



Take Home Solution

Prof. Krishna R. Pattipati

**Dept. of Electrical and Computer Engineering
University of Connecticut**

Contact: krishna@engr.uconn.edu (860) 486-2890

ECE 6141
Neural Networks for Classification and Optimization



Problem 1:

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

$$p(\mu) = \lambda \exp(-\lambda\mu); \mu \geq 0; \lambda > 0$$

$$p(\mu | \{x_n\}_{n=1}^N) \propto p(\{x_n\}_{n=1}^N | \mu) p(\mu) = L(\mu)$$

$$\min_{\mu \geq 0} J = -\ln L(\mu) = \frac{1}{2\sigma^2} \left(\sum_{n=1}^N (x_n - \mu)^2 \right) + \lambda\mu + \text{constant}$$

$$\frac{\partial J}{\partial \mu} = \lambda - \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu)$$

$$\Rightarrow \hat{\mu}_{MAP} = \max\left(\frac{\sum_{n=1}^N x_n - \sigma^2 \lambda}{N}, 0\right)$$

$$\frac{\partial^2 J}{\partial \mu^2} = H = \frac{N}{\sigma^2}$$

so, Laplace approximation of posterior : $p(\mu | \{x_n\}_{n=1}^N) \sim N(\hat{\mu}_{MAP}, \frac{\sigma^2}{N})$



Problem 2

$$a) \underline{x} \in R^n$$

$$\hat{\underline{\mu}}_i = \frac{1}{n_i} \sum_{k \in C_i} \underline{x}_k; i = 1, 2; C_i = \text{samples from class } i; |C_i| = n_i$$

$$\hat{\Sigma} = \frac{1}{N-2} \left(\sum_{i=1}^2 \sum_{k \in C_i} (\underline{x}_k - \hat{\underline{\mu}}_i)(\underline{x}_k - \hat{\underline{\mu}}_i)^T \right)$$

$$\hat{\pi}_i = \frac{n_i}{N}; i = 1, 2$$

$$g(\underline{x}) = g_2(\underline{x}) - g_1(\underline{x})$$

$$= (\hat{\underline{\mu}}_2 - \hat{\underline{\mu}}_1)^T \hat{\Sigma}^{-1} \underline{x} - \left[\frac{1}{2} (\hat{\underline{\mu}}_2^T \hat{\Sigma}^{-1} \hat{\underline{\mu}}_2 - \hat{\underline{\mu}}_1^T \hat{\Sigma}^{-1} \hat{\underline{\mu}}_1) - \ln \frac{\hat{\pi}_2}{\hat{\pi}_1} \right]$$

$$\Rightarrow \underline{x}^T \hat{\Sigma}^{-1} (\hat{\underline{\mu}}_2 - \hat{\underline{\mu}}_1) > \frac{1}{2} \hat{\underline{\mu}}_2^T \hat{\Sigma}^{-1} \hat{\underline{\mu}}_2 - \frac{1}{2} \hat{\underline{\mu}}_1^T \hat{\Sigma}^{-1} \hat{\underline{\mu}}_1 + \ln \frac{n_1}{n_2}$$



Problem 1

$$b) J = \sum_{i=1}^2 \sum_{k \in C_i} (z_i - \underline{w}_a^T \underline{x}_{ak})^2; \underline{x}_{ak} = [1, \underline{x}_k^T]^T$$

$$\left(\sum_{i=1}^2 \sum_{k \in C_i} \underline{x}_{ak} \underline{x}_{ak}^T \right) \underline{\hat{w}}_a = \left(\sum_{i=1}^2 z_i \left(\sum_{k \in C_i} \underline{x}_{ak} \right) \right) = N(\underline{\hat{\mu}}_{2a} - \underline{\hat{\mu}}_{1a})$$

$$\text{where } \underline{\hat{\mu}}_{ia} = [1, \underline{\hat{\mu}}_i^T]^T$$

$$\text{For } \underline{w}_0 \text{ term: } N \underline{\hat{w}}_0 = -\underline{\hat{w}}^T \sum_{i=1}^2 \sum_{k \in C_i} \underline{x}_k = -\underline{\hat{w}}^T \sum_{i=1}^2 n_i \underline{\mu}_i \Rightarrow \underline{\hat{w}}_0 = -\underline{\hat{w}}^T \underline{\mu}$$

$$\text{For } \underline{w} \text{ term: } N \underline{\hat{w}}_0 + \left[(N-2) \hat{\Sigma} + \sum_{i=1}^2 n_i \underline{\mu}_i \underline{\mu}_i^T \right] \underline{\hat{w}} = N(\underline{\hat{\mu}}_2 - \underline{\hat{\mu}}_1)$$

$$\Rightarrow \left[(N-2) \hat{\Sigma} + \sum_{i=1}^2 n_i \underline{\mu}_i \underline{\mu}_i^T - N \underline{\mu} \underline{\mu}^T \right] \underline{\hat{w}} = N(\underline{\hat{\mu}}_2 - \underline{\hat{\mu}}_1)$$

$$\left[(N-2) \hat{\Sigma} + \sum_{i=1}^2 n_i \underline{\mu}_i \underline{\mu}_i^T - \frac{1}{N} \left(\sum_{i=1}^2 n_i \underline{\mu}_i \right) \left(\sum_{j=1}^2 n_j \underline{\mu}_j \right)^T \right] \underline{\hat{w}} = N(\underline{\hat{\mu}}_2 - \underline{\hat{\mu}}_1)$$

$$\left((N-2) \hat{\Sigma} + \frac{n_1 n_2}{N} \hat{\Sigma}_B \right) \underline{w} = N(\underline{\hat{\mu}}_2 - \underline{\hat{\mu}}_1)$$

where

$$\hat{\Sigma}_B = (\underline{\hat{\mu}}_2 - \underline{\hat{\mu}}_1) (\underline{\hat{\mu}}_2 - \underline{\hat{\mu}}_1)^T$$



