

CISC 863 Statistical Machine Learning

Assignment One: Probability Theory and Distributions

Show that the variance of a sum is $\text{var}[X + Y] = \text{var}[X] + \text{var}[Y] + 2\text{cov}[X, Y]$.

$$\begin{aligned} \text{Ans: } \text{var}[X + Y] &= E[(X + Y - \mu_X - \mu_Y)^2] & \left. \begin{array}{l} \therefore \text{we know that,} \\ \text{var}[X] = E(X - \mu)^2 \end{array} \right\} \\ &= E[(X - \mu_X) + (Y - \mu_Y)]^2 \\ &= E[(X - \mu_X)^2 + (Y - \mu_Y)^2 + 2(X - \mu_X)(Y - \mu_Y)] \\ &= E[(X - \mu_X)^2] + E[(Y - \mu_Y)^2] + E[2(X - \mu_X)(Y - \mu_Y)] \\ &= \text{var}[X] + \text{var}[Y] + 2\text{cov}[X, Y] \quad (\text{showed}) \end{aligned}$$

Suppose $\theta \sim \text{Beta}(\alpha, \beta)$, derive the mean, mode and variance.

Mean: From given eqn of expected value of mean we can write that,

$$\begin{aligned}
 E[X] &= \int_0^1 x \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)} dx \dots \dots 1 \\
 &= \frac{1}{B(\alpha, \beta)} \left[\frac{x^{\alpha} (-1) (1-x)^{\beta}}{\beta} \Big|_0^1 - \int_0^1 \alpha x^{\alpha-1} (-1) \frac{(1-x)^{\beta}}{\beta} dx \right] \\
 &= \frac{1}{B(\alpha, \beta)} \int_0^1 \frac{\alpha}{\beta} x^{\alpha-1} (1-x)^{\beta-1} (1-x) dx \\
 &= \frac{1}{B(\alpha, \beta)} \left[\int_0^1 \frac{\alpha}{\beta} x^{\alpha-1} (1-x)^{\beta-1} dx - \int_0^1 x \frac{\alpha}{\beta} x^{\alpha-1} (1-x)^{\beta-1} dx \right] \\
 &= \frac{\alpha}{\beta} (1 - E[X])
 \end{aligned}$$

$$\therefore E[X] = \frac{\alpha}{\alpha + \beta} = \text{Expected value of mean.}$$

Mode: In order to calculate mode, we have to differentiate the Beta distribution and equal it to zero.

$$\text{mode} = \frac{d}{dx} \text{Beta}(\alpha, \beta) = 0 \Rightarrow \frac{d}{dx} (x^{\alpha-1} (1-x)^{\beta-1}) = 0$$

$$\Rightarrow \cancel{(\alpha-1)} x^{\alpha-2} \cdot \cancel{(\beta-1)} (1-x)^{\beta-2} \Rightarrow (\alpha-1) x^{\alpha-2} (1-x)^{\beta-1} - x^{\alpha-1} (\beta-1) (1-x)^{\beta-2} = 0$$

$$\Rightarrow (\alpha-1) (1-x) - x (\beta-1) = 0 \Rightarrow (\alpha-1) - x (\alpha-1 + \beta-1) = 0$$

$$\Rightarrow x = \frac{\alpha-1}{\alpha + \beta - 2}$$

Variance: We know that,

$$\text{var}[X] = E[X^2] - E[X]^2$$

Thus we need $E[X^2]$ to calculate.

$$E[X^2] = \int_0^1 x^2 \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)} dx$$

$$= \frac{1}{B(\alpha, \beta)} \left[-x^{\alpha+1} (1-x)^{\beta} \Big|_0^1 - \int_0^1 -(\alpha+1) x^{\alpha} (1-x)^{\beta} dx \right]$$

$$= \frac{\alpha+1}{\beta} \int_0^1 \frac{1}{B(\alpha, \beta)} x x^{\alpha-1} (1-x)^{\beta}$$

$$= \frac{\alpha+1}{\beta} \left[\int_0^1 \frac{1}{B(\alpha, \beta)} x x^{\alpha-1} (1-x)^{\beta-1} - \int_0^1 \frac{1}{B(\alpha, \beta)} x^2 x^{\alpha-1} (1-x)^{\beta-1} \right]$$

$$= \frac{\alpha+1}{\beta} \left[E[X] - E[X^2] \right] = \frac{\alpha+1}{\beta} \left[\frac{\alpha}{\alpha+\beta} - E[X^2] \right]$$

$$\therefore E[X^2] = \frac{\alpha+1}{\alpha+\beta+1} \cdot \frac{\alpha}{\alpha+\beta}$$

