

中国作家协会

STAT 414. Assignment 1.

#1.

- Solution: $\text{cov}(xy) = E[(x - E(x))(y - E(y))]$
 $= E[xy - E(x)y - xE(y) + E(x)E(y)]$
 $= E(xy) - E(x)E(y)$

(93)

(7) Since $x \perp y$, $P(xy) = P(x)P(y)$

so $E(xy) = \sum \sum p(xy) xy$
 $= \sum \sum p(x)p(y) xy$
 $= E(x)E(y)$

Therefore $\text{cov}(x,y) = 0$.

- Proof: $E(x) = \mu$, $\text{Var}(x) = \sigma^2$. Want to show $E(x_n x_m) = \mu^2 + I_{nm} \sigma^2$

Since the expectation $E(x_n^2) = \mu^2 + \sigma^2$

① Suppose $x_n = x_m$, then $E(x_n^2) = \mu^2 + \sigma^2$
 We know that the second moment of a Gaussian distribution is

$$E(x^2) = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2 \quad (*)$$

② Suppose $x_n \neq x_m$, $x_n \neq x_m$, so $E(x_n x_m) = E(x_n)E(x_m)$ b/c independence

Then $E(x_n x_m) = \mu \cdot \mu = \mu^2$ $\quad (**)$

Consider (*) & (**) together, we now have

$$E(x_n x_m) = \mu^2 + I_{nm} \sigma^2 \text{ with } I_{nm} = \begin{cases} 1 & \text{if } n=m \\ 0 & \text{otherwise.} \end{cases}$$

- Proof: joint distribution $p(x,y)$

To prove equation one, the key is to expand expectation to the form of integrals.

$$\begin{aligned} E_y E_x[x|y] &= E_y \left(\int x p(x|y) dx \right) = \int \left(\int x p(x|y) dx \right) p(y) dy \\ &= \iint x p(x,y) dx dy \quad \text{by product rule:} \\ &\quad p(x|y)p(y) \\ &= E(x) \end{aligned}$$

中 国 作 家 协 会

To prove the second equation,
we write:

$$\text{Var}(x) = E(x^2) - [E(x)]^2$$

and we apply the result of first equation:

~~$E_y E_x[x|y]$~~

$$E_y [\text{Var}_x[x|y]] = E_y [E_x[x^2|y] - E_x[x|y]^2] \quad \textcircled{1}$$

$$\text{Var}_y [E_x[x|y]] = E_y [E_x[x|y]^2] - E_y [E_x[x|y]]^2 \quad \textcircled{2}$$

$$\begin{aligned} \textcircled{1} + \textcircled{2} &= E[x^2] - E_y [E_x[x|y]^2] + E_y [E_x[x|y]]^2 - (E_x[x])^2 \\ &= \text{Var}(x) \end{aligned}$$

#2. $p(\text{mistake}) = \int_{R_1} p(x, C_2) dx + \int_{R_2} p(x, C_1) dx$

Proof: Want to show $p(\text{mistake}) \leq \int [p(x, C_1) p(x, C_2)]^{\frac{1}{2}} dx$

First we might want to use the product rule of probability to
rewrite this:

$$p(\text{mistake}) = \int_{R_1} p(C_2|x) p(x) dx + \int_{R_2} p(C_1|x) p(x) dx$$

Since we want to minimize the probability of mistake, i.e. the value of x
should make $p(C_2|x) \leq p(C_1|x)$ in R_1 and vice versa in R_2 .

Then $p(\text{mistake}) = \int_{R_1} p(C_2|x) p(x) dx + \int_{R_2} p(C_1|x) p(x) dx$

$$\leq \int_{R_1} [p(C_2|x) p(C_1|x)]^{\frac{1}{2}} p(x) dx + \int_{R_2} [p(C_1|x) p(C_2|x)]^{\frac{1}{2}} p(x) dx$$

$$= \int_{R_1+R_2} [p(C_2|x) p(x) p(C_1|x) p(x)]^{\frac{1}{2}} dx \quad \# \text{ same integrand}$$

$$= \int [p(x, C_2) p(x, C_1)]^{\frac{1}{2}} dx \quad \# \text{ use product rule again}$$

中 国 作 家 协 会

#3.