Problem 1

$$c) (\hat{\Sigma}_B) \hat{w} = (\hat{\mu}_2 - \hat{\mu}_1) (\hat{\mu}_2 - \hat{\mu}_1)^T \hat{w} \propto (\hat{\mu}_2 - \hat{\mu}_1)$$

$$\text{so, } \hat{w} \propto \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1)$$

d) for any distinct coding, $z_1 \neq z_2$

$$\text{RHS: } \left(\sum_{i=1}^2 z_i \left(\sum_{k \in C_i} x_{ak} \right) \right) = \sum_{i=1}^2 z_i n_i \mu_{ai} = \begin{bmatrix} \sum_{i=1}^2 z_i n_i \\ \sum_{i=1}^2 z_i n_i \mu_i \end{bmatrix}$$

$$\text{For } w_0 \text{ term: } N \hat{w}_0 = -\hat{w}^T \sum_{i=1}^2 \sum_{k \in C_i} x_k + \sum_{i=1}^2 z_i n_i = -\hat{w}^T \sum_{i=1}^2 n_i \mu_i + \frac{1}{N} \sum_{i=1}^2 z_i n_i \Rightarrow \hat{w}_0 = \sum_{i=1}^2 \hat{\pi}_i (z_i - \hat{w}^T \mu_i) = \sum_{i=1}^2 \hat{\pi}_i z_i - \hat{w}^T \mu$$

$$\text{For } w \text{ term: } N \underline{w} \hat{w}_0 + \left[(N-2) \hat{\Sigma} + \sum_{i=1}^2 n_i \underline{\mu}_i \underline{\mu}_i^T \right] \hat{w} = \sum_{i=1}^2 z_i n_i \underline{\mu}_i$$

$$\Rightarrow \left[(N-2) \hat{\Sigma} + \sum_{i=1}^2 n_i \underline{\mu}_i \underline{\mu}_i^T - N \underline{\mu} \underline{\mu}^T \right] \hat{w} = \sum_{i=1}^2 z_i n_i \underline{\mu}_i - N \left(\sum_{i=1}^2 \hat{\pi}_i z_i \right) \underline{\mu} = \sum_{i=1}^2 n_i \underline{\mu}_i z_i - \frac{1}{N} \left(\sum_{i=1}^2 n_i z_i \right) \left(\sum_{j=1}^2 n_j \underline{\mu}_j \right)$$

$$\left[(N-2) \hat{\Sigma} + \sum_{i=1}^2 n_i \underline{\mu}_i \underline{\mu}_i^T - \frac{1}{N} \left(\sum_{i=1}^2 n_i \underline{\mu}_i \right) \left(\sum_{j=1}^2 n_j \underline{\mu}_j \right)^T \right] \hat{w} = \frac{n_1 n_2}{N} (z_2 - z_1) (\hat{\mu}_2 - \hat{\mu}_1)$$

$$\left((N-2) \hat{\Sigma} + \frac{n_1 n_2}{N} \hat{\Sigma}_B \right) \underline{w} = \frac{n_1 n_2}{N} (z_2 - z_1) (\hat{\mu}_2 - \hat{\mu}_1)$$

where

$$\hat{\Sigma}_B = (\hat{\mu}_2 - \hat{\mu}_1) (\hat{\mu}_2 - \hat{\mu}_1)^T$$

$$\text{Evidently, when } z_2 = \frac{N}{n_2} \text{ and } z_1 = -\frac{N}{n_1}, \frac{n_1 n_2}{N} (z_2 - z_1) = N$$

$$e) \text{ Decision rule: } \sum_{i=1}^2 \hat{\pi}_i z_i + \hat{w}^T (x_i - \underline{\mu}) \begin{matrix} \text{class 2} \\ \text{class 1} \end{matrix} > 0 \Rightarrow \text{Need } \sum_{i=1}^2 \hat{\pi}_i z_i = 0$$

$$\Rightarrow \sum_{i=1}^2 n_i z_i = 0 \Rightarrow \text{possible codings: } z_1 = -\frac{N}{n_1}, z_2 = \frac{N}{n_2} \dots \text{no need for equal number of samples}$$

$$z_1 = -1, z_2 = 1 \text{ \& } n_1 = n_2 \text{ etc.}$$



Problem 3

$$z \sim N(z; \mu, \sigma^2)$$

$$z = \mu + \sigma t \text{ where } t \sim N(0,1)$$

$$E\left[z \mid z \geq c\right] = E\left[\mu + \sigma t \mid t \geq \frac{c - \mu}{\sigma}\right]$$

$$\text{Recall } \phi(t \mid t \geq \frac{c - \mu}{\sigma}) = \frac{\phi(t)}{1 - \Phi(\frac{c - \mu}{\sigma})}$$

$$\text{so, } E\left[\mu + \sigma t \mid t \geq \frac{c - \mu}{\sigma}\right] = \mu + \frac{1}{1 - \Phi(\frac{c - \mu}{\sigma})} \frac{\sigma}{\sqrt{2\pi}} \int_{\frac{c - \mu}{\sigma}}^{\infty} t e^{-t^2/2} dt$$

$$= \mu - \frac{\sigma}{1 - \Phi(\frac{c - \mu}{\sigma})} \int_{\frac{c - \mu}{\sigma}}^{\infty} \phi'(t) dt = \mu - \frac{\sigma}{1 - \Phi(\frac{c - \mu}{\sigma})} \phi(t) \Big|_{\frac{c - \mu}{\sigma}}^{\infty} = \mu + \frac{\sigma \phi(\frac{c - \mu}{\sigma})}{1 - \Phi(\frac{c - \mu}{\sigma})} = \mu + \sigma H\left(\frac{c - \mu}{\sigma}\right)$$

