

1.1. Variance and covariance - 10 pts. Let  $X, Y$  be two independent random vectors in  $\mathbb{R}^m$ .

(a) Find their covariance.

(b) For a constant matrix  $A \in \mathbb{R}^{m \times m}$ , show the following two properties:

$$\mathbb{E}(X + AY) = \mathbb{E}(X) + A\mathbb{E}(Y)$$

$$\text{Var}(X + AY) = \text{Var}(X) + A\text{Var}(Y)A^T$$

(c) Using part (b), show that if  $X \sim \mathcal{N}(\mu, \Sigma)$ , then  $AX \sim \mathcal{N}(A\mu, A\Sigma A^T)$ .

a)  $\text{Cov}(X, Y) = E(XY) - \mu_X \mu_Y$

But, since  $X, Y$  are independent R.V then  $\text{Cov}(X, Y) = 0$

b) Let  $A \in \mathbb{R}^{m \times m}$  be a constant matrix

$$E(X + AY) = E(X) + E(AY)$$

$$= E(X) + AE(Y) \quad \text{since } A \text{ is constant matrix}$$

$$\text{Var}(X + AY) = \text{Var}(X) + \text{Var}(AY) + \text{Cov}(X, Y)$$

$$= \text{Var}(X) + \text{Var}(AY)$$

$$= \text{Var}(X) + A\text{Var}(Y)A^T$$

Note:  $\text{Var}(AY) = E[(A(Y - \mu_Y))(A(Y - \mu_Y))^T]$

$$= E(A(Y - \mu_Y)(Y - \mu_Y)^T A^T)$$

$$= A\text{Var}(Y)A^T$$

c) Assume  $X \sim \mathcal{N}(\mu, \Sigma)$

$$\begin{aligned} \text{Then } E(AX) &= AE(X) \quad \text{and} \quad \text{Var}(AX) = A\text{Var}(X)A^T \\ &= A\mu \quad \quad \quad = A\Sigma A^T \end{aligned}$$

$$\text{Then } AX \sim \mathcal{N}(A\mu, A\Sigma A^T)$$

■

1.2. Calculus - 10 pts. Let  $x, y \in \mathbb{R}^m$  and  $A \in \mathbb{R}^{m \times m}$ . In vector notation, what is

- (a) the gradient with respect to  $x$  of  $x^T y$ ?
- (b) the gradient with respect to  $x$  of  $x^T x$ ?
- (c) the gradient with respect to  $x$  of  $\frac{1}{2} x^T A x$ ?
- (d) the gradient with respect to  $x$  of  $\exp(x^T A x)$ ?

$$\text{Let } x^T = [x_1, \dots, x_m] \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

$$a) \quad x^T y = x_1 y_1 + x_2 y_2 + \dots + x_m y_m$$

$$\frac{\partial}{\partial x} x^T y = \begin{bmatrix} \frac{\partial}{\partial x_1} x^T y \\ \vdots \\ \frac{\partial}{\partial x_m} x^T y \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} (x_1 y_1 + \dots + x_m y_m) \\ \vdots \\ \frac{\partial}{\partial x_m} (x_1 y_1 + \dots + x_m y_m) \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = y$$

$$b) \quad x^T x = x_1^2 + x_2^2 + \dots + x_m^2$$

$$\frac{\partial}{\partial x} x^T x = \begin{bmatrix} \frac{\partial}{\partial x_1} x^T x \\ \vdots \\ \frac{\partial}{\partial x_m} x^T x \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} (x_1^2 + \dots + x_m^2) \\ \vdots \\ \frac{\partial}{\partial x_m} (x_1^2 + \dots + x_m^2) \end{bmatrix} = \begin{bmatrix} 2x_1 \\ \vdots \\ 2x_m \end{bmatrix} = 2x$$

$$\begin{aligned} c) \quad \frac{\partial}{\partial x} \frac{1}{2} x^T A x &= \frac{1}{2} \frac{\partial}{\partial x} \sum_{k=1}^m \sum_{j=1}^m a_{kj} x_j x_k \\ &= \frac{1}{2} \sum_{j=1}^m a_{ij} x_j + \sum_{k=1}^m a_{ki} x_k \\ &= \frac{1}{2} (Ax + A^T x) \\ &= \frac{1}{2} (A + A^T) x \end{aligned}$$

$$\text{if } A \text{ is symmetric } \Rightarrow A = A^T$$

$$\text{so, } \frac{\partial}{\partial x} \frac{1}{2} x^T A x = Ax$$

$$\begin{aligned} d) \quad \frac{\partial}{\partial x} \exp \{ x^T A x \} &= \frac{\partial}{\partial x} x^T A x \times \exp \{ x^T A x \} \\ &= 2Ax \exp \{ x^T A x \} \end{aligned}$$