Proof: pdf: $f = \binom{N}{m} p^m (1-p)^{N-m}$

$$\begin{aligned}
 E(m) &= \sum_{m=0}^N m \binom{N}{m} p^m (1-p)^{N-m} \\
 &= \sum_{m=1}^N m \binom{N}{m} p^m (1-p)^{N-m} \\
 &= \sum_{m=1}^N N \binom{N-1}{m-1} p^m (1-p)^{N-m} \quad b/c \quad a \binom{b}{a} = b \binom{b-1}{a-1} \\
 &= Np \sum_{m=1}^N \binom{N-1}{m-1} p^{m-1} (1-p)^{(N-1)-(m-1)} \\
 &= Np (p + (1-p))^{N-1} \quad \text{by binomial theorem.} \\
 &= Np
 \end{aligned}$$

$$\text{Var}(m) = E\left(\frac{m^2}{m}\right) - [E\left(\frac{m}{m}\right)]^2$$

$$\begin{aligned}
 E(m^2) &= \sum_{m=0}^N m^2 \binom{N}{m} p^m (1-p)^{N-m} \\
 &= \sum_{m=0}^N mN \binom{N-1}{m-1} p^m (1-p)^{N-m} \\
 &= Np \sum_{m=0}^N m \binom{N-1}{m-1} p^{m-1} (1-p)^{N-m} \quad \text{let } N-1 = s \\
 &= Np \left(\sum_{t=0}^s t \binom{s}{t} p^t (1-p)^{s-t} \right) \\
 &= Np \left(\sum_{t=0}^s t \binom{s}{t} p^t (1-p)^{s-t} + \sum_{t=0}^s \binom{s}{t} p^t (1-p)^{s-t} \right) \\
 &= Np \left((s-1) \sum_{t=1}^s \binom{s-1}{t-1} p^{t-1} (1-p)^{(s-1)-(t-1)} + \sum_{t=0}^s \binom{s}{t} p^t (1-p)^{s-t} \right) \\
 &= Np \left((N-1) p (p+1-p)^{N-1} + (p+1-p)^N \right) \\
 &= Np (Np - p + 1) \\
 &= N^2 p^2 - Np^2 + Np
 \end{aligned}$$

$$\text{So } \text{Var}(m) = E(m^2) - [E(m)]^2 = N^2 p^2 - Np^2 + Np - N^2 p^2 = Np(1-p)$$

中 国 作 家 协 会

#4. Proof: Sps Σ is $n \times n$ matrix. (i.e. $D = n$)

Σ is symmetric, we can write it as a decomposition:

$$\Sigma = U \Lambda U^T \text{ where } U \text{ orthogonal, } \Lambda \text{ diagonal with } \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

$$\Sigma x = U \Lambda U^T x = U \Lambda \begin{pmatrix} u_1^T x \\ \vdots \\ u_n^T x \end{pmatrix} = U \begin{pmatrix} \lambda_1 u_1^T x \\ \vdots \\ \lambda_n u_n^T x \end{pmatrix}$$

essentially
we want to prove

$$\begin{aligned} &= \sum_{i=1}^n \lambda_i u_i^T x \cdot u_i \\ &= \left(\sum_{i=1}^n \lambda_i u_i u_i^T \right) x \end{aligned}$$

Hence $\Sigma = \sum_{i=1}^D \lambda_i u_i u_i^T$

Invert $\Sigma = U \Lambda U^T$ from both sides, note that U is orthogonal

i.e. $U^{-1} = U^T$

so $\Sigma^{-1} = U \Lambda^{-1} U^T$

$$= \sum_{i=1}^n \frac{1}{\lambda_i} u_i u_i^T$$

b/c Λ^{-1} is a diagonal matrix whose diagonal elements are inverses of their counterparts in diagonal matrix Λ .

中 国 作 家 协 会

#5.