$$\text{where } H(u) = \frac{\phi(u)}{1 - \Phi(u)}$$

$$E\left[z^2 \mid z \geq c\right] = E\left[(\mu + \sigma t)^2 \mid t \geq \frac{c - \mu}{\sigma}\right] = \mu^2 + 2\sigma\mu H\left(\frac{c - \mu}{\sigma}\right) + \sigma^2 E[t^2 \mid t \geq \frac{c - \mu}{\sigma}]$$

$$E[t^2 \mid t \geq \frac{c - \mu}{\sigma}] = \left[\frac{1}{1 - \Phi(\frac{c - \mu}{\sigma})} \right] \left[\int_{\frac{c - \mu}{\sigma}}^{\infty} t^2 \phi(t) dt \right]$$

$$\int_{\frac{c - \mu}{\sigma}}^{\infty} t^2 \phi(t) dt = - \int_{\frac{c - \mu}{\sigma}}^{\infty} t \phi'(t) dt = \left(\frac{c - \mu}{\sigma}\right) \phi\left(\frac{c - \mu}{\sigma}\right) + \int_{\frac{c - \mu}{\sigma}}^{\infty} \phi(t) dt = \left(\frac{c - \mu}{\sigma}\right) \phi\left(\frac{c - \mu}{\sigma}\right) + 1 - \Phi\left(\frac{c - \mu}{\sigma}\right)$$

$$\text{so, } E\left[z^2 \mid z \geq c\right] = \mu^2 + 2\sigma\mu H\left(\frac{c - \mu}{\sigma}\right) + \sigma(c - \mu) H\left(\frac{c - \mu}{\sigma}\right) + \sigma^2$$
$$= \mu^2 + \sigma^2 + \sigma(c + \mu) H\left(\frac{c - \mu}{\sigma}\right)$$



Problem 4:

$$\underline{w}^{(n+1)} = \underline{w}^{(n)} + \eta(z^n - \underline{w}^{(n)T} \underline{x}^n) \underline{x}^n = \underline{w}^{(n)} + \eta \underbrace{(z^n - \underline{w}^{*T} \underline{x}^n)}_{e^{*n}} - (\underline{w}^{(n)} - \underline{w}^*)^T \underline{x}^n \underline{x}^n$$

$$\underline{v}^{(n+1)} = [I - \eta \underline{x}^n \underline{x}^{nT}] \underline{v}^{(n)} + \eta e^{*n} \underline{x}^n; \underline{v}^{(n)} = \underline{w}^{(n)} - \underline{w}^*$$

Let $\Sigma_n = E\{\underline{v}^{(n)} \underline{v}^{(n)T}\}$; $R_x = E[\underline{x}^n \underline{x}^{nT}] \sim$ **Correlation matrix of data**; $E[(e^{*n})^2] = \sigma_e^2$

$$\begin{aligned} \Sigma_{n+1} = E\{\underline{v}^{(n+1)} \underline{v}^{(n+1)T}\} &= E\left\{ [I - \eta \underline{x}^n \underline{x}^{nT}] \underline{v}^{(n)} \underline{v}^{(n)T} [I - \eta \underline{x}^n \underline{x}^{nT}]^T \right\} + \eta E\{e^{*n} [I - \eta \underline{x}^n \underline{x}^{nT}] \underline{v}^{(n)} \underline{x}^{nT}\} + \\ &+ \eta E\{e^{*n} \underline{x}^n \underline{v}^{(n)T} [I - \eta \underline{x}^n \underline{x}^{nT}]^T\} + \eta^2 E\{e^{*n2} \underline{x}^n \underline{x}^{nT}\} \end{aligned}$$

Using LMS assumption and the orthogonality of error and the weight estimate, we have

$$\begin{aligned} \Sigma_{n+1} &= \Sigma_n - \eta R_x \Sigma_n - \eta \Sigma_n R_x + \eta^2 E\{\underline{x}^n \underline{x}^{nT} \underline{v}^{(n)} \underline{v}^{(n)T} \underline{x}^n \underline{x}^{nT}\} + \eta^2 E\{(e^{*n})^2 \underline{x}^n \underline{x}^{nT}\} \\ &= \Sigma_n - \eta R_x \Sigma_n - \eta \Sigma_n R_x + 2\eta^2 R_x \Sigma_n R_x + \eta^2 R_x \text{tr}\{\Sigma_n R_x\} + \eta^2 \sigma_e^2 R_x \end{aligned}$$

Note: Use $E\{x_1 x_2 x_3 x_4\} = E\{x_1 x_2\} E\{x_3 x_4\} + E\{x_1 x_3\} E\{x_2 x_4\} + E\{x_1 x_4\} E\{x_2 x_3\}$

$$\begin{aligned} \text{Vector case: } E\{\underline{x}^n \underline{x}^{nT} \underline{v}^{(n)} \underline{v}^{(n)T} \underline{x}^n \underline{x}^{nT}\} &= E\{\underline{x}^n \underline{x}^{nT}\} E\{\underline{v}^{(n)} \underline{v}^{(n)T} \underline{x}^n \underline{x}^{nT}\} + E\{\underline{x}^n \underline{x}^{nT} \underline{v}^{(n)} \underline{v}^{(n)T}\} E\{\underline{x}^n \underline{x}^{nT}\} \\ &+ E\{\underline{x}^{nT} \underline{v}^{(n)} \underline{v}^{(n)T} \underline{x}^n\} E\{\underline{x}^n \underline{x}^{nT}\} \end{aligned}$$