2.1. Linear regression - 20 pts. Suppose that  $\Phi \in \mathbb{R}^{n \times m}$  with  $n \geq m$  and  $\mathbf{t} \in \mathbb{R}^n$ , and that  $\mathbf{t} | (\Phi, \mathbf{w}) \sim \mathcal{N}(\Phi \mathbf{w}, \sigma^2 \mathbf{I})$ . We know that the maximum likelihood estimate  $\hat{\mathbf{w}}$  of  $\mathbf{w}$  is given by

$$\hat{\mathbf{w}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}.$$

- (a) Write the log-likelihood implied by the model above, and compute its gradient w.r.t.  $\mathbf{w}$ . By setting it equal to 0, derive the above estimator  $\hat{\mathbf{w}}$ .  
 (b) Find the distribution of  $\hat{\mathbf{w}}$ , its expectation and covariance matrix.

a) given  $\Phi, \mathbf{w}$  we can write the log likelihood:

$$\begin{aligned} p(\mathbf{t} | \Phi \mathbf{w}, \sigma^2 \mathbf{I}) &= \prod_{i=1}^N \mathcal{N}(t_i | \Phi \mathbf{w}, \sigma^2 \mathbf{I}) \\ &= \prod_{i=1}^N \frac{1}{(2\pi)^{D/2}} \frac{1}{|\sigma^2 \mathbf{I}|^{1/2}} \exp \left\{ -\frac{1}{2} (t_i - \Phi \mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (t_i - \Phi \mathbf{w}) \right\} \end{aligned}$$

$$\begin{aligned} \ln(p(\mathbf{t} | \Phi \mathbf{w}, \sigma^2 \mathbf{I})) &= \sum_{i=1}^N \ln \left( \frac{1}{(2\pi)^{D/2}} \frac{1}{|\sigma^2 \mathbf{I}|^{1/2}} \right) + \left( -\frac{1}{2} (t_i - \Phi \mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (t_i - \Phi \mathbf{w}) \right) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\sigma^2 \mathbf{I}| - \frac{1}{2} \sum_{i=1}^N (t_i - \Phi \mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (t_i - \Phi \mathbf{w}) \\ &\propto \frac{1}{2\sigma^2 \mathbf{I}} \sum_{i=1}^N (t_i - \Phi \mathbf{w})^T (t_i - \Phi \mathbf{w}) \quad \text{constant terms removed} \end{aligned}$$

Computing the gradient wrt  $\mathbf{w}$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{t} | \Phi \mathbf{w}, \sigma^2 \mathbf{I}) &= \frac{1}{2\sigma^2 \mathbf{I}} \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}} (t_i - \Phi \mathbf{w})^T (t_i - \Phi \mathbf{w}) \\ &= \frac{1}{2\sigma^2 \mathbf{I}} \times \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}} (t_i^T t_i - t_i^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T t_i + \mathbf{w}^T \Phi^T \mathbf{w}) \\ &= \frac{1}{2\sigma^2 \mathbf{I}} \times 2(\Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}) \\ &= \frac{\Phi^T (\mathbf{t} - \Phi \mathbf{w})}{\sigma^2} \end{aligned}$$

