

Notes on Probability Theory

Paul F. Roysdon, Ph.D.

I. DEFINITIONS

The *expected value*, or *mean*, of a discrete random variable (r.v. hereafter) \mathbf{X} :

$$\begin{aligned}\mathbf{m}_{\mathbf{X}} &= \mathbb{E} \langle \mathbf{X} \rangle \\ &= \sum_{k=1}^{\infty} \mathbf{x}_k P_{\mathbf{X}}(\mathbf{x}_k)\end{aligned}$$

The *variance* of the r.v. \mathbf{X} :

$$\begin{aligned}\sigma_{\mathbf{X}}^2 &= \text{var} \langle \mathbf{X} \rangle \\ &= \mathbb{E} \langle (\mathbf{X} - \mathbf{m}_{\mathbf{X}})^2 \rangle \\ &= \sum_{k=1}^{\infty} (\mathbf{x}_k - \mathbf{m}_{\mathbf{X}})^2 P_{\mathbf{X}}(\mathbf{x}_k)\end{aligned}$$

which can also be expressed as

$$\begin{aligned}\text{var} \langle \mathbf{X} \rangle &= \mathbb{E} \langle (\mathbf{X} - \mathbf{m}_{\mathbf{X}})^2 \rangle \\ &= \mathbb{E} \langle \mathbf{X}^2 - 2\mathbf{m}_{\mathbf{X}}\mathbf{X} + \mathbf{m}_{\mathbf{X}}^2 \rangle \\ &= \mathbb{E} \langle \mathbf{X}^2 \rangle - 2\mathbf{m}_{\mathbf{X}}\mathbb{E} \langle \mathbf{X} \rangle + \mathbf{m}_{\mathbf{X}}^2 \\ &= \mathbb{E} \langle \mathbf{X}^2 \rangle - \mathbf{m}_{\mathbf{X}}^2\end{aligned}$$

where $\mathbb{E} \langle \mathbf{X}^2 \rangle$ is called the second moment of \mathbf{X} .

The *standard deviation* of the r.v. \mathbf{X} :

$$\begin{aligned}\sigma_{\mathbf{X}} &= \text{std} \langle \mathbf{X} \rangle \\ &= \sqrt{\text{var} \langle \mathbf{X} \rangle}\end{aligned}$$

For the discrete r.v. $\mathbf{x} \in \mathbb{R}^{m \times 1}$ and scalar variable c , the variance has the following properties, each resulting in a *scalar random variable*:

- 1) Adding a constant c , does not affect the variance of \mathbf{X}

$$\begin{aligned}\text{var} \langle \mathbf{x} + c \rangle &= \mathbb{E} \langle (\mathbf{x} + c - (\mathbb{E} \langle \mathbf{x} \rangle + c))^2 \rangle \\ &= \mathbb{E} \langle (\mathbf{x} - \mathbb{E} \langle \mathbf{x} \rangle)^2 \rangle \\ &= \text{var} \langle \mathbf{x} \rangle\end{aligned}$$

- 2) Multiplying by a constant c , scales the variance of \mathbf{x} by c^2

$$\begin{aligned}\text{var} \langle c\mathbf{x} \rangle &= \mathbb{E} \langle (c\mathbf{x} - c\mathbb{E} \langle \mathbf{x} \rangle)^2 \rangle \\ &= \mathbb{E} \langle c^2(\mathbf{x} - \mathbb{E} \langle \mathbf{x} \rangle)^2 \rangle \\ &= c^2 \text{var} \langle \mathbf{x} \rangle\end{aligned}$$

- 3) Let $\mathbf{x} = c$, then the constant r.v. has zero variance

$$\begin{aligned}\text{var} \langle \mathbf{x} \rangle &= \mathbb{E} \langle (\mathbf{x} - c)^2 \rangle \\ &= \mathbb{E} \langle 0 \rangle \\ &= 0\end{aligned}$$