Let $R_x = Q \Lambda_x Q^T$ and let $\hat{\Sigma}_{n+1} = Q^T \Sigma_{n+1} Q$

$$\begin{aligned} Q^T \Sigma_{n+1} Q &= Q^T \Sigma_n Q - \eta Q^T R_x Q Q^T \Sigma_n Q - \eta Q^T \Sigma_n Q Q^T R_x Q + 2\eta^2 Q^T R_x Q Q^T \Sigma_n Q Q^T R_x Q \\ &+ \eta^2 R_x \text{tr}\{Q^T \Sigma_n Q Q^T R_x Q\} + \eta^2 \sigma_e^2 Q^T R_x Q \end{aligned}$$

$$\Rightarrow \hat{\Sigma}_{n+1} = \hat{\Sigma}_n - \eta \Lambda_x \hat{\Sigma}_n - \eta \hat{\Sigma}_n \Lambda_x + 2\eta^2 \Lambda_x \hat{\Sigma}_n \Lambda_x + \eta^2 \Lambda_x \text{tr}\{\hat{\Sigma}_n \Lambda_x\} + \eta^2 \sigma_e^2 \Lambda_x$$



Problem 4:

$$\hat{\Sigma}_{n+1} = \hat{\Sigma}_n - \eta \Lambda_x \hat{\Sigma}_n - \eta \hat{\Sigma}_n \Lambda_x + 2\eta^2 \Lambda_x \hat{\Sigma}_n \Lambda_x + \eta^2 \Lambda_x \text{tr}\{\hat{\Sigma}_n \Lambda_x\} + \eta^2 \sigma_e^2 \Lambda_x$$

Now consider the diagonal elements of $\hat{\Sigma}_{n+1}$ and represent them as a vector \underline{s}_{n+1}

$$\underline{s}_{n+1} = (I_{p+1} - 2\eta \Lambda_x + 2\eta^2 \Lambda_x^2 + \eta^2 \underline{\lambda} \underline{\lambda}^T) \underline{s}_n + \eta^2 \sigma_e^2 \underline{\lambda}$$

where $\underline{\lambda} = [\lambda_1 \quad \lambda_2 \quad \dots \quad \lambda_{p+1}]^T$

Need eigen values of $(I_{p+1} - 2\eta \Lambda_x + 2\eta^2 \Lambda_x^2 + \eta^2 \underline{\lambda} \underline{\lambda}^T)$ inside unit circle.

$$\Rightarrow |\lambda_i \underbrace{\{(I_{p+1} - \eta \Lambda_x)^2 + \eta^2 (\Lambda_x^2 + \underline{\lambda} \underline{\lambda}^T)\}}_{\Phi}| < 1 \forall i = 1, 2, \dots, p+1$$

$$\Phi_{ij} = \begin{cases} (1 - \eta \lambda_i)^2 + 2\eta^2 \lambda_i^2; i = j \\ \eta^2 \lambda_i \lambda_j; i \neq j \end{cases}$$

$$\|\Phi\|_{\infty} < 1 \Rightarrow \max_i \left\{ (1 - \eta \lambda_i)^2 + \eta^2 \lambda_i \sum_{j=1}^{p+1} \lambda_j \right\} < 1$$

$$\Rightarrow -2\eta \lambda_i + 2\eta^2 \lambda_i^2 + \eta^2 \lambda_i \text{Tr}(R_x) < 0 \forall i$$

$$\Rightarrow 2\eta \lambda_i + \eta \text{Tr}(R_x) < 2 \text{ \& } \eta > 0$$

$$\Rightarrow 0 < \eta < \min_i \frac{2}{\text{tr}(R_x) + 2\lambda_i} < \frac{2}{\text{tr}(R_x) + 2\lambda_{\max}} < \frac{2}{\text{tr}(R_x)} \quad (\text{in practice: } 0 < \eta < \frac{2}{\text{tr}(R_x) + 2\lambda_{\max}})$$

Sum min g over all i

$$2\eta \text{Tr}(R_x) + (p+1)\eta \text{Tr}(R_x) < 2(p+1) \text{ \& } \eta > 0 \Rightarrow 0 < \eta < \frac{2}{\text{Tr}(R_x) [1 + \frac{2}{p+1}]} < \frac{2}{\text{tr}(R_x)}$$

Let $f(\lambda) = -2\eta\lambda + 2\eta^2\lambda^2 + \eta^2\lambda\text{Tr}(R_x)$

Max of f at: $\lambda^* = \frac{2 - \eta\text{Tr}(R_x)}{4\eta}$

$$f(\lambda^*) = -\frac{\eta}{2} \left(1 - \frac{\eta}{2} \text{Tr}(R_x)\right)^2 \Rightarrow 0 < \eta < \frac{2}{\text{Tr}(R_x)}$$



Problem 5

Consider a general regularized least squares regression problem.

$$J = \frac{1}{N} \|\underline{z} - X\underline{w}\|_2^2 + \frac{\lambda}{N} \underline{w}^T \Gamma^T \Gamma \underline{w}; \underline{z} \in R^N; X \in R^{N \times (p+1)}$$

where $\underline{z} = X\underline{w} + \underline{v}$; $v_n \sim N(0, \sigma^2) \forall n = 1, 2, \dots, N$

Let $\hat{\underline{w}}(0, \Gamma) = (X^T X)^{-1} X^T \underline{z}$, least squares solution when $\lambda = 0$.