Setting it to 0 then finding  $\hat{\mathbf{w}}$ :

$$\frac{\Phi^T (\mathbf{t} - \Phi \mathbf{w})}{\sigma^2} = 0$$

$$\Rightarrow \Phi^T \mathbf{t} - \Phi^T \Phi \hat{\mathbf{w}} = 0$$

$$\Phi^T \Phi \hat{\mathbf{w}} = \Phi^T \mathbf{t}$$

$$\hat{\mathbf{w}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

■

$$b) E(\hat{w}) = E[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}]$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\mathbf{t})$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} \\ = \mathbf{w}$$

$$\text{Var}(\hat{w}) = \text{Var}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}]$$

$$= [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}] \text{Var}(\mathbf{t}) [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}]^T$$

$$= [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}] \sigma^2 \mathbf{I} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}]^T$$

$$= \sigma^2 [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}] \mathbf{I} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}]^T$$

$$= \sigma^2$$

Then we have  $\hat{w} \sim N(\mathbf{w}, \sigma^2)$

because linear transformation of a Gaussian R.V is Gaussian again. ■

2.2. Ridge regression and MAP - 20 pts. Suppose that we have  $\mathbf{t} | (\Phi, \mathbf{w}) \sim \mathcal{N}(\Phi \mathbf{w}, \sigma^2 \mathbf{I})$  and we place a normal prior on  $\mathbf{w} | \Phi$ , i.e.,  $\mathbf{w} \sim \mathcal{N}(0, \tau^2 \mathbf{I})$ . Recall from the first lecture (also in preliminaries.pdf) that MAP estimate of  $\mathbf{w}$  is given as the maximum of the posterior density

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \{ p(\mathbf{w} | \Phi, \mathbf{t}) \propto p(\mathbf{t} | \Phi, \mathbf{w}) p(\mathbf{w} | \Phi) \}.$$

Here,  $\propto$  notation means *proportional to*, and is used since we dropped the term  $p(\mathbf{t} | \Phi)$  in the denominator as it doesn't have  $\mathbf{w}$  in it, thus it doesn't contribute to the maximization problem.

Show that the MAP estimate of  $\mathbf{w}$  given  $(\mathbf{t}, \Phi)$  in this context is

$$(2.1) \quad \hat{\mathbf{w}}_{MAP} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t}$$

where  $\lambda = \sigma^2 / \tau^2$ .

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \{ p(\mathbf{w} | \Phi, \mathbf{t}) \propto p(\mathbf{t} | \Phi, \mathbf{w}) p(\mathbf{w} | \Phi) \}$$

We know the distribution of  $p(\mathbf{t} | \Phi, \mathbf{w})$  and  $p(\mathbf{w} | \Phi)$  then,

$$\begin{aligned} \log(p(\mathbf{t} | \Phi, \mathbf{w}) p(\mathbf{w} | \Phi)) &= \log \left[ \frac{1}{(2\pi)^{n/2} \cdot |\tau^2 \mathbf{I}|^{1/2}} e^{-\frac{\mathbf{w}^T \mathbf{w}}{2\tau^2}} \right] \times \frac{1}{(2\pi)^{n/2} \cdot |\sigma^2 \mathbf{I}|^{1/2}} e^{-\frac{(\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})}{2\sigma^2}} \\ &= -\frac{\mathbf{w}^T \mathbf{w}}{2\tau^2} - \frac{(\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})}{2\sigma^2} + C \end{aligned}$$

\* all constant forms bundled to C

Then taking gradient with respect to  $\mathbf{w}$ :

$$\begin{aligned} \frac{d \log(p(\mathbf{t} | \Phi, \mathbf{w}) p(\mathbf{w} | \Phi))}{d\mathbf{w}} &= \frac{d}{d\mathbf{w}} \left( -\frac{\mathbf{w}^T \mathbf{w}}{2\tau^2} - \frac{\mathbf{t}^T \mathbf{t} - \mathbf{t}^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}}{2\sigma^2} \right) \\ &= -\frac{2\mathbf{w}}{2\tau^2} + \frac{2\Phi^T (\mathbf{t} - \Phi \mathbf{w})}{2\sigma^2} \end{aligned}$$

Setting it to equal to 0:

$$\frac{\mathbf{w}}{\tau^2} - \frac{\Phi^T (\mathbf{t} - \Phi \mathbf{w})}{\sigma^2} = 0$$

$$\frac{\sigma^2}{\tau^2} \mathbf{w} - \Phi^T \mathbf{t} + \Phi^T \Phi \mathbf{w} = 0$$

$$\left( \frac{\sigma^2}{\tau^2} + \Phi^T \Phi \right) \mathbf{w} = \Phi^T \mathbf{t}$$

$$\hat{\mathbf{w}}_{MAP} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t} \quad \text{where } \lambda = \frac{\sigma^2}{\tau^2}$$

■