For the discrete r.v. $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m \times 1}$, the covariance has the following property, resulting in a *matrix of random variables*:

$$\begin{aligned}\text{cov} \langle \mathbf{x}, \mathbf{y} \rangle &= \mathbb{E} \langle \mathbf{x}\mathbf{y} - \mathbf{x}\mathbb{E} \langle \mathbf{y} \rangle - \mathbf{y}\mathbb{E} \langle \mathbf{x} \rangle + \mathbb{E} \langle \mathbf{x} \rangle \mathbb{E} \langle \mathbf{y} \rangle \rangle \\ &= \mathbb{E} \langle \mathbf{x}\mathbf{y} \rangle - 2\mathbb{E} \langle \mathbf{x} \rangle \mathbb{E} \langle \mathbf{y} \rangle + \mathbb{E} \langle \mathbf{x} \rangle \mathbb{E} \langle \mathbf{y} \rangle \\ &= \mathbb{E} \langle \mathbf{x}\mathbf{y} \rangle - \mathbb{E} \langle \mathbf{x} \rangle \mathbb{E} \langle \mathbf{y} \rangle\end{aligned}$$

If either \mathbf{x} or \mathbf{y} have zero mean, i.e. $\mathbb{E} \langle \mathbf{x} \rangle = 0$ or $\mathbb{E} \langle \mathbf{y} \rangle = 0$, then $\text{cov} \langle \mathbf{x}, \mathbf{y} \rangle = \mathbb{E} \langle \mathbf{x}\mathbf{y} \rangle$. If the correlation of \mathbf{x} and \mathbf{y} is zero, i.e. $\mathbb{E} \langle \mathbf{x}\mathbf{y} \rangle = 0$ then \mathbf{x} and \mathbf{y} are orthogonal.

For the discrete r.v. $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m \times 1}$, the correlation coefficient has the following property

$$\begin{aligned}\rho_{\mathbf{x}, \mathbf{y}} &= \frac{\text{cov} \langle \mathbf{x}, \mathbf{y} \rangle}{\sigma_{\mathbf{x}} \sigma_{\mathbf{y}}} \\ &= \frac{\mathbb{E} \langle \mathbf{x}\mathbf{y} \rangle - \mathbb{E} \langle \mathbf{x} \rangle \mathbb{E} \langle \mathbf{y} \rangle}{\sigma_{\mathbf{x}} \sigma_{\mathbf{y}}}\end{aligned}$$

where $\sigma_{\mathbf{x}} = \sqrt{\text{var} \langle \mathbf{x} \rangle}$, $\sigma_{\mathbf{y}} = \sqrt{\text{var} \langle \mathbf{y} \rangle}$, and $-1 \leq \rho_{\mathbf{x}, \mathbf{y}} \leq 1$. If $\rho_{\mathbf{x}, \mathbf{y}} = 0$ then \mathbf{x} and \mathbf{y} are said to be *uncorrelated*.

II. LINEAR COMBINATIONS

A. Linear Forms

Assume \mathbf{X} and \mathbf{x} to be a matrix and a vector of random variables. Then

$$\begin{aligned}\mathbb{E} \langle \mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C} \rangle &= \mathbf{A} \mathbb{E} \langle \mathbf{X} \rangle \mathbf{B} + \mathbf{C} \\ \text{var} \langle \mathbf{A}\mathbf{x} \rangle &= \mathbf{A} \text{var} \langle \mathbf{x} \rangle \mathbf{A}^{\top} \\ \text{cov} \langle \mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y} \rangle &= \mathbf{A} \text{cov} \langle \mathbf{x}, \mathbf{y} \rangle \mathbf{B}^{\top}\end{aligned}$$

Assume \mathbf{x} to be a stochastic vector with mean \mathbf{m} , then

$$\begin{aligned}\mathbb{E} \langle \mathbf{A}\mathbf{x} + \mathbf{b} \rangle &= \mathbf{A}\mathbf{m} + \mathbf{b} \\ \mathbb{E} \langle \mathbf{A}\mathbf{x} \rangle &= \mathbf{A}\mathbf{m} \\ \mathbb{E} \langle \mathbf{x} + \mathbf{b} \rangle &= \mathbf{m} + \mathbf{b}\end{aligned}$$

B. Quadratic Forms

Assume \mathbf{A} is symmetric, $\mathbf{c} = \mathbb{E} \langle \mathbf{x} \rangle$ and $\mathbf{\Sigma} = \text{var} \langle \mathbf{x} \rangle$. Assume also that all coordinates x_i are independent, have the same central moments $\mu_1, \mu_2, \mu_3, \mu_4$ and denote $\mathbf{a} = \text{diag}(\mathbf{A})$. Then

$$\begin{aligned}\mathbb{E} \langle \mathbf{x}^{\top} \mathbf{A} \mathbf{x} \rangle &= \text{Tr}(\mathbf{A}\mathbf{\Sigma}) + \mathbf{c}^{\top} \mathbf{A} \mathbf{c} \\ \text{var} \langle \mathbf{x}^{\top} \mathbf{A} \mathbf{x} \rangle &= 2\mu_2^2 \text{Tr}(\mathbf{A}^2) + 4\mu_2 \mathbf{c}^{\top} \mathbf{A}^2 \mathbf{c} \\ &\quad + 4\mu_3 \mathbf{c}^{\top} \mathbf{A} \mathbf{a} + (\mu_4 - 3\mu_2^2) \mathbf{a}^{\top} \mathbf{a}\end{aligned}$$