$$(a) \nabla_{\underline{w}} J = \frac{2}{N} [-X^T (\underline{z} - X\underline{w}) + \lambda \Gamma^T \Gamma \underline{w}] = \underline{0}$$

$$\begin{aligned} \Rightarrow \hat{\underline{w}}(\lambda, \Gamma) &= (X^T X + \lambda \Gamma^T \Gamma)^{-1} X^T \underline{z} = (X^T X + \lambda \Gamma^T \Gamma)^{-1} X^T (X\underline{w} + \underline{v}) \\ &= (X^T X + \lambda \Gamma^T \Gamma)^{-1} [X^T X \underline{w} + X^T \underline{v}] \\ &= (X^T X + \lambda \Gamma^T \Gamma)^{-1} [(X^T X + \lambda \Gamma^T \Gamma) \underline{w} - \lambda \Gamma^T \Gamma \underline{w} + X^T \underline{v}] \\ &= \underline{w} - \lambda (X^T X + \lambda \Gamma^T \Gamma)^{-1} \Gamma^T \Gamma \underline{w} + (X^T X + \lambda \Gamma^T \Gamma)^{-1} X^T \underline{v} \end{aligned}$$

$$\Gamma = I_{p+1} \Rightarrow \hat{\underline{w}}(\lambda, I_{p+1}) = \underline{w} - \lambda (X^T X + \lambda I_{p+1})^{-1} \underline{w} + (X^T X + \lambda I_{p+1})^{-1} X^T \underline{v}$$

$$\Gamma = X \Rightarrow \hat{\underline{w}}(\lambda, X) = \frac{\underline{w}}{1+\lambda} + \frac{1}{1+\lambda} (X^T X)^{-1} X^T \underline{v} = \frac{\hat{\underline{w}}(0, X)}{1+\lambda} = \frac{\hat{\underline{w}}(0, \Gamma)}{1+\lambda} = \frac{\hat{\underline{w}}(0)}{1+\lambda}$$

It scales the least squares solution by $1/(1+\lambda)$



Problem 5

(b)

$$\begin{aligned}\underline{b} &= \underline{w} - E\{\hat{\underline{w}}(\lambda, \Gamma)\} = \underline{w} - \underline{w} + \lambda(X^T X + \lambda\Gamma^T \Gamma)^{-1} \Gamma^T \Gamma \underline{w} - E\{(X^T X + \lambda\Gamma^T \Gamma)^{-1} X^T \underline{v}\} \\ &= \lambda(X^T X + \lambda\Gamma^T \Gamma)^{-1} \Gamma^T \Gamma \underline{w}\end{aligned}$$

$$\Gamma = I_{p+1} \Rightarrow \underline{b} = \lambda(X^T X + \lambda I_{p+1})^{-1} \underline{w}$$

$$\Gamma = X \Rightarrow \underline{b} = \frac{\lambda}{1+\lambda} \underline{w}$$

(c)

$$\begin{aligned}r &= z - \hat{z} = \underline{x}^T \underline{w} + v - \underline{x}^T \hat{\underline{w}}(\lambda, \Gamma) \\ &= \lambda \underline{x}^T (X^T X + \lambda\Gamma^T \Gamma)^{-1} \Gamma^T \Gamma \underline{w} + v - \underline{x}^T (X^T X + \lambda\Gamma^T \Gamma)^{-1} X^T \underline{v}\end{aligned}$$

$$\Gamma = I_{p+1} \Rightarrow \underline{r} = \lambda \underline{x}^T (X^T X + \lambda I_{p+1})^{-1} \underline{w} + v - \underline{x}^T (X^T X + \lambda I_{p+1})^{-1} X^T \underline{v}$$

$$\Gamma = X \Rightarrow \underline{r} = \frac{\lambda}{1+\lambda} (\underline{x}^T \underline{w} + v) + \frac{1}{1+\lambda} [v - \underline{x}^T (X^T X)^{-1} X^T \underline{v}]$$



Problem 5

(d)

$$r = \lambda \underline{x}^T (X^T X + \lambda \Gamma^T \Gamma)^{-1} \Gamma^T \Gamma \underline{w} + v - \underline{x}^T (X^T X + \lambda \Gamma^T \Gamma)^{-1} X^T \underline{v}$$

$$\text{bias}^2(\lambda, \Gamma) = E(r)^2 \approx \lambda^2 \underline{w}^T \Gamma^T \Gamma (N \Sigma_x + \lambda \Gamma^T \Gamma)^{-1} \Sigma_x (N \Sigma_x + \lambda \Gamma^T \Gamma)^{-1} \Gamma^T \Gamma \underline{w}$$

$$\text{where } \Sigma_x = E_{\underline{x}}(\underline{x} \underline{x}^T) \approx \frac{X^T X}{N}$$

$$\Gamma = I_{p+1} \Rightarrow \text{bias}^2(\lambda, I_{p+1}) = \lambda^2 \underline{w}^T (N \Sigma_x + \lambda I_{p+1})^{-1} \Sigma_x (N \Sigma_x + \lambda I_{p+1})^{-1} \underline{w}$$

$$\text{If } \Sigma_x = I_{p+1} \text{ also, } \text{bias}^2(\lambda, I_{p+1}) = \left(\frac{\lambda}{N + \lambda} \right)^2 \underline{w}^T \underline{w}$$

$$\Gamma = X \Rightarrow \text{bias}^2(\lambda, X) = \left(\frac{\lambda}{\lambda + 1} \right)^2 \underline{w}^T \Sigma_x \underline{w}$$

$$\text{If } \Sigma_x = I_{p+1} \text{ also, } \text{bias}^2(\lambda, X) = \left(\frac{\lambda}{\lambda + 1} \right)^2 \underline{w}^T \underline{w}$$



