

CS 189 HW01

Hanze Yao
SID: 3033083286

1. Properties of Gaussians

(a) $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}$

$$\begin{aligned} E(e^{\lambda x}) &= \int_{-\infty}^{\infty} e^{\lambda x} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}} dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x+\lambda\sigma^2)^2 - \lambda^2\sigma^4}{2\sigma^2}} dx \\ &= e^{\frac{\lambda^2\sigma^2}{2}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x+\lambda\sigma^2)^2}{2\sigma^2}} dx \\ &= e^{\frac{\lambda^2\sigma^2}{2}} = e^{\frac{\sigma^2\lambda^2}{2}} \end{aligned}$$

(b) Use Markov inequality

$$P(X \geq a) \leq \frac{E(X)}{a} \quad (\text{when } X \text{ is nonnegative and } a > 0)$$

$$P(X \geq t) \leq \frac{E(X)}{t}$$

$$P(X \geq t) = P\left(e^{\frac{t}{\sigma^2} X} \geq e^{\frac{t}{\sigma^2} X t}\right) \leq \frac{E(e^{\frac{t}{\sigma^2} X})}{e^{\frac{t}{\sigma^2} X t}} = \frac{e^{\frac{(\frac{t}{\sigma^2})^2 \times \sigma^2}{2}}}{e^{\frac{t^2}{\sigma^2}}} = \frac{e^{\frac{t^2}{2\sigma^2}}}{e^{\frac{t^2}{\sigma^2}}} = e^{-\frac{t^2}{2\sigma^2}}$$

This is because $\frac{t}{\sigma^2} > 0$, then $e^{\frac{t}{\sigma^2} X}$ is an increasing function.

$$\because X \sim N(0, \sigma^2) \therefore f(x) = f(-x)$$

$$\therefore P(X \geq t) = P(X \leq -t) =$$

$$\therefore P(X \leq -t \text{ or } X \geq t) \leq 2P(X \geq t) = 2e^{-\frac{t^2}{2\sigma^2}}$$

$$\Downarrow$$

$$P(|X| \geq t) \leq 2e^{-\frac{t^2}{2\sigma^2}}$$

(c) $\because X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2) \therefore \frac{1}{n} \sum_{i=1}^n X_i \sim N(0, \frac{n\sigma^2}{n^2})$

$$P\left(\frac{1}{n} \sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{t^2}{2 \times \frac{\sigma^2}{n}}\right) = \exp\left(-\frac{nt^2}{2\sigma^2}\right) \quad N\left(0, \frac{\sigma^2}{n}\right)$$

$$\text{as } n \rightarrow \infty \quad P\left(\frac{1}{n} \sum_{i=1}^n X_i \geq t\right) \rightarrow 0$$

(d) First, need to prove that U_X and V_X are uncorrelated.

$$\text{Need to prove that } E(U_X V_X) = E(U_X) E(V_X)$$

$$U_X = \langle U, X \rangle = U^T X \sim N(0, \|U\|_2^2)$$

$$V_X = \langle V, X \rangle = V^T X \sim N(0, \|V\|_2^2)$$

$$E(U_X V_X) = E((U_1 X_1 + U_2 X_2 + \dots + U_d X_d)(V_1 X_1 + V_2 X_2 + \dots + V_d X_d))$$

$$= E \left(U_1 V_1 X_1^2 + U_2 V_2 X_2^2 + \dots + U_d V_d X_d^2 + \sum_{i \neq j} U_i V_j X_i X_j \right)$$

$$\because X_1, \dots, X_d \text{ iid } N(0, 1)$$

$$\therefore E(X_i^2) = E(X_j^2) = \dots = E(X_d^2) = \text{Var}(X) - (E(X))^2 = 1 - 0^2 = 1$$

$$E \left(\sum_{i \neq j} U_i V_j X_i X_j \right) = \sum_{i \neq j} U_i V_j E(X_i) E(X_j) = 0$$

$$\therefore E(U_x V_x) = U_1 V_1 + U_2 V_2 + \dots + U_d V_d = \langle U, V \rangle$$

$$\because U \perp V \therefore E(U_x V_x) = 0 = \langle U, V \rangle$$

$$\therefore E(U_x) = E(V_x) = 0$$

$$\therefore E(U_x V_x) = E(U_x) E(V_x)$$

$$\therefore U_x \text{ and } V_x \text{ are uncorrelated}$$

$$\therefore \text{jointly normal random variables are independent iff they are uncorrelated}$$