Also, assume \mathbf{x} to be a stochastic vector with mean \mathbf{m} , and covariance \mathbf{M} . Then

$$\begin{aligned}
\mathbb{E} \langle (\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{B}\mathbf{x} + \mathbf{b})^\top \rangle &= \mathbf{A}\mathbf{M}\mathbf{B}^\top + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{B}\mathbf{m} + \mathbf{b})^\top \\
\mathbb{E} \langle \mathbf{x}\mathbf{x}^\top \rangle &= \mathbf{M} + \mathbf{m}\mathbf{m}^\top \\
\mathbb{E} \langle \mathbf{x}\mathbf{a}^\top \rangle &= (\mathbf{M} + \mathbf{m}\mathbf{m}^\top)\mathbf{a} \\
\mathbb{E} \langle \mathbf{x}^\top \mathbf{a}\mathbf{x} \rangle &= \mathbf{a}^\top (\mathbf{M} + \mathbf{m}\mathbf{m}^\top) \\
\mathbb{E} \langle (\mathbf{A}\mathbf{x})(\mathbf{A}\mathbf{x})^\top \rangle &= \mathbf{A}(\mathbf{M} + \mathbf{m}\mathbf{m}^\top)\mathbf{A}^\top \\
\mathbb{E} \langle (\mathbf{x} + \mathbf{a})(\mathbf{x} + \mathbf{a})^\top \rangle &= \mathbf{M} + (\mathbf{m} + \mathbf{a})(\mathbf{m} + \mathbf{a})^\top \\
\mathbb{E} \langle (\mathbf{A}\mathbf{x} + \mathbf{a})^\top (\mathbf{B}\mathbf{x} + \mathbf{b}) \rangle &= \text{Tr}(\mathbf{A}\mathbf{M}\mathbf{B}^\top) + (\mathbf{A}\mathbf{m} + \mathbf{a})^\top (\mathbf{B}\mathbf{m} + \mathbf{b}) \\
\mathbb{E} \langle \mathbf{x}^\top \mathbf{x} \rangle &= \text{Tr}(\mathbf{M}) + \mathbf{m}^\top \mathbf{m} \\
\mathbb{E} \langle \mathbf{x}^\top \mathbf{A}\mathbf{x} \rangle &= \text{Tr}(\mathbf{A}\mathbf{M}) + \mathbf{m}^\top \mathbf{A}\mathbf{m} \\
\mathbb{E} \langle (\mathbf{A}\mathbf{x})^\top (\mathbf{A}\mathbf{x}) \rangle &= \text{Tr}(\mathbf{A}\mathbf{M}\mathbf{A}^\top) + (\mathbf{A}\mathbf{m})^\top (\mathbf{A}\mathbf{m}) \\
\mathbb{E} \langle (\mathbf{x} + \mathbf{a})^\top (\mathbf{x} + \mathbf{a}) \rangle &= \text{Tr}(\mathbf{M}) + (\mathbf{m} + \mathbf{a})^\top (\mathbf{m} + \mathbf{a})
\end{aligned}$$

C. Cubic Forms

Assume \mathbf{x} to be a stochastic vector with independent coordinates, mean \mathbf{m} , covariance \mathbf{M} and central moments $\mathbf{v}_3 = \mathbb{E} \langle (\mathbf{x} - \mathbf{m})^3 \rangle$. Then