$$\therefore \text{var}[X] = \frac{\alpha+1}{\alpha+\beta+1} \cdot \frac{\alpha}{\alpha+\beta} - \frac{\alpha^2}{(\alpha+\beta)^2}$$

$$= \frac{\alpha^2 + \alpha}{(\alpha+\beta)(\alpha+\beta+1)} - \frac{\alpha^2}{(\alpha+\beta)^2}$$

$$= \frac{(\alpha^2 + \alpha)(\alpha+\beta) - \alpha^2(\alpha+\beta+1)}{(\alpha+\beta+1)(\alpha+\beta)^2}$$

$$= \frac{\alpha^3 + \alpha^2 + \alpha^2\beta + \alpha\beta - \alpha^3 - \alpha^2\beta - \alpha^2}{(\alpha+\beta+1)(\alpha+\beta)^2}$$

$$= \frac{\alpha\beta}{(\alpha+\beta+1)(\alpha+\beta)^2}$$

Since a positive definite matrix Σ can be defined as the quadratic form $U^T \Lambda U$, show that a necessary and sufficient condition for Σ to be positive definite is that all the eigenvalues λ_i of Λ are positive.

As we know that,

$\Sigma = U^T \Lambda U$ will be positive definite if $U^T \Lambda U > 0$

$$\text{Now, } U^T \Lambda U = [x_1 \dots x_N] \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}$$

$$= x_1^2 \lambda_1 + x_2^2 \lambda_2 + \dots + x_N^2 \lambda_N \dots \text{--- (1)}$$

① eqn needs to be ^{greater} zero to make Σ positive definite.

As we know x_i^2 must be greater than zero.

Thus we need to make sure that all eigen values λ_i must be ~~need~~ to be ~~zero~~ positive to make Σ positive definite.

Proof by contradiction:

Say one eigenvalue λ_i is zero then there will exist at least one eigenvector e such that $\Sigma e = 0$. So $e^T \Sigma e = 0$ and condition ~~$x^T M x > 0$~~ . $U^T \Lambda U > 0$ doesn't hold when $x = e$.

Another condition we may assume that if the matrix is negative i.e. $\lambda < 0$. Then there will exist at least one eigenvector e such that $\Sigma e = \lambda e$. So $e^T \Sigma e = e^T \lambda e = \lambda e^T e = \lambda |e|^2$.

Since $\lambda < 0$, condition $x^T M x > 0$ does not hold when $x = e$ so, it is necessary for the eigenvalues of a positive definite matrix to be positive.

Derive the maximum likelihood solutions for the mean and the variance of a univariate Gaussian distribution by maximize the log likelihood function with respect to μ and Σ .

Maximum likelihood solution for mean;

We know that, likelihood is,

$$P(x_1, x_2, \dots, x_N | \mu) = \prod_{i=1}^N P(x_i | \mu)$$
$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-(x_i - \mu)^2 / 2\sigma^2} \quad [\because \text{this is a gaussian distribution.}]$$

Now since this is a gaussian univariate distribution, we can maximize the log likelihood.

$$\cancel{P(x|\mu)} \quad \log(P(x_1, x_2, \dots, x_N | \mu)) = \sum_{i=1}^N \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{(x_i - \mu)^2}{2\sigma^2} \quad \text{--- ①}$$

we take the derivative with respect to, μ ,

$$\frac{d \log(P(x_1, x_2, \dots, x_N | \mu))}{d\mu} = \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2}.$$

Since we want to maximize assume that left side is '0';

$$0 = \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2}$$

$$\Rightarrow 0 = \sum_{i=1}^N (x_i - \mu)$$

$$\Rightarrow \sum_{i=1}^N \mu = \sum_{i=1}^N x_i$$

$$\Rightarrow N\mu = \sum_{i=1}^N x_i$$

$$\therefore \hat{\mu} = \frac{\sum_{i=1}^N x_i}{N}. \quad \text{Ans.}$$

In the same way we can take derivative of ① with respect to σ^2 ,

$$\frac{d \log(p(x_1, x_2, \dots, x_N | \mu))}{d \sigma^2} =$$

$$\frac{d}{d \sigma^2} \left[\sum_{i=1}^N \left\{ \log \left(\frac{1}{\sqrt{2\pi} \sigma} \right) - \frac{(x_i - \mu)^2}{2\sigma^2} \right\} \right]$$

$$= \frac{d}{d \sigma^2} \left\{ -\frac{N}{2} \log \left(\frac{1}{\sqrt{2\pi} \sigma} \right) - \frac{1}{2} \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^2} \right\}$$

$$= -\frac{1}{2} \log 2\pi - \frac{N}{2} \frac{d}{d \sigma^2} \log \sigma + \frac{1}{2} \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^3}$$

$$= -\frac{1}{2} \log 2\pi - \frac{1}{2} \frac{N}{\sigma} + \frac{1}{2} \sum_{i=1}^N \frac{(x_i - \mu)^2}{(\sigma^3)^2}$$