$$\therefore U_x \text{ and } V_x \text{ are independent}$$

2. Identities with expectation

$$(a) f(x) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$E(x^k) = \int_0^{\infty} x^k \lambda e^{-\lambda x} dx = -x^k e^{-\lambda x} \Big|_0^{\infty} + \int_0^{\infty} k x^{k-1} e^{-\lambda x} dx$$

$$= 0 - \frac{1}{\lambda} e^{-\lambda x} k x^{k-1} \Big|_0^{\infty}$$

$$= \underbrace{-\frac{1}{\lambda} e^{-\lambda x} k x^{k-1} \Big|_0^{\infty}}_{=0} + \int_0^{\infty} \frac{1}{\lambda} e^{-\lambda x} k(k-1) x^{k-2} dx$$

$$= \int_0^{\infty} \frac{1}{\lambda} e^{-\lambda x} k(k-1) x^{k-2} dx$$

When this integration by part ^{process} is over,

$$E(x^k) = \int_0^{\infty} \frac{1}{\lambda^{k+1}} e^{-\lambda x} k! dx$$

$$= \frac{k!}{\lambda^{k+1}} \int_0^{\infty} e^{-\lambda x} dx$$

$$= \frac{k!}{\lambda^{k+1}} \times \left(-\frac{1}{\lambda} e^{-\lambda x} \Big|_0^{\infty} \right) = \frac{k!}{\lambda^k}$$

$$(b) X = \int_0^x dt = \int_0^{\infty} 1\{x \geq t\} dt$$

$$E(x) = E\left[\int_0^{\infty} 1\{x \geq t\} dt\right] = \int_0^{\infty} E[1\{x \geq t\}] dt$$

$$= \int_0^{\infty} P(x \geq t) dt$$

$$(c) P(X > 0) = E(1\{X > 0\})$$

Use Cauchy-Schwarz inequality

$$(E(\underbrace{X \times 1\{X > 0\}}_{\downarrow}))^2 \leq E(x^2) E(\underbrace{(1\{X > 0\})^2}_{\downarrow})$$

$$\therefore X \geq 0$$

$$\therefore X 1\{X > 0\} = X, X > 0$$

$$X 1\{X > 0\} = 0, X = 0$$

$$\therefore X 1\{X > 0\} = X$$

$$(1\{X > 0\})^2$$

$\therefore 1\{X > 0\}$ can only be

0 or 1

$\therefore E((1\{X > 0\})^2)$ is the same as

$$E(1\{X > 0\})$$

$$\therefore (E(x))^2 \leq E(x^2) E(1\{X > 0\}) \Rightarrow (E(x))^2 \leq E(x^2) P(X > 0) \Rightarrow P(X > 0) \geq \frac{(E(x))^2}{E(x^2)}$$

3. Gradients and Norms

$$(a) \|x\|_2 = \sqrt{\sum_{j=1}^d |x_j|^2} \quad \|x\|_2^2 = \sum_{j=1}^d |x_j|^2$$

$$\|x\|_1 = \sum_{j=1}^d |x_j| \quad \|x\|_1^2 = \sum_{j=1}^d |x_j|^2 + 2 \sum_{i < j} |x_i| |x_j| \geq \|x\|_2^2$$

$\therefore \|x\|_2 \geq 0$ and $\|x\|_1 \geq 0$ Can prove $\|x\|_2 \leq \|x\|_1$ by proving $\|x\|_2^2 \leq \|x\|_1^2$

Now use Cauchy-Schwarz inequality

$$\|x\|_1 = \langle x, \underbrace{1}_{\substack{\text{is a } \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \text{ vector} \\ d \times 1}} \rangle \leq \sqrt{\langle x, x \rangle} \sqrt{\langle 1, 1 \rangle} = \|x\|_2 \times \sqrt{d} = \sqrt{d} \|x\|_2$$

$$(b) (i) \alpha = \sum_{j=1}^d y_j \ln p_j$$

$$\frac{\partial \alpha}{\partial p_j} = \left(\frac{\partial \alpha}{\partial p_1} \quad \frac{\partial \alpha}{\partial p_2} \quad \dots \quad \frac{\partial \alpha}{\partial p_d} \right) = \left(\frac{y_1}{p_1} \quad \frac{y_2}{p_2} \quad \dots \quad \frac{y_d}{p_d} \right)$$

$$(ii) \begin{pmatrix} \frac{\partial p_1}{\partial y_1} & \frac{\partial p_1}{\partial y_2} & \dots & \frac{\partial p_1}{\partial y_d} \\ \frac{\partial p_2}{\partial y_1} & \frac{\partial p_2}{\partial y_2} & \dots & \frac{\partial p_2}{\partial y_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial p_d}{\partial y_1} & \frac{\partial p_d}{\partial y_2} & \dots & \frac{\partial p_d}{\partial y_d} \end{pmatrix} = \begin{pmatrix} \cosh(y_1) & & & 0 \\ & \cosh(y_2) & & \\ & & \ddots & \\ 0 & & & \cosh(y_d) \end{pmatrix}$$