$$\begin{aligned}
&\mathbb{E} \langle (\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{B}\mathbf{x} + \mathbf{b})^\top (\mathbf{C}\mathbf{x} + \mathbf{c}) \rangle \\
&= \mathbf{A} \text{diag}(\mathbf{B}^\top \mathbf{C}) \mathbf{v}_3 \\
&\quad + \text{Tr}(\mathbf{B}\mathbf{M}\mathbf{C}^\top)(\mathbf{A}\mathbf{m} + \mathbf{a}) \\
&\quad + \mathbf{A}\mathbf{M}\mathbf{C}^\top (\mathbf{B}\mathbf{m} + \mathbf{b}) \\
&\quad + (\mathbf{A}\mathbf{M}\mathbf{B}^\top + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{B}\mathbf{m} + \mathbf{b})^\top)(\mathbf{C}\mathbf{m} + \mathbf{c}) \\
&\mathbb{E} \langle \mathbf{x}\mathbf{x}^\top \mathbf{x} \rangle \\
&= \mathbf{v}_3 + 2\mathbf{M}\mathbf{m} + (\text{Tr}(\mathbf{M}) + \mathbf{m}^\top \mathbf{m})\mathbf{m} \\
&\mathbb{E} \langle (\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{A}\mathbf{x} + \mathbf{a})^\top (\mathbf{A}\mathbf{x} + \mathbf{a}) \rangle \\
&= \mathbf{A} \text{diag}(\mathbf{A}^\top \mathbf{A}) \mathbf{v}_3 \\
&\quad + [2\mathbf{A}\mathbf{M}\mathbf{A}^\top + (\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{A}\mathbf{x} + \mathbf{a})^\top](\mathbf{A}\mathbf{m} + \mathbf{a}) \\
&\quad + \text{Tr}(\mathbf{A}\mathbf{M}\mathbf{A}^\top)(\mathbf{A}\mathbf{m} + \mathbf{a}) \\
&\mathbb{E} \langle (\mathbf{A}\mathbf{x} + \mathbf{a})\mathbf{b}^\top (\mathbf{C}\mathbf{x} + \mathbf{c})(\mathbf{D}\mathbf{x} + \mathbf{d})^\top \rangle \\
&= (\mathbf{A}\mathbf{x} + \mathbf{a})\mathbf{b}^\top (\mathbf{C}\mathbf{M}\mathbf{D}^\top + (\mathbf{C}\mathbf{m} + \mathbf{c})(\mathbf{D}\mathbf{m} + \mathbf{d})^\top) \\
&\quad + (\mathbf{A}\mathbf{M}\mathbf{C}^\top + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{C}\mathbf{m} + \mathbf{c})^\top)\mathbf{b}(\mathbf{D}\mathbf{m} + \mathbf{d})^\top \\
&\quad + \mathbf{b}^\top (\mathbf{C}\mathbf{m} + \mathbf{c})(\mathbf{A}\mathbf{M}\mathbf{D}^\top - (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{D}\mathbf{m} + \mathbf{d})^\top)
\end{aligned}$$

III. GAUSSIANS

The following section contains some useful relations from statistics theory.

A. Density and normalization

The density of $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{\Sigma})$ is

$$p(\mathbf{x}) = \frac{1}{\sqrt{\det(2\pi\mathbf{\Sigma})}} \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{m})^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{m}) \right]$$

Note that if \mathbf{x} is d -dimensional, then $\det(2\pi\mathbf{\Sigma}) = (2\pi)^d \det(\mathbf{\Sigma})$.

Integration and normalization is given by:

$$\begin{aligned}
&\int \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{m})^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{m}) \right] d\mathbf{x} \\
&= \sqrt{\det(2\pi\mathbf{\Sigma})} \\
&\int \exp \left[-\frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x} + \mathbf{b}^\top \mathbf{x} \right] d\mathbf{x} \\
&= \sqrt{\det(2\pi\mathbf{A}^{-1})} \exp \left[\frac{1}{2}\mathbf{b}^\top \mathbf{A}^{-1}\mathbf{b} \right] \\
&\int \exp \left[-\frac{1}{2}\text{Tr}(\mathbf{S}^\top \mathbf{A}\mathbf{S}) + \text{Tr}(\mathbf{B}^\top \mathbf{S}) \right] d\mathbf{S} \\
&= \sqrt{\det(2\pi\mathbf{A}^{-1})} \exp \left[\frac{1}{2}\text{Tr}(\mathbf{B}^\top \mathbf{A}^{-1}\mathbf{B}) \right]
\end{aligned}$$

The derivatives of the density are

$$\begin{aligned}
\frac{\partial p(\mathbf{x})}{\partial \mathbf{x}} &= -p(\mathbf{x})\mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{m}) \\
\frac{\partial^2 p}{\partial \mathbf{x} \partial \mathbf{x}^\top} &= p(\mathbf{x}) \left(\mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^\top \mathbf{\Sigma}^{-1} - \mathbf{\Sigma}^{-1} \right)
\end{aligned}$$