To maximize we can make left hand side zero,

$$0 = -\frac{1}{2} \frac{N}{\sigma} + \frac{1}{2} \sum_{i=1}^N \frac{(x_i - \mu)^2}{(\sigma^3)^2}$$

$$\Rightarrow \frac{N}{\sigma} = \frac{1}{\sigma^3} \sum_{i=1}^N (x_i - \mu)^2 \frac{1}{(\sigma^3)^2}$$

$$\Rightarrow \frac{\sigma^2}{N} = \frac{(\sigma^2)^2}{\sum_{i=1}^N (x_i - \mu)^2}$$

$$\Rightarrow \frac{1}{N} = \frac{\sigma^2}{\sum_{i=1}^N (x_i - \mu)^2} \quad [\because \sigma^2 \neq 0]$$

$$\Rightarrow \sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} \quad \text{Ans.}$$

Plot Gaussian likelihoods with unknown means, conjugate priors of the means, and their corresponding posterior distributions with tools you are comfortable with (e.g., Matlab, R, Octave) and different parameter settings

- Prior: $N(\mu|0,6)$
- Likelihood: $N(D|\mu, 10)$ (data file in mycourses)
- Posterior: $N(\mu|D, 0,6)$
- Predictive: $N(x^* = 2.4|D)$

For posterior we know,

$$\delta_N^v = \frac{1}{\left(\frac{N}{\delta^v}\right) + \left(\frac{1}{\delta_0^v}\right)} = \frac{1}{\left(\frac{20}{10}\right) + \left(\frac{1}{6}\right)} = \left(\frac{1}{2 + \frac{1}{6}}\right) = \frac{6}{13}$$

$$\mu_N = \left(\frac{N\bar{x}}{\delta^v} + \frac{\mu_0}{\delta_0^v}\right) \cdot \delta_N^v = \left(\frac{20}{6^v} \cdot \bar{x} + 0\right) \cdot \frac{6}{13} = \left(\frac{20}{10} \bar{x}\right) \cdot \frac{6}{13}$$

For predictive posterior,

$$\delta_{NP}^v = \delta_N^v + \delta^v = \frac{6}{13} + 10 = \frac{136}{13}$$

$\mu_N =$ same as posterior.

Scalar QDA: Consider the following training set of heights x (in inches) and gender y (male/female) of some college students: $x = \{67, 79, 71, 68, 67, 60\}$, and $y = \{m, m, m, f, f, f\}$.

- a. Fit a Bayes classifier to this data, using Maximum Likelihood estimation (MLE), i.e., estimate the parameters of the class conditional likelihoods

$$p(x | y = c) = N(x | \mu_c, \sigma_c)$$

And the class prior

$$p(y = c) = \pi_c$$

What are your values of μ_c, σ_c, π_c for $c = 'm', 'f'$? Show your work (so you can get partial credit if you make an arithmetic error).

- b. Compute $p(y = 'm' | x, \hat{\theta})$, where $x = 72$, and $\hat{\theta}$ are the MLE parameters. (This is called a plug-in prediction.)

Note: solve this exercise by hand AND using a computer (Matlab, Python, Octave, whatever). Show your work (derivation and code).

Hint: refer to the Probability Theory slides on QDA: 99-104

Ans: Here given, heights, $x = \{67, 79, 71, 68, 67, 60\}$
 $y = \{m, m, m, f, f, f\}$

$$p(x = m) = 0.5 = \pi_m \quad \mu_m = 72.3$$

$$p(x = f) = 0.5 = \pi_f \quad \mu_f = 65.0$$

As we know, variance,

$$\sigma_m^2 = \frac{\sum (x_m - \mu_m)^2}{\sum m} = 4.98$$

$$\sigma_f^2 = \frac{\sum (x_f - \mu_f)^2}{\sum f} = 3.56$$

$$P(y = 'm' | x = 72, \hat{\theta}) = \frac{p(x = 72 | y = m) \times \text{prior}(x = m)}{p(x = 72 | y = m) \times \text{prior}(x = m) + p(x = 72 | y = f) \times \text{prior}(y = f)}$$

$$= 0.83. \text{ Aug.}$$