

**Problem 2:** Let's assume we have a dataset where each datapoint,  $(x_i, y_i)$  is weighted by a scalar factor which we will call  $t_i$ . We will assume that  $t_i > 0$  for all  $i$ . This makes the sum of squares error function look like the following:

$$E_{\text{weighted}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N t_i [\mathbf{w}^T \phi(x_i) - y_i]^2$$

Find the equation for the value of  $\mathbf{w}$  that minimizes this error function.

Furthermore, explain how this weighting factor,  $t_i$ , can be interpreted in terms of

- 1) the variance of the noise on the data and
- 2) data points for which there are exact copies in the dataset.

$$E = \frac{1}{2} (\Phi \mathbf{w} - \mathbf{y})^T \mathbf{T} (\Phi \mathbf{w} - \mathbf{y})$$

$$\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T \mathbf{A} \mathbf{w} = 2 \mathbf{A} \mathbf{w} \leftarrow$$

$$\mathbf{T} = \begin{bmatrix} t_1 & & & \\ & t_2 & & \\ & & \ddots & \\ & & & t_N \end{bmatrix}$$

$$E = \frac{1}{2} [(\mathbf{w}^T \Phi^T - \mathbf{y}^T) \mathbf{T} (\Phi \mathbf{w} - \mathbf{y})]$$

$$= \frac{1}{2} [\mathbf{w}^T \Phi^T \mathbf{T} \Phi \mathbf{w} - 2 \mathbf{w}^T \Phi^T \mathbf{T} \mathbf{y} + \mathbf{y}^T \mathbf{T} \mathbf{y}]$$

$$\frac{\partial E}{\partial \mathbf{w}} = \frac{1}{2} [2 \Phi^T \mathbf{T} \Phi \mathbf{w} - 2 \Phi^T \mathbf{T} \mathbf{y}] \stackrel{!}{=} 0$$

$$\Phi^T \mathbf{T} \Phi \mathbf{w}_* = \Phi^T \mathbf{T} \mathbf{y} \Rightarrow \mathbf{w}_* = (\Phi^T \mathbf{T} \Phi)^{-1} \Phi^T \mathbf{T} \mathbf{y}$$

$$\text{if } \mathbf{T} = \mathbf{I}_{N \times N} \Rightarrow \mathbf{w}_* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \text{ (standard formula)}$$

Lemma:

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} & & \\ & A & \\ & & \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \sum_{i,j} A_{ij} x_i x_j$$

$$[w^T \phi(x_1) - y_1, \dots, w^T \phi(x_N) - y_N] \quad \begin{array}{|c|} \hline 1 \\ \hline \end{array} \quad \begin{bmatrix} w^T \phi(x_1) - y_1 \\ \vdots \\ w^T \phi(x_N) - y_N \end{bmatrix}$$

$$y_i \sim \mathcal{N}(y_i | w^T x_i, \beta^{-1})$$

Approach ①  $y_i \sim \mathcal{N}(y_i | w^T x_i, t_i^{-1})$

$$\log P(D|w) = \prod_i \log \mathcal{N}(y_i | w^T x_i, t_i^{-1})$$

$$= \sum_i -\frac{t_i}{2} (y_i - w^T x_i)^2$$

$$= w_{MLE} = \underset{w}{\operatorname{argmax}} \left[ -\sum_i \frac{t_i}{2} (y_i - w^T x_i)^2 \right]$$

$$= \underset{w}{\operatorname{argmin}} \left[ \sum_i t_i (y_i - w^T x_i)^2 \right]$$

Approach ②  $y_i \sim \mathcal{N}(y_i | w^T x_i, t_i \beta^{-1})$

Interpretation 2:

$$\begin{array}{l} x_1 \rightsquigarrow y_1 \\ x_2 \rightsquigarrow y_2 \\ x_3 \rightsquigarrow y_3 \end{array} \Rightarrow \begin{cases} x_1 \rightsquigarrow y_1 \\ x_1 \rightsquigarrow y_2 \\ x_2 \rightsquigarrow y_1 \\ x_2 \rightsquigarrow y_2 \\ x_3 \rightsquigarrow y_1 \\ x_3 \rightsquigarrow y_2 \\ \vdots \end{cases}$$

$$\Rightarrow E = \sum_i t_i (y_i - w^T x_i)^2$$

Works only if  $t_i \in \mathbb{N}$

$$\left| \begin{array}{l} x_3 \leadsto y_3 \\ \vdots \\ a_3 \leadsto y_3 \end{array} \right.$$

**Problem 3:** Show that the following holds: The ridge regression estimates can be obtained by ordinary least squares regression on an augmented dataset: Augment the design matrix  $\Phi \in \mathbb{R}^{N \times M}$  with  $M$  additional rows  $\sqrt{\lambda} I_{M \times M}$  and augment  $y$  with  $M$  zeros.

