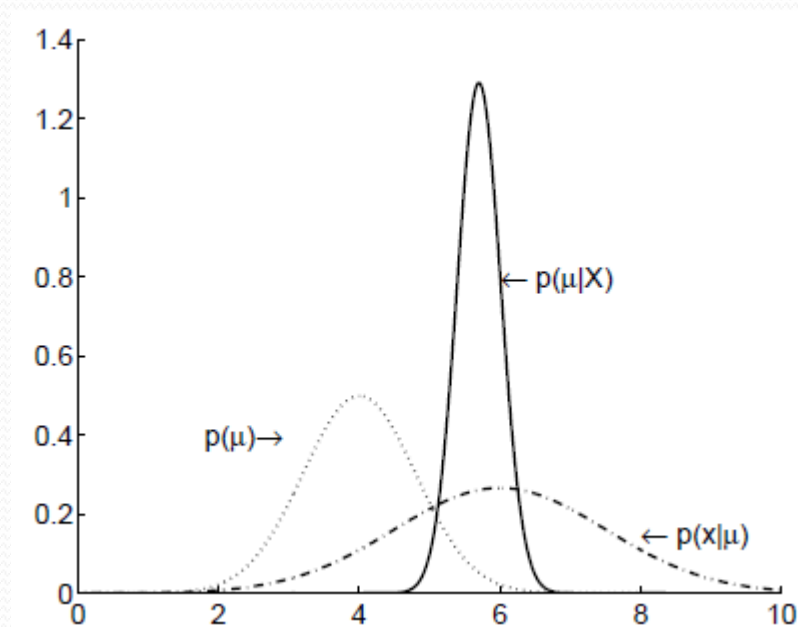$$= Dirichlet(\mathbf{q} \mid \boldsymbol{\alpha} + \mathbf{n})$$

- Dirichlet is a conjugate prior
- With $K$=2, Dirichlet reduced to Beta

# Estimating the Parameters of a Distribution: Continuous case

- $p(x^t) \sim N(\mu, \sigma^2)$
- Gaussian prior for $\mu$, $p(\mu) \sim N(\mu_0, \sigma_0^2)$
- Posterior is also Gaussian $p(\mu|X) \sim N(\mu_N, \sigma_N^2)$ where

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

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

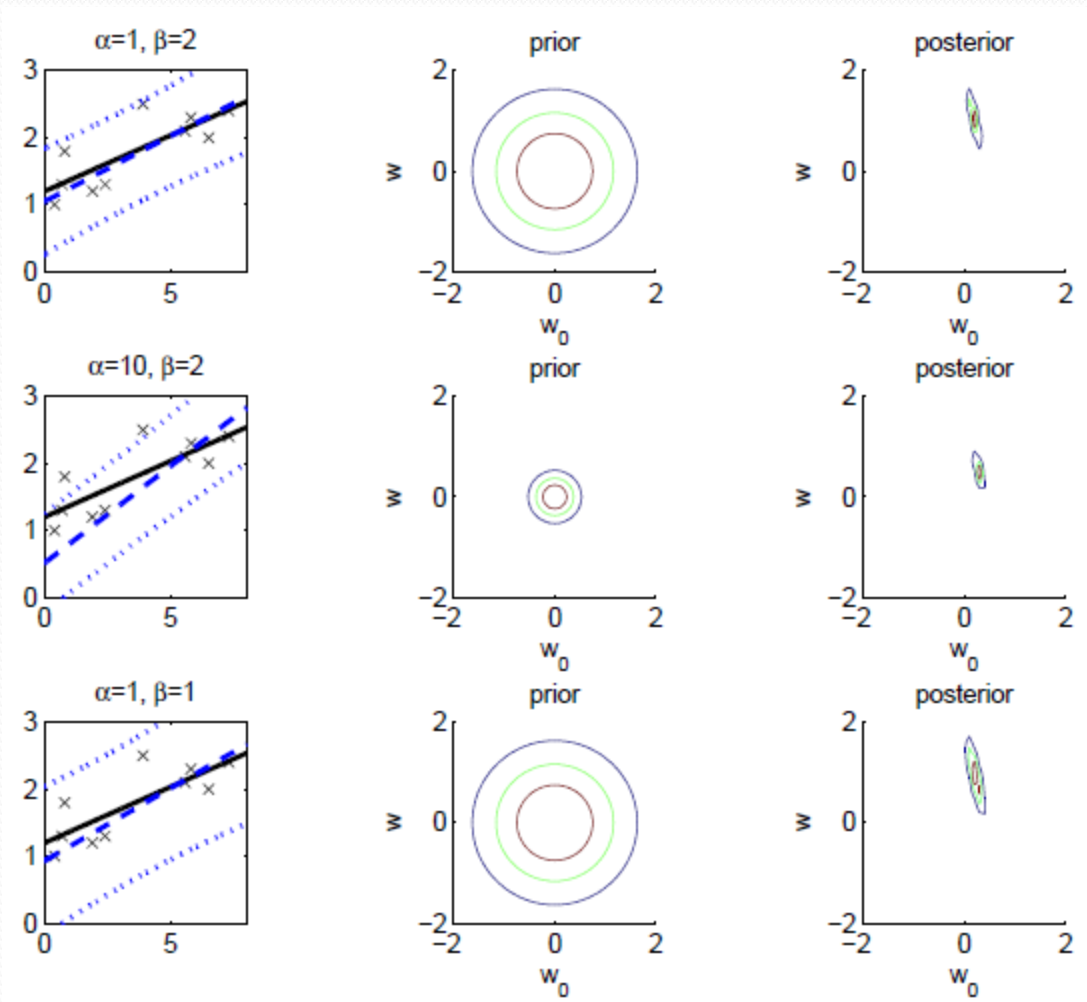# Estimating the Parameters of a Function: Regression