This is a diagonal matrix

$$(iii) \text{ Let } A = \begin{pmatrix} A_{11} & A_{12} & \dots & A_{1m} \\ A_{21} & A_{22} & \dots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{d1} & A_{d2} & \dots & A_{dm} \end{pmatrix} \quad p = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_m \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_d \end{pmatrix}$$

$$y = \begin{pmatrix} \sum_{j=1}^m A_{1j} p_j + b_1 \\ \sum_{j=1}^m A_{2j} p_j + b_2 \\ \vdots \\ \sum_{j=1}^m A_{dj} p_j + b_d \end{pmatrix} \quad \therefore \begin{pmatrix} \frac{\partial y_1}{\partial p_1} & \frac{\partial y_1}{\partial p_2} & \dots & \frac{\partial y_1}{\partial p_m} \\ \frac{\partial y_2}{\partial p_1} & \frac{\partial y_2}{\partial p_2} & \dots & \frac{\partial y_2}{\partial p_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_d}{\partial p_1} & \frac{\partial y_d}{\partial p_2} & \dots & \frac{\partial y_d}{\partial p_m} \end{pmatrix} = \begin{pmatrix} A_{11} & A_{12} & \dots & A_{1m} \\ A_{21} & A_{22} & \dots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{d1} & A_{d2} & \dots & A_{dm} \end{pmatrix}$$

$$(c) \quad X, A \in \mathbb{R}^{d \times d}$$

$$\nabla_x \text{Tr}(A^T X)$$

$$f(x) = \text{Tr}(A^T X)$$

$$f(x+\Delta) = \text{Tr}(A^T (x+\Delta))$$

$$= \text{Tr}(A^T x + A^T \Delta) = \text{Tr}(A^T x) + \text{Tr}(A^T \Delta)$$

$$= f(x) + \text{Tr}(A^T \Delta)$$

$$\therefore f(x+\Delta) = f(x) + \text{Tr}\left(\frac{\partial f}{\partial x} \Delta\right)$$

$$\therefore \frac{\partial f}{\partial x} = A^T \quad \therefore \nabla_x \text{Tr}(A^T X) = \left(\frac{\partial f}{\partial x}\right)^T = (A^T)^T = A$$

$$(d) \quad Xw = y$$

$$\|y - Xw\|_2^2 = (y - Xw)^T (y - Xw)$$

$$= (y^T - w^T X^T) (y - Xw)$$

$$= y^T y - y^T Xw - w^T X^T y + w^T X^T X w$$

$$= y^T y - 2y^T Xw + w^T X^T X w$$

$$\nabla_w \|y - Xw\|_2^2 = (-2y^T X)^T + (X^T X + X^T X) w$$

$$= 2X^T X w - 2X^T y$$

$$\therefore w^* = \arg\min_w \|y - Xw\|_2^2$$

$$\therefore \nabla_w \|y - Xw^*\|_2^2 = 0 \Rightarrow 2X^T X w^* - 2X^T y = 0$$

$$2X^T X w^* = 2X^T y$$

$$X^T X w^* = X^T y$$

$$\therefore X \text{ is full rank} \therefore X^T X \text{ is full rank as well}$$

$$\therefore w^* = (X^T X)^{-1} X^T y$$

4. Linear Algebra Review

(a) First: prove (i) \Rightarrow (ii): $X^T A X = X^T (U \Lambda U^T) X = (X^T U) \Lambda (U^T X)$

Let $Z = U^T X$, $\|X\|_2^2 = Z^T \Lambda Z = \sum_{i=1}^d \lambda_i z_i^2 \geq 0$

Then all λ_i must be non negative. Which is: All eigenvalues of A are non-negative.