Problem 5

(e)

$$\begin{aligned}\text{var}(\lambda, \Gamma) &= E\{[r - E(r)]^2\} = \\ &\sigma^2 + [E_{\underline{x}, \underline{v}}\{\underline{x}^T (X^T X + \lambda \Gamma^T \Gamma)^{-1} X^T \underline{v} \underline{v}^T X (X^T X + \lambda \Gamma^T \Gamma)^{-1} \underline{x}\}] \\ &\approx \sigma^2 (1 + N \cdot \text{tr}([\Sigma_x (N \Sigma_x + \lambda \Gamma \Gamma^T)^{-1}]^2))\end{aligned}$$

$$\text{where } \Sigma_x = E_{\underline{x}}(\underline{x} \underline{x}^T) \approx \frac{X^T X}{N}$$

$$\Gamma = I_{p+1} \Rightarrow \text{var}(\lambda, I_{p+1}) = \sigma^2 (1 + N \cdot \text{tr}([\Sigma_x (N \Sigma_x + \lambda I_{p+1})^{-1}]^2))$$

$$\text{If } \Sigma_x = I_{p+1} \text{ also, } \text{var}(\lambda, X) = \sigma^2 \left[1 + \frac{N(p+1)}{(N+\lambda)^2}\right]$$

$$\Gamma = X \Rightarrow \text{var}(\lambda, X) = \sigma^2 \left[1 + \frac{(p+1)}{N(1+\lambda)^2}\right] \dots \text{independent of } \Sigma_x$$



Problem 5

(f)

$$\Gamma = I_{p+1}; \Sigma_x = I_{p+1}$$

$$MSE = bias^2(\lambda, I_{p+1}) + var(\lambda, I_{p+1}) = \left(\frac{\lambda}{N + \lambda} \right)^2 \underline{w}^T \underline{w} + \sigma^2 \left[1 + \frac{N(p+1)}{(N + \lambda)^2} \right]$$

$$\frac{dMSE}{d\lambda} = \left(\frac{2\lambda}{(N + \lambda)^2} - \frac{2\lambda^2}{(N + \lambda)^3} \right) \underline{w}^T \underline{w} - \frac{2\sigma^2 N(p+1)}{(N + \lambda)^3} = 0$$

$$\Rightarrow \frac{2N\lambda}{(N + \lambda)^3} \underline{w}^T \underline{w} - \frac{2\sigma^2 N(p+1)}{(N + \lambda)^3} = 0 \Rightarrow \lambda^* = \frac{\sigma^2(p+1)}{\underline{w}^T \underline{w}}$$

$$\text{Note: } \frac{dMSE}{d\lambda} = 2N \underline{w}^T \underline{w} \left(\frac{\lambda - \lambda^*}{(N + \lambda)^3} \right)$$

$$\Rightarrow \frac{d^2MSE}{d\lambda^2} = 2N \underline{w}^T \underline{w} \left(\frac{1}{(N + \lambda)^3} - \frac{3(\lambda - \lambda^*)}{(N + \lambda)^4} \right) = \frac{2N \underline{w}^T \underline{w}}{(N + \lambda^*)^3} > 0 \text{ at } \lambda = \lambda^*$$



Problem 5

(f)

$$\Gamma = X; \Sigma_x = I_{p+1}$$

$$MSE = bias^2(\lambda, X) + var(\lambda, X) = \left(\frac{\lambda}{\lambda+1} \right)^2 \underline{w}^T \underline{w} + \sigma^2 \left[1 + \frac{(p+1)}{N(1+\lambda)^2} \right]$$

$$\frac{dMSE}{d\lambda} = \left(\frac{2\lambda}{(\lambda+1)^2} - \frac{2\lambda^2}{(\lambda+1)^3} \right) \underline{w}^T \underline{w} - \frac{2\sigma^2(p+1)}{N(1+\lambda)^3} = 0$$

$$\Rightarrow \frac{2\lambda}{(1+\lambda)^3} \underline{w}^T \underline{w} - \frac{2\sigma^2(p+1)}{N(1+\lambda)^3} = 0 \Rightarrow \lambda^* = \frac{\sigma^2(p+1)}{N \cdot \underline{w}^T \underline{w}}$$

$$\text{Note: } \frac{dMSE}{d\lambda} = 2 \underline{w}^T \underline{w} \left(\frac{\lambda - \lambda^*}{(1+\lambda)^3} \right)$$

$$\Rightarrow \frac{d^2MSE}{d\lambda^2} = 2 \underline{w}^T \underline{w} \left(\frac{1}{(1+\lambda)^3} - \frac{3(\lambda - \lambda^*)}{(1+\lambda)^4} \right) = \frac{2 \underline{w}^T \underline{w}}{(1+\lambda^*)^3} > 0 \text{ at } \lambda = \lambda^*$$



Problem 6:

$$C_1 : \left\{ \underline{x}^1 = \begin{bmatrix} 1 \\ 5 \end{bmatrix}; \underline{x}^2 = \begin{bmatrix} -2 \\ -4 \end{bmatrix} \right\}$$

$$C_2 : \left\{ \underline{x}^3 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}; \underline{x}^4 = \begin{bmatrix} -1 \\ 5 \end{bmatrix} \right\}$$

$$\underline{\Phi}(\underline{x}) = [1 \sqrt{2}x_1 \sqrt{2}x_2 \sqrt{2}x_1x_2 \ x_1^2 \ x_2^2]^T$$

$$\underline{\Phi}_1 = [1 \sqrt{2} \sqrt{2}x_5 \sqrt{2}x_5 \ 1 \ 25]^T; \underline{\Phi}_2 = [1 - \sqrt{2}x_2 - \sqrt{2}x_4 \sqrt{2}x_8 \ 4 \ 16]^T; z^1 = z^2 = -1$$

$$\underline{\Phi}_3 = [1 \sqrt{2}x_2 \sqrt{2}x_3 \sqrt{2}x_6 \ 4 \ 9]^T; \underline{\Phi}_2 = [1 - \sqrt{2} \sqrt{2}x_5 - \sqrt{2}x_5 \ 1 \ 25]^T; z^3 = z^4 = 1$$

$$Dual : \max_{\underline{\lambda}} q(\underline{\lambda}) = \sum_{i=1}^4 \lambda_i - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \lambda_i \lambda_j z^i z^j \underline{\Phi}(\underline{x}^i)^T \underline{\Phi}(\underline{x}^j)$$