- $r = \mathbf{w}^T \mathbf{x} + \varepsilon$ where $p(\varepsilon) \sim N(0, 1/\beta)$, and $p(r^t | x^t, \mathbf{w}, \beta) \sim N(\mathbf{w}^T \mathbf{x}^t, 1/\beta)$

- Log likelihood

$$L(\mathbf{r} | \mathbf{X}, \mathbf{w}, \beta) = \log \prod_t p(r^t | \mathbf{x}^t, \mathbf{w}, \beta)$$

$$= -N \log\left(\sqrt{2\pi}\right) + N \log \beta - \frac{\beta}{2} \sum_t \left(r^t - \mathbf{w}^T \mathbf{x}^t\right)$$

- ML solution $\quad \mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r}$

- Gaussian conjugate prior: $p(\mathbf{w}) \sim N(0, 1/\alpha)$

- Posterior: $p(\mathbf{w}|\mathbf{X}) \sim N(\boldsymbol{\mu}_N, \Sigma_N)$ where

$$\boldsymbol{\mu}_N = \beta \boldsymbol{\Sigma}_N \mathbf{X}^T \mathbf{r}$$

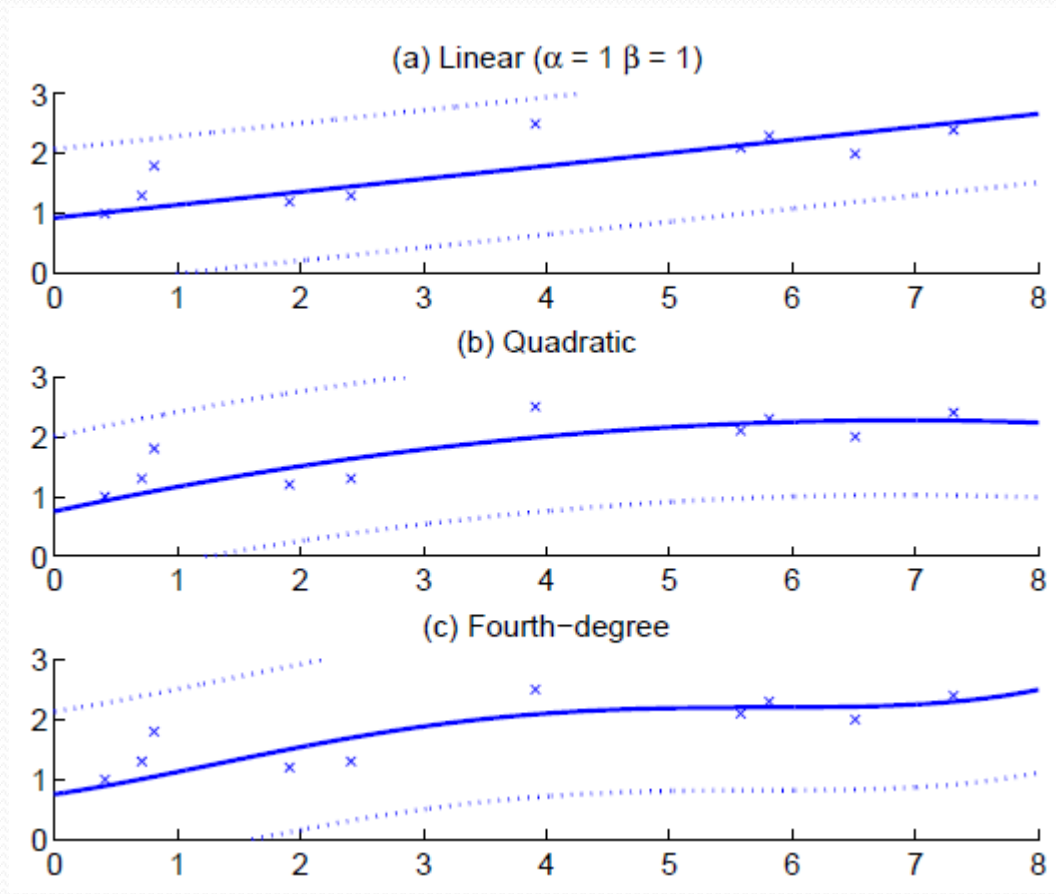$$\boldsymbol{\Sigma}_N = (\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})^{-1}$$

# Basis/Kernel Functions

- For new $\mathbf{x}'$, the estimate r' is calculated as

$$r' = (\mathbf{x}')^T$$

$$= \beta(\mathbf{x}')^T \mathbf{\Sigma}_N \mathbf{X}^T \mathbf{r}$$

$$= \sum_t \beta(\mathbf{x}')^T \mathbf{\Sigma}_N \mathbf{x}^t r^t \qquad \text{Dual representation}$$

- Linear kernel