$$E = \frac{1}{2} \left( \sum_i (y_i - \Phi(x_i)^T w)^2 \right) + \frac{\lambda}{2} \|w\|_2^2 \Rightarrow \text{Ridge Regression}$$

extra sample 1  $\Rightarrow [\sqrt{\lambda}, 0, \dots, 0]^T$   $\leadsto \lambda w_1^2$

extra sample  $N \Rightarrow [\dots, \sqrt{\lambda}]^T$   $\leadsto \lambda w_M^2$

$\leadsto \dots + w_M^2$

$$\begin{bmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_N) \\ \sqrt{\lambda} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda} \end{bmatrix} w - \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} (\Phi(x_1)^T w - y_1) \\ \vdots \\ (\Phi(x_N)^T w - y_N) \\ w_1 \sqrt{\lambda} \\ \vdots \\ w_M \sqrt{\lambda} \end{bmatrix}$$

$\Rightarrow E = \text{cost function Ridge Regression}$

$$= \sum_i (\Phi(x_i)^T w - y_i)^2 + \lambda w_1^2 + \lambda w_2^2 + \dots + w_M^2 \lambda$$

**Problem 4:** It turns out that the conjugate prior for the situation when we have an unknown mean and unknown precision is a normal-gamma distribution (See section 2.3.6 in Bishop). This is also true when we have a conditional Gaussian distribution of the linear regression model. This means that if our likelihood is as follows:

$$p(\mathbf{y} | \Phi, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | \mathbf{w}^T \phi(x_i), \beta^{-1})$$

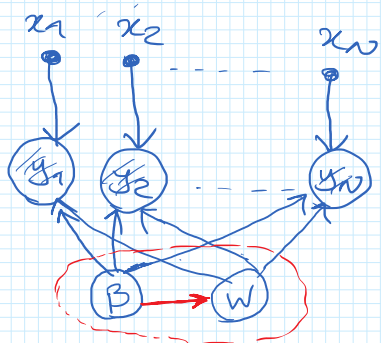
Then the conjugate prior for both  $\mathbf{w}$  and  $\beta$  is

$$p(\mathbf{w}, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \beta^{-1} \mathbf{S}_0) \text{Gamma}(\beta | a_0, b_0)$$

Show that the posterior distribution takes the same form as the prior, i.e.

$$p(\mathbf{w}, \beta | \mathcal{D}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \beta^{-1} \mathbf{S}_N) \text{Gamma}(\beta | a_N, b_N)$$

Also be sure to give the expressions for  $\mathbf{m}_N$ ,  $\mathbf{S}_N$ ,  $a_N$ , and  $b_N$ .



Normal Distribution

$$\mathcal{N}(\mathbf{x} | \mu, \Sigma) = (2\pi)^{-\frac{M}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

$$\mathcal{N}_{\text{can}}(\mathbf{x} | \eta, \Lambda) \propto \exp\left\{ \eta^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \Lambda \mathbf{x} \right\}$$

for example  $\Rightarrow p(\mathbf{w}) \propto \exp\left\{ \underbrace{[1.0, 0.5, 0.3]}_{\eta} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} - \frac{1}{2} \mathbf{x}^T \underbrace{\mathbf{I}}_{\Lambda} \mathbf{x} \right\}$

How to convert canonical form to standard form  
 $\mu = \Lambda^{-1} \eta$        $\Sigma = \Lambda^{-1}$

$$p(\mathbf{w} | \mathcal{D}, \beta) = \frac{p(\mathbf{w}, \beta | \mathcal{D})}{p(\beta | \mathcal{D})}$$

$$\propto p(\mathbf{w}, \beta | \mathcal{D})$$

$$\underbrace{p(\mathbf{w} | \beta) p(\beta)}_{= p(\mathbf{w}, \beta)}$$

$$\propto p(\mathcal{D} | \mathbf{w}, \beta) p(\mathbf{w}, \beta)$$

$\eta^T \mathbf{x} \rightarrow \text{neglect}$

$\mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 \rightarrow \text{neglect}$

$$\propto \exp\left\{ -\frac{\beta}{2} (\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w}) \right\} \propto \beta^{\frac{N}{2}} \exp\left\{ -\frac{\beta}{2} \mathbf{y}^T \mathbf{I} \mathbf{y} \right\}$$

$$\propto \exp\left\{-\frac{\beta}{2}(\underline{y}-\Phi\underline{w})^T(\underline{y}-\Phi\underline{w})\right\} \exp\left\{-\frac{\beta}{2}(\underline{w}-\underline{m}_0)^T \underline{S}_0^{-1}(\underline{w}-\underline{m}_0)\right\}$$