Proof: Note that the coefficient before exponential part $\exp(-\frac{|x|^8}{20^2})$ must "normalize" the distribution.

$$\text{So we only calculate } \int_{-\infty}^{\infty} \exp\left(-\frac{|x|^8}{20^2}\right) dx$$

$$= 2 \int_0^{\infty} \exp\left(-\frac{|x|^8}{20^2}\right) dx$$

$$= 2 \int_0^{\infty} \frac{20^2 (20^2 u)^{\frac{1-8}{8}}}{8} \exp(-u) du$$

$$= \cancel{\frac{2(20^2)^{\frac{1}{8}}}{8}} \int_0^{\infty} u^{\frac{1-8}{8}} \exp(u) du = \frac{2(20^2)^{\frac{1}{8}} \Gamma(\frac{1}{8})}{8}$$

$$\text{let } u = \frac{x^8}{20^2} \\ \text{so } du = \frac{1}{20^2} 8x^{8-1}$$

$$\text{because } P(x) = \int_0^{\infty} u^{x-1} e^{-u} du$$

Hence this times the term $\frac{8}{2(20^2)^{\frac{1}{8}} \Gamma(\frac{1}{8})}$ is equal to 1.

Normalization proved. Also $p > 0$ because exponential is nonnegative with integral part = 1. So it's a valid distribution.

- When $q=2$, $P(x|\sigma^2, 2) = \frac{1}{(20^2)^{\frac{1}{2}} \Gamma(\frac{1}{2})} \exp\left(-\frac{(x-0)^2}{20^2}\right)$

$$= \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-0)^2}{20^2}\right)$$

$$\text{since } \Gamma(\frac{1}{2}) = \sqrt{\pi}$$

It's Gaussian.

- target variable $t = y(x, \omega) + \varepsilon$. ε is random noise

Log likelihood function over ω, σ^2 ?

observed data $X = \{x_1, \dots, x_N\}$

$$t = (t_1, \dots, t_N)^T$$

(this can be proved by applying a transformation to change variables to polar expressions).

(see next page)

中国作家协会

~~$y(x, \omega)$~~ $\xi = t - y(x, \omega)$
 Given noise variable, the conditional distribution of target variable t
 and input x, ω

$$\text{is } p(t|x, \omega, \sigma^2) = \frac{\frac{g}{\theta}}{2(2\sigma^2)^{\frac{1}{\theta}} \Gamma(\frac{1}{\theta})} \exp\left(-\frac{|t-y(x, \omega)|^{\theta}}{2\sigma^2}\right)$$

For the derivation of Log likelihood function, we take logarithm on both sides of the product of such conditional probability.
 The process is quite tedious with long terms of constants.

$$\begin{aligned} L &= \ln \prod_{i \in \{1, \dots, N\}} \left\{ \frac{\frac{g}{\theta}}{2(2\sigma^2)^{\frac{1}{\theta}} \Gamma(\frac{1}{\theta})} \exp\left(-\frac{|t_i - y(x_i, \omega)|^{\theta}}{2\sigma^2}\right) \right\} \\ &= \ln \left\{ \left(\frac{\frac{g}{\theta}}{2(2\sigma^2)^{\frac{1}{\theta}} \Gamma(\frac{1}{\theta})} \right)^N \exp\left(\sum_{i=1}^N \frac{|t_i - y(x_i, \omega)|^{\theta}}{2\sigma^2}\right) \right\} \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^N |t_i - y(x_i, \omega)|^{\theta} + N \ln \left\{ \frac{\frac{g}{\theta}}{2(2\sigma^2)^{\frac{1}{\theta}} \Gamma(\frac{1}{\theta})} \right\} \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^N |t_i - y(x_i, \omega)|^{\theta} - \underbrace{\frac{N}{\theta} \ln 2\sigma^2 + N \left(\ln \frac{g}{\theta} - \frac{1}{\theta} \ln 2 - \ln \Gamma\left(\frac{1}{\theta}\right) \right)}_{\text{with fixed } g.} \end{aligned}$$

this term is a constant real number, has nothing to do with any of x , ω or σ^2 .

中国作家协会

6. expected L_g loss function

$$E[L_g] = \iint |y(x) - t|^g p(x, t) dt dx$$

i). Proof: Let $y(x) \perp x$ so that $p(x, t) = p(t|x)$

$$\text{so } \frac{\partial}{\partial y(x)} E[L_g] = \frac{\partial}{\partial y(x)} \iint |y(x) - t|^g p(t|x) dt dx \\ = \int \cancel{g} |y(x) - t|^{g-1} p(t|x) dt$$

$$= 0 \quad (\text{since minimized, and to simplify. I didn't check the 2nd derivative})$$

$$\Rightarrow \int_{-\infty}^{\infty} |y(x) - t|^{g-1} p(t|x) dt = 0 \quad (*)$$

Since we want to check the median (and its relation with probability mass), so we split (*) into two parts: ~~at the same time~~

$$\cancel{\int_{-\infty}^{\infty} p(t|x) dt} = \int_{-\infty}^{y(x)} p(t|x) dt$$

$$\int_{y(x)}^{\infty} |y(x) - t|^{g-1} p(t|x) dt = \int_{-\infty}^{y(x)} |y(x) - t|^{g-1} p(t|x) dt$$

where left part for $t \geq y(x)$, right part for $t < y(x)$.