$$r' = \sum_t \beta(\mathbf{x}')^T \mathbf{\Sigma}_N \mathbf{x}^t r^t \sum_t \beta K(\mathbf{x}', \mathbf{x}^t) r^t$$

- For any other $\phi(\mathbf{x})$, we can write $K(\mathbf{x}', \mathbf{x}) = \phi(\mathbf{x}')^T \phi(\mathbf{x})$

# Kernel Functions

# Gaussian Processes

- Assume Gaussian prior $p(\mathbf{w}) \sim N(0, 1/\alpha)$
- $\mathbf{y} = \mathbf{X}\mathbf{w}$, where $E[\mathbf{y}] = 0$ and $\text{Cov}(\mathbf{y}) = \mathbf{K}$ with $\mathbf{K}_{ij} = (\mathbf{x}^i)^T \mathbf{x}^i$
- $\mathbf{K}$ is the covariance function, here linear
- With basis function $\phi(\mathbf{x})$, $\mathbf{K}_{ij} = (\phi(\mathbf{x}^i))^T \phi(\mathbf{x}^i)$
- $r \sim N_N(\mathbf{0}, C_N)$ where $C_N = (1/\beta)\mathbf{I} + \mathbf{K}$
- With new $\mathbf{x}'$ added as $\mathbf{x}_{N+1}$, $r_{N+1} \sim N_{N+1}(0, C_{N+1})$

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k} & c \end{bmatrix}$$

where $\mathbf{k} = [K(\mathbf{x}', \mathbf{x}^t)_t]^T$ and $c = K(\mathbf{x}', \mathbf{x}') + 1/\beta$.
$p(r' | \mathbf{x}', \mathbf{X}, \mathbf{r}) \sim N(\mathbf{k}^T \mathbf{C}_{N-1} \mathbf{r}, c - \mathbf{k}^T \mathbf{C}_{N-1} \mathbf{k})$

(a) Linear ($\alpha = 1$ $\beta = 5$)

(b) Quadratic

(c) Gaussian