$$subject\ to : \sum_{i=1}^N \lambda_i z^i = 0 \text{ and } 0 \leq \lambda_i \Rightarrow \lambda_1 + \lambda_2 = \lambda_3 + \lambda_4$$

$$\underline{\lambda} = [0.0154 \ 0.0067 \ 0.0126 \ 0.0095]^T$$

$$\underline{w}^T = [0.0000 \ 0.0192 \ 0.0494 \ -0.1453 \ 0.0178 \ -0.1415]$$

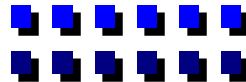
$$w_0 = -3.1713$$

$$g(\underline{x}) = \underline{w}^T \underline{\Phi}(\underline{x}) - w_0 = 0.0272x_1 + 0.0699x_2 - 0.2037x_1x_2 + 0.0178x_1^2 - 0.1415x_2^2 + 3.1713 = 0$$

Support vectors: All four are support vectors!

$$\frac{1}{\rho} = \|\underline{w}\|_2 = \frac{1}{\text{margin}} = \sqrt{\sum_{i=1}^N \lambda_i} = 0.2104$$

$$\text{margin} = \rho = 4.7529$$





■ Problem 7:

$$N = 11; p = 2$$

$$\hat{w}_{ML} = (X^T X)^{-1} X^T \underline{z} = \begin{bmatrix} -3.2564 \\ 0.0427 \end{bmatrix}$$

$$\hat{\sigma}^2 = \frac{1}{N - p} \|\underline{z} - X \hat{w}\|^2 = 0.0170$$

$$\Sigma_{ML} = \hat{\sigma}^2 (X^T X)^{-1} = \begin{bmatrix} 0.1323 & -0.0012 \\ -0.0012 & 0.00001142 \end{bmatrix}$$

$$\bar{w}_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}; \Sigma_0 = \begin{bmatrix} 10^{10} & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma_{MAP} = \left[\Sigma_{ML}^{-1} + \Sigma_0^{-1} \right]^{-1} = \begin{bmatrix} 0.132323933270331 & -0.001222250333231 \\ -0.001222250333231 & 0.000011422900311 \end{bmatrix}$$

$$\hat{w}_{MAP} = \Sigma_{MAP} \left[\Sigma_{ML}^{-1} \hat{w}_{ML} + \Sigma_0^{-1} \bar{w}_0 \right] = \begin{bmatrix} -3.256376353285603 \\ 0.042650925986531 \end{bmatrix}$$

$$w_1 \sim N(w_1; 0.042650925986531, 0.000011422900311)$$

When the prior is non-informative, ML and MAP give virtually the same estimates



Problem 8

$$J = -\ln p(\theta_1, \theta_2 | D) = N\theta_2 + \frac{e^{-2\theta_2}}{2} [Ns^2 + N(\bar{z} - \theta_1)^2]$$

where \bar{z} is the sample mean and s^2 is the sample variance

$$\nabla_{\underline{\theta}} J = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix} = \begin{bmatrix} -Ne^{-2\theta_2}(\bar{z} - \theta_1) \\ N - [Ns^2 + N(\bar{z} - \theta_1)^2]e^{-2\theta_2} \end{bmatrix}$$

$$\nabla_{\underline{\theta}} J |_{\hat{\underline{\theta}}_{MAP}} = \underline{0} \Rightarrow \hat{\theta}_{1,MAP} = \bar{z}; \hat{\theta}_{2,MAP} = \ln s$$

$$\nabla_{\underline{\theta}}^2 J = \begin{bmatrix} \frac{\partial^2 J}{\partial \theta_1^2} & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2 J}{\partial \theta_1 \partial \theta_2} & \frac{\partial^2 J}{\partial \theta_2^2} \end{bmatrix} = \begin{bmatrix} Ne^{-2\theta_2} & 2Ne^{-2\theta_2}(\bar{z} - \theta_1) \\ 2Ne^{-2\theta_2}(\bar{z} - \theta_1) & 2Ns^2e^{-2\theta_2} \end{bmatrix}$$

$$\nabla_{\underline{\theta}}^2 J |_{\hat{\underline{\theta}}_{MAP}} = H = \begin{bmatrix} \frac{N}{s^2} & 0 \\ 0 & 2N \end{bmatrix} \Rightarrow \Sigma = \begin{bmatrix} \frac{s^2}{N} & 0 \\ 0 & \frac{1}{2N} \end{bmatrix}$$

$$p(\theta_1, \theta_2 | D) = N(\theta_1; \bar{z}, \frac{s^2}{N})N(\theta_2; \ln s, \frac{1}{2N}) = \frac{1}{\sqrt{2\pi}} \frac{N}{s} e^{-N \left[\frac{(\theta_1 - \bar{z})^2}{2s^2} + (\theta_2 - \ln s)^2 \right]}$$



Problem 9:

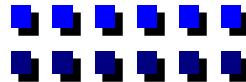
$$P(\underline{v}, \underline{h}) = \frac{1}{Z} \exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{i=1}^m b_i h_i + \sum_{j=1}^n c_j v_j\right)$$

$$\text{where } Z = \sum_{\underline{v}} \sum_{\underline{h}} \exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{i=1}^m b_i h_i + \sum_{j=1}^n c_j v_j\right)$$

$$\begin{aligned} P(\underline{h} | \underline{v}) &= \frac{P(\underline{h}, \underline{v})}{P(\underline{v})} = \frac{P(\underline{h}, \underline{v})}{\sum_{\underline{h}'} P(\underline{h}', \underline{v})} = \frac{\exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{i=1}^m b_i h_i + \sum_{j=1}^n c_j v_j\right)}{\sum_{\underline{h}'} \exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{i=1}^m b_i h_i + \sum_{j=1}^n c_j v_j\right)} \\ &= \frac{\exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{i=1}^m b_i h_i\right)}{\sum_{\underline{h}'} \exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{i=1}^m b_i h_i\right)} = \prod_{i=1}^m \frac{\exp\left(\sum_{j=1}^n d_{ij} v_j + b_i\right) h_i}{1 + \exp\left(\sum_{j=1}^n d_{ij} v_j + b_i\right)}; h_i \in \{0, 1\} \end{aligned}$$

and consequently

$$P(h_i = 1 | \underline{v}) = \frac{\exp\left(\sum_{j=1}^n d_{ij} v_j + b_i\right)}{\left[1 + \exp\left(\sum_{j=1}^n d_{ij} v_j + b_i\right)\right]} = g\left(\sum_{j=1}^n d_{ij} v_j + b_i\right) \dots \text{sigmoid function}$$