B. Marginal Distribution

Assume $\mathbf{x} \sim \mathcal{N}_x(\boldsymbol{\mu}, \mathbf{\Sigma})$ where

$$\begin{aligned}
\mathbf{x} &= \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix} \\
\boldsymbol{\mu} &= \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix} \\
\mathbf{\Sigma} &= \begin{bmatrix} \mathbf{\Sigma}_a & \mathbf{\Sigma}_c \\ \mathbf{\Sigma}_c^\top & \mathbf{\Sigma}_b \end{bmatrix}
\end{aligned}$$

then

$$\begin{aligned}
p(\mathbf{x}_a) &= \mathcal{N}_{\mathbf{x}_a}(\boldsymbol{\mu}_a, \mathbf{\Sigma}_a) \\
p(\mathbf{x}_b) &= \mathcal{N}_{\mathbf{x}_b}(\boldsymbol{\mu}_b, \mathbf{\Sigma}_b)
\end{aligned}$$

C. Conditional Distribution

Assume $\mathbf{x} \sim \mathcal{N}_x(\boldsymbol{\mu}, \mathbf{\Sigma})$ where

$$\begin{aligned}
\mathbf{x} &= \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix} \\
\boldsymbol{\mu} &= \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix} \\
\mathbf{\Sigma} &= \begin{bmatrix} \mathbf{\Sigma}_a & \mathbf{\Sigma}_c \\ \mathbf{\Sigma}_c^\top & \mathbf{\Sigma}_b \end{bmatrix}
\end{aligned}$$

then

$$\begin{aligned}
p(\mathbf{x}_a|\mathbf{x}_b) &= \mathcal{N}_{\mathbf{x}_a}(\hat{\boldsymbol{\mu}}_a, \hat{\mathbf{\Sigma}}_a) \\
p(\mathbf{x}_b|\mathbf{x}_a) &= \mathcal{N}_{\mathbf{x}_b}(\hat{\boldsymbol{\mu}}_b, \hat{\mathbf{\Sigma}}_b)
\end{aligned}$$

where

$$\begin{aligned}
\hat{\boldsymbol{\mu}}_a &= \boldsymbol{\mu}_a + \mathbf{\Sigma}_c \mathbf{\Sigma}_b^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\
\hat{\mathbf{\Sigma}}_a &= \mathbf{\Sigma}_a - \mathbf{\Sigma}_c \mathbf{\Sigma}_b^{-1} \mathbf{\Sigma}_c^\top \\
\hat{\boldsymbol{\mu}}_b &= \boldsymbol{\mu}_b + \mathbf{\Sigma}_c \mathbf{\Sigma}_a^{-1}(\mathbf{x}_a - \boldsymbol{\mu}_a) \\
\hat{\mathbf{\Sigma}}_b &= \mathbf{\Sigma}_b - \mathbf{\Sigma}_c \mathbf{\Sigma}_a^{-1} \mathbf{\Sigma}_c^\top
\end{aligned}$$

D. Linear combination

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \Sigma_x)$ and $\mathbf{y} \sim \mathcal{N}(\mathbf{m}_y, \Sigma_y)$

$$\begin{aligned} \mathbf{Ax} + \mathbf{By} + \mathbf{c} \\ \sim \mathcal{N}(\mathbf{Am}_x + \mathbf{Bm}_y + \mathbf{c}, \mathbf{A}\Sigma_x\mathbf{A}^\top + \mathbf{B}\Sigma_y\mathbf{B}^\top) \end{aligned}$$

E. Rearranging Means

$$\begin{aligned} \mathcal{N}_{\mathbf{Ax}}[\mathbf{m}, \Sigma] &= \frac{\sqrt{\det(2\pi(\mathbf{A}^\top \Sigma^{-1} \mathbf{A})^{-1})}}{\sqrt{\det(2\pi \Sigma)}} \\ &\times \mathcal{N}_{\mathbf{x}}[\mathbf{A}^{-1}\mathbf{m}, (\mathbf{A}^\top \Sigma^{-1} \mathbf{A})^{-1}] \end{aligned}$$