Second: prove (ii) \Rightarrow (iii): If all A 's eigenvalues are non negative. Then $A = V \Lambda V^T = V \Lambda^{\frac{1}{2}} \Lambda^{\frac{1}{2}} V^T$

$\therefore \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$ and $\lambda_1, \lambda_2, \dots, \lambda_d \geq 0$

proves: $\Lambda^{\frac{1}{2}} = \text{diag}(\lambda_1^{\frac{1}{2}}, \lambda_2^{\frac{1}{2}}, \dots, \lambda_d^{\frac{1}{2}})$

$\therefore \Lambda^{\frac{1}{2}}$ is a diagonal matrix $\therefore (\Lambda^{\frac{1}{2}})^T = \Lambda^{\frac{1}{2}}$

$\therefore A = (V(\Lambda^{\frac{1}{2}})^T)(\Lambda^{\frac{1}{2}} V^T) = (V(\Lambda^{\frac{1}{2}})^T)(V(\Lambda^{\frac{1}{2}})^T)^T$

$\therefore U = V(\Lambda^{\frac{1}{2}})^T$

Third: prove (iii) \Rightarrow (i): According to eigen decomposition, $A = V \Lambda V^T$

$\therefore A = U U^T = V(\Lambda^{\frac{1}{2}})^T (V(\Lambda^{\frac{1}{2}})^T)^T = V(\Lambda^{\frac{1}{2}})^T \Lambda^{\frac{1}{2}} V$

$\therefore V(\Lambda^{\frac{1}{2}})^T \Lambda^{\frac{1}{2}} V = A = V \Lambda V^T \therefore (\Lambda^{\frac{1}{2}})^T = \Lambda^{\frac{1}{2}} \therefore \Lambda$ is a diagonal matrix with A 's eigenvalues on the diagonal $\therefore \Lambda^{\frac{1}{2}}$ must be

$\text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_d}) \therefore \lambda_1, \lambda_2, \dots, \lambda_d \geq 0 \therefore \sum_{i=1}^d \lambda_i x_i^2 \geq 0$

Which is $X^T A X \geq 0$ for all $X \in \mathbb{R}^d$

(b)

(i) $\therefore A$ and B are PSD

$\therefore A$ and B 's eigenvalues λ_i and λ_j are all non-negative

$\therefore 2A + 3B$'s eigenvalues $= 2\lambda_i + 3\lambda_j$

$\therefore 2A + 3B$'s eigenvalues are all non-negative as well

$\therefore 2A + 3B$ is PSD

(ii) According to a(iii), There exists a matrix $U \in \mathbb{R}^{d \times d}$, $A = U U^T$

Assume $U = \begin{pmatrix} -u_1 & & \\ & -u_2 & \\ & & -u_d \end{pmatrix} \therefore A = \begin{pmatrix} -u_1 & & \\ & -u_2 & \\ & & -u_d \end{pmatrix} \begin{pmatrix} | & | & | \\ u_1 & u_2 & \dots & u_d \\ | & | & | \end{pmatrix}$

Then diagonal entries of A are $u_1^T u_1, u_2^T u_2, \dots, u_d^T u_d$

Which are $\|u_1\|_2^2, \|u_2\|_2^2, \dots, \|u_d\|_2^2$.

Obviously they are non-negative

(iii) According to a(i), for all $X \in \mathbb{R}^d$, $X^T A X \geq 0$

if $X = \begin{pmatrix} 1 \\ i \\ 2 \\ \vdots \\ d-1 \end{pmatrix}$ Then $X^T A X = (1 \dots 1) A \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \geq 0$

$\sum_{i=1}^d \sum_{j=1}^d A_{ij} \geq 0$

Because $\text{Tr}(AB) = \text{Tr}(BA)$

$$(iv) \text{Tr}(AB) = \text{Tr}(AB^{\frac{1}{2}}B^{\frac{1}{2}}) = \text{Tr}(B^{\frac{1}{2}}AB^{\frac{1}{2}}) \\ = \text{Tr}(B^{\frac{1}{2}}A^{\frac{1}{2}}A^{\frac{1}{2}}B^{\frac{1}{2}})$$

$\therefore A, B$ are PSD matrix $\therefore A, B$ are symmetric

$\therefore A^{\frac{1}{2}}, B^{\frac{1}{2}}$ are symmetric

$$\therefore \text{Tr}(AB) = \text{Tr}(B^{\frac{1}{2}}A^{\frac{1}{2}}A^{\frac{1}{2}}B^{\frac{1}{2}}) = \text{Tr}(B^{\frac{1}{2}}A^{\frac{1}{2}}(B^{\frac{1}{2}}A^{\frac{1}{2}})^T)$$