Problem 9:

$$\text{So, } P(\underline{v} | \underline{h}) = \frac{\exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{j=1}^n c_j v_j\right)}{\sum_{\underline{v}} \exp\left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} h_i v_j + \sum_{j=1}^n c_j v_j\right)} = \prod_{j=1}^n \frac{\exp\left(\sum_{i=1}^m d_{ij} h_i + c_j\right) v_j}{1 + \exp\left(\sum_{i=1}^m d_{ij} h_i + c_j\right)}; v_j \in \{0, 1\}$$

and consequently

$$P(v_j = 1 | \underline{h}) = \frac{\exp\left(\sum_{i=1}^m d_{ij} h_i + c_j\right)}{\left[1 + \exp\left(\sum_{i=1}^m d_{ij} h_i + c_j\right)\right]} = g\left(\sum_{i=1}^m d_{ij} h_i + c_j\right)$$



■ Problem 10:

Complete Data Likelihood

$$D_c = \{(x^1, \underline{z}^1), (x^2, \underline{z}^2), (x^3, \underline{z}^3)\}; \underline{z} \in \{[1 \ 0], [0 \ 1]\}; x^1 = 1; x^2 = 10; x^3 = 20$$

$$\Rightarrow -\ln p(D_c | \underline{\theta}) = \sum_{n=1}^3 \sum_{j=1}^2 z_j^n \left\{ -\ln P_j + \frac{1}{2} \ln 2\pi + \ln \sigma_j + \frac{1}{2\sigma_j^2} (x^n - \mu_j)^2 \right\}; \underline{\theta} = \{P_1, P_2, \sigma_1, \sigma_2, \mu_1, \mu_2\}$$

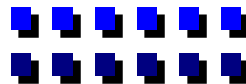
EM minimizes the expected value of the negative complete data log likelihood:

$$E_{\underline{z}} \{-\ln p(D_c | \underline{\theta})\} = \sum_{n=1}^3 \sum_{j=1}^2 \gamma_j^n \left\{ -\ln P_j + \frac{1}{2} \ln 2\pi + \ln \sigma_j + \frac{1}{2\sigma_j^2} (x^n - \mu_j)^2 \right\}$$

$$\Gamma = \begin{bmatrix} 1 & 0 \\ 0.4 & 0.6 \\ 0 & 1 \end{bmatrix}; \text{Recall } N_j = \sum_{n=1}^N \gamma_j^n; j = 1, 2, \dots, M \Rightarrow N_1 = 1.4; N_2 = 1.6$$

M-step: Re-estimate the parameters using the current responsibilities

$$\mu_j^{new} = \frac{1}{N_j} \sum_{n=1}^N \gamma_j^n x^n \Rightarrow \mu_1^{new} = \frac{1 + 0.4 * 10}{1.4} = 3.5714; \mu_2^{new} = \frac{0.6 * 10 + 20}{1.6} = 16.25$$





■ Problem 10:

M-step: Re-estimate the parameters using the current responsibilities

$$\mu_j^{new} = \frac{1}{N_j} \sum_{n=1}^N \gamma_j^n x^n \Rightarrow \mu_1^{new} = \frac{1 + 0.4 * 10}{1.4} = 3.5714; \mu_2^{new} = \frac{0.6 * 10 + 20}{1.6} = 16.25$$

$$\sigma_j^{2,new} = \frac{1}{N_j} \sum_{n=1}^N \gamma_j^n (x^n - \mu_j^{new})^2 \Rightarrow \sigma_1^{2,new} = 16.5306; \sigma_2^{2,new} = 23.4375$$

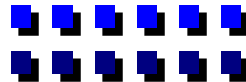
$$P_j^{new} = \frac{N_j}{N} \Rightarrow P_1^{new} = 0.4667; P_2^{new} = 0.5333$$

E-step:

$$\gamma_j^n = \frac{P_j N(\underline{x}^n; \underline{\mu}_j, \Sigma_j)}{\sum_{k=1}^M P_k N(\underline{x}^n; \underline{\mu}_k, \Sigma_k)}; j = 1, 2, \dots, M; n = 1, 2, \dots, N$$

$$\Gamma = \begin{bmatrix} 0.9919 & 0.0081 \\ 0.4072 & 0.5928 \\ 0.0004 & 0.9996 \end{bmatrix} \text{ etc. After 5 iterations: } \Gamma = \begin{bmatrix} 0.5373 & 0.4627 \\ 0.5373 & 0.4627 \\ 0.5373 & 0.4627 \end{bmatrix}; \text{NLL}=11.7250$$

All converge to same mean of 10.333 and variance of 60.222.





■ Problem 10:

With $\Gamma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}$, it converges to means 1 and 15 and variances of 0 and 25. NLL = -9.6366 setting var = max(0, eps)

With $\Gamma = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$, it converges to means 5.5 and 20 and variances of 20.25 and 0. NLL = -9.8473 setting var = max(0, eps)

Different initializations will give different answers!!!