$$\text{Plug in } g=1, \int_{y(x)}^{\infty} p(t|x) dt = \int_{-\infty}^{y(x)} p(t|x) dt$$

so $y(x)$ is the conditional median for ~~sure~~ sure.



(8)

中 国 作 家 协 会

b). Proof: as $q \rightarrow 0$, $|y(x)-t|^3 \rightarrow 1$ except when $y(x)-t=0$,
since 0^0 doesn't exist.

Therefore $\int_{-\infty}^{\infty} |y(x)-t|^3 p(t|x) dt \rightarrow \int_{-\infty}^{\infty} p(t|x) dt \rightarrow 1$

except some points near $y(x)-t=0$.

At these points, the integral above is less than 1
because these points should be excluded from the
~~integration~~ation.

Thus, in order to minimize the value of this integral,
we need to have $y(x)=t$ as frequent as possible.
That is to say, to have this with the largest value
of $p(t|x)$, i.e. the most possible scenario.
Therefore, it's given by the conditional mode.

(# problem 7 see the
back of this page)

7.

Proof: $\forall v \in \mathbb{R}^m$

$$\begin{aligned} & (\Phi(\Phi^T\Phi)^{-1}\Phi^T)^2 v \\ &= \Phi(\Phi^T\Phi)^{-1}\Phi^T\Phi(\Phi^T\Phi)^{-1}\Phi^T v \\ &= \Phi(\Phi^T\Phi)^{-1}(\Phi^T\Phi)(\Phi^T\Phi)^{-1}\Phi^T v \\ &= \Phi(\Phi^T\Phi)^{-1}\Phi^T v \in \text{Image } \Phi(\Phi) \end{aligned}$$

So the matrix takes any v and projects it onto the space spanned by the columns of Φ .

Since $w^* = (\Phi^T\Phi)^{-1}\Phi^T t$

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

Check for orthogonality:

Suppose ϕ_i is i th column of Φ .

$$\Phi(\Phi^T\Phi)^{-1}\Phi^T \phi_i = \phi_i = I \phi_i$$

$$\begin{aligned} (\Phi w^* - t)^T \phi_i &= (\Phi(\Phi^T\Phi)^{-1}\Phi^T - I)^T \phi_i \\ &\stackrel{t^T}{=} t^T(I - I)\phi_i \\ &= 0 \end{aligned}$$

So such matrix is an orthogonal projection.

中 国 作 家 协 会

#8.

Proof: Know that $p(x|y) = h(x) g(\underline{y}) \exp(\eta^T u(x))$

$$-\nabla \ln g(y) = E(u(x))$$

$$= g(y) \int h(x) \exp(\eta^T u(x)) u(x) dx$$

$$= \cancel{g(y)} \int p(x|y) u(x) dx$$

So similarly,

$$-\nabla \nabla \ln g(y) = \cancel{\int} \{ h(x) g(y) \exp(\eta^T u(x)) u(x) u(x)^T \}$$

$$+ h(x) \nabla g(y) \exp(\eta^T u(x)) u(x) dx \quad \begin{matrix} \text{split into} \\ \downarrow \text{two parts} \end{matrix}$$

$$= \int h(x) g(y) \exp(\eta^T u(x)) u(x) u(x)^T dx + \int h(x) \nabla g(y) \exp(\eta^T u(x)) u(x) dx$$

$$= g(y) \int h(x) \exp(\eta^T u(x)) u(x) u(x)^T dx + \nabla g(y) \int h(x) \exp(\eta^T u(x)) u(x) dx$$

$$= E(u(x) u(x)^T) \stackrel{*}{=} E(u(x)) E(u(x)^T) \quad \text{by equations above} \quad \uparrow$$

$$= \text{cov}(u(x), u(x)^T)$$

$$= \text{cov}(u(x))$$

→ it's a vector here.