According to Q(ii) $B^{\frac{1}{2}}A^{\frac{1}{2}}(B^{\frac{1}{2}}A^{\frac{1}{2}})^T$ is a PSD matrix

$$\therefore \text{Tr}(AB) = \text{Tr}(B^{\frac{1}{2}}A^{\frac{1}{2}}(B^{\frac{1}{2}}A^{\frac{1}{2}})^T) = \text{sum of eigenvalues of } B^{\frac{1}{2}}A^{\frac{1}{2}}(B^{\frac{1}{2}}A^{\frac{1}{2}})^T \\ \geq 0$$

(V) According to (iv) $B^{\frac{1}{2}}AB^{\frac{1}{2}}$ is a PSD matrix,

$\therefore B^{\frac{1}{2}}AB^{\frac{1}{2}}$ is a symmetric matrix, so it is diagonalizable

If $\text{Tr}(AB) = 0$, then $\text{Tr}(B^{\frac{1}{2}}AB^{\frac{1}{2}}) = 0$

\therefore All eigenvalues of $B^{\frac{1}{2}}AB^{\frac{1}{2}}$ are non-negative

\therefore All eigenvalues of $B^{\frac{1}{2}}AB^{\frac{1}{2}}$ are 0.

$$\therefore B^{\frac{1}{2}}AB^{\frac{1}{2}} = P \begin{pmatrix} 0 \end{pmatrix} P^{-1}$$

\hookrightarrow a zero matrix.

$$\therefore AB = B^{-\frac{1}{2}}(B^{\frac{1}{2}}AB^{\frac{1}{2}})B^{\frac{1}{2}} = B^{-\frac{1}{2}}P \begin{pmatrix} 0 \end{pmatrix} PB^{\frac{1}{2}} = 0$$

\hookrightarrow a zero matrix

When $AB = 0$, of course $\text{Tr}(AB) = 0$

\therefore If A and B are PSD, then $\text{Tr}(AB) = 0$ if and only if $AB = 0$

(C) According to the definition of eigenvalue

$$\lambda_A X = AX \quad \because \lambda_A \text{ is a scalar}$$

$$\therefore X\lambda_A = AX$$

$$X^T X \lambda_A = X^T A X$$

$$\text{if } \|X\|_2 = 1 \text{ then } X^T X = I$$

$$\therefore \lambda_A = X^T A X$$

$$\therefore \lambda_{\max}(A) = \max_{\|X\|_2 = 1} X^T A X$$

5. Covariance Practice

$$C = E[(Z - \mu)(Z - \mu)^T]$$

For $X \in \mathbb{R}^d$

$$E(\|Z - \mu\|^2) = E([Z - \mu]^T X) [Z - \mu]^T X) \geq 0$$

$$E(X^T (Z - \mu)(Z - \mu)^T X) \geq 0$$

$$X^T [E(Z - \mu)(Z - \mu)^T] X \geq 0$$

$$X^T C X \geq 0 \text{ for all } X \in \mathbb{R}^d$$

$\therefore C$ is PSD \Rightarrow The covariance matrix is always positive semi-definite.

6. A Simple Classification Approach.

(b) The residual error is: $\|Xw - y\|_2^2 = 422.75085$

The first 20 entries of w are as follows:

-0.33077663, 0.39177105, 0.14818856, -0.16060987
0.103271045, -0.019703118, -0.12769328, 0.009455555
-0.017149584, -0.005674782, -0.0046893246, -0.011267236
-0.0057097357, 0.004589835, 0.01787667, -0.03010068
0.010038383, -0.069282594, -0.023043357, -0.026301404

(c) Training set percentage of correctly classified: 0.9976

Test set percentage of correctly classified: 0.8981

(d) The performance evaluated on a separate test set is to avoid over-fitting problem.

Performance is similar in our case is because there are only two classes and we want to separate them using linear regression. The handwritten numbers 0 and 1 are very easy to classify.

(e) With 0/1 target

Training set percentage of correctly classified: 0.9897

Test set ————— : 0.9915

After adding bias term

With -1/1 target

Training = 0.9941

Test = 0.9962

With 0/1 target

Training = 0.9941

Test = 0.9962

After adding this bias term means the linear model has an interception.

For 0/1 target, the hyperplane actually needs an interception so after adding the bias term its accuracy increases.

For -1/1 target, the hyperplane does not need an interception, so after adding the bias term its accuracy decreases a little bit.

After adding the bias term, the result of the two kinds of target become the same.