$\gamma \rightarrow \text{neglect}$        $m_0^T S_0 m_0 \rightarrow \text{neglect}$

$$\propto \exp\left\{-\frac{\beta}{2}\left[\underline{w}^T \Phi^T \Phi \underline{w} - 2\underline{y}^T \Phi \underline{w}\right]\right\} \exp\left\{-\frac{\beta}{2}\left[\underline{w}^T \underline{S}_0^{-1} \underline{w} - 2\underline{m}_0^T \underline{S}_0^{-1} \underline{w}\right]\right\}$$

$$\propto \exp\left\{\left[\beta \Phi^T \underline{y} + \beta \underline{S}_0^{-1} \underline{m}_0\right]^T \underline{w} - \frac{1}{2} \underline{w}^T \left[\beta \Phi \Phi^T + \beta \underline{S}_0^{-1}\right] \underline{w}\right\}$$

$$\mathcal{N}_{\text{can}}(\underline{w} \mid \beta \Phi^T \underline{y} + \beta \underline{S}_0^{-1} \underline{m}_0, \beta \Phi \Phi^T + \beta \underline{S}_0^{-1})$$

convenient  $\rightarrow \mathcal{N}\left(\underline{w} \mid \left(\Phi \Phi^T + \underline{S}_0^{-1}\right)^{-1} (\Phi^T \underline{y} + \underline{S}_0^{-1} \underline{m}_0), \frac{1}{\beta} \left(\Phi \Phi^T + \underline{S}_0^{-1}\right)^{-1}\right)$

$\Rightarrow \underline{S}_N = \Phi \Phi^T + \underline{S}_0^{-1}, \underline{m}_N = \underline{S}_N (\Phi^T \underline{y} + \underline{S}_0^{-1} \underline{m}_0)$

$$P(\beta \mid D) = \frac{P(\underline{w}, \beta \mid D)}{P(\underline{w} \mid \beta, D)} = \frac{P(\underline{w}, \beta) P(D \mid \underline{w}, \beta)}{P(\underline{w} \mid \beta, D)}$$

$$= \frac{\mathcal{N}(\underline{w} \mid \underline{m}_0, \underline{S}_0^{-1}) \text{Gamma}(\beta \mid a_0, b_0) \prod_n \mathcal{N}(y_n \mid x_n, \underline{w})}{\mathcal{N}(\underline{w} \mid \underline{m}_N, \underline{S}_N^{-1})}$$

$$= \frac{|\underline{S}_0^{-1}|^{-\frac{1}{2}} \exp\left\{-\frac{\beta}{2}(\underline{w}-\underline{m}_0)^T \underline{S}_0^{-1}(\underline{w}-\underline{m}_0)\right\} e^{-b_0 \beta} \beta^{a_0-1}}{|\underline{S}_N^{-1}|^{-\frac{1}{2}} \exp\left\{-\frac{\beta}{2}(\underline{w}-\underline{m}_N)^T \underline{S}_N^{-1}(\underline{w}-\underline{m}_N)\right\}} \times \beta^{\frac{N}{2}}$$

$$|\underline{S}_N^{-1}|^{-\frac{1}{2}} \exp\left\{-\frac{\beta}{2}(\underline{w}-\underline{m}_N)^T \underline{S}_N^{-1}(\underline{w}-\underline{m}_N)\right\}$$

$$\frac{M}{R^2} \times \frac{a_0-1}{R} \times \frac{N}{R^2}$$

$$c \quad R \Gamma(1) \quad \dots \quad T-1 \quad \dots$$

$$= \frac{\beta^{\frac{M}{2}}}{\beta^{\frac{M}{2}}} \times \beta^{\alpha_0 - 1} \times \beta^{\frac{N}{2}} \exp \left\{ -\beta \left[ \frac{1}{2} (w - m_0)^T \bar{S}_0^{-1} (w - m_0) + b_0 + \frac{1}{2} (y - \Phi w)^T (y - \Phi w) - \frac{1}{2} (w - m_N)^T S_N^{-1} (w - m_N) \right] \right\}$$

-  $p(\beta | D)$  is not a function of  $w$ , but we see some terms contain  $w$ .

According to the values for  $m_N, S_N$  (computed in previous part), the above expression becomes independent from  $w$ .

Because:

$$\text{The terms containing } w = \underbrace{w^T (\bar{S}_0^{-1} + \Phi^T \Phi - S_N^{-1}) w}_{0} - 2 \underbrace{(m_0^T \bar{S}_0^{-1} + y^T \Phi - m_N^T S_N^{-1}) w}_{0}$$

$$= \beta^{\alpha_0 + \frac{N}{2} - 1} \exp \left\{ -\beta \left[ \frac{1}{2} m_0^T \bar{S}_0^{-1} m_0 - \frac{1}{2} m_N^T S_N^{-1} m_N + b_0 + \frac{1}{2} y^T y \right] \right\}$$

$$= \text{Gamma} \left( \beta \mid \alpha_0 + \frac{N}{2}, b_0 + \frac{1}{2} (m_0^T \bar{S}_0^{-1} m_0 - m_N^T S_N^{-1} m_N + y^T y) \right)$$