F. Rearranging into squared form

If \mathbf{A} is symmetric, then

$$\begin{aligned} -\frac{1}{2}\mathbf{x}^\top \mathbf{Ax} + \mathbf{b}^\top \mathbf{x} \\ = -\frac{1}{2}(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b})^\top \mathbf{A}(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}) + \frac{1}{2}\mathbf{b}^\top \mathbf{A}^{-1}\mathbf{b} \\ -\frac{1}{2}\text{Tr}(\mathbf{X}^\top \mathbf{AX}) + \text{Tr}(\mathbf{B}^\top \mathbf{X}) \\ = -\frac{1}{2}[(\mathbf{X} - \mathbf{A}^{-1}\mathbf{B})^\top \mathbf{A}(\mathbf{X} - \mathbf{A}^{-1}\mathbf{B})] + \frac{1}{2}\text{Tr}(\mathbf{B}^\top \mathbf{A}^{-1}\mathbf{B}) \end{aligned}$$

G. Sum of two squared forms

In vector form (assuming Σ_1, Σ_2 are symmetric)

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^\top \Sigma_1^{-1}(\mathbf{x} - \mathbf{m}_1) \\ -\frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^\top \Sigma_2^{-1}(\mathbf{x} - \mathbf{m}_2) \\ = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_c)^\top \Sigma_c^{-1}(\mathbf{x} - \mathbf{m}_c) + C \\ \Sigma_c^{-1} = \Sigma_1^{-1} + \Sigma_2^{-1} \\ \mathbf{m}_c = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}(\Sigma_1^{-1}\mathbf{m}_1 + \Sigma_2^{-1}\mathbf{m}_2) \\ C = \frac{1}{2}(\mathbf{m}_1^\top \Sigma_1^{-1} + \mathbf{m}_2^\top \Sigma_2^{-1})(\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \\ (\Sigma_1^{-1}\mathbf{m}_1 + \Sigma_2^{-1}\mathbf{m}_2) - \frac{1}{2}(\mathbf{m}_1^\top \Sigma_1^{-1}\mathbf{m}_1 + \mathbf{m}_2^\top \Sigma_2^{-1}\mathbf{m}_2) \end{aligned}$$

In a trace form (assuming Σ_1, Σ_2 are symmetric)

$$\begin{aligned} -\frac{1}{2}\text{Tr}[(\mathbf{X} - \mathbf{M}_1)^\top \Sigma_1^{-1}(\mathbf{X} - \mathbf{M}_1)] \\ -\frac{1}{2}\text{Tr}[(\mathbf{X} - \mathbf{M}_2)^\top \Sigma_2^{-1}(\mathbf{X} - \mathbf{M}_2)] \\ = -\frac{1}{2}\text{Tr}[(\mathbf{X} - \mathbf{M}_c)^\top \Sigma_c^{-1}(\mathbf{X} - \mathbf{M}_c)] + C \\ \Sigma_c^{-1} = \Sigma_1^{-1} + \Sigma_2^{-1} \\ \mathbf{M}_c = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}(\Sigma_1^{-1}\mathbf{M}_1 + \Sigma_2^{-1}\mathbf{M}_2) \\ C = \frac{1}{2}\text{Tr}[(\mathbf{M}_1^\top \Sigma_1^{-1} + \mathbf{M}_2^\top \Sigma_2^{-1})(\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \\ (\Sigma_1^{-1}\mathbf{M}_1 + \Sigma_2^{-1}\mathbf{M}_2)] - \frac{1}{2}(\mathbf{M}_1^\top \Sigma_1^{-1}\mathbf{M}_1 \\ + \mathbf{M}_2^\top \Sigma_2^{-1}\mathbf{M}_2) \end{aligned}$$

H. Product of Gaussian densities

Let $\mathcal{N}_{\mathbf{x}}(\mathbf{m}, \Sigma)$ denote a density of \mathbf{x} , then

$$\mathcal{N}_{\mathbf{x}}(\mathbf{m}_1, \Sigma_1) \cdot \mathcal{N}_{\mathbf{x}}(\mathbf{m}_2, \Sigma_2) = c_c \mathcal{N}_{\mathbf{x}}(\mathbf{m}_c, \Sigma_c)$$

$$\begin{aligned} c_c &= \mathcal{N}_{\mathbf{m}_1}(\mathbf{m}_1, (\Sigma_1 + \Sigma_2)) \\ &= \frac{1}{\sqrt{\det(2\pi(\Sigma_1 + \Sigma_2))}} \\ &\quad \exp\left[-\frac{1}{2}(\mathbf{m}_1 - \mathbf{m}_2)^\top (\Sigma_1 + \Sigma_2)^{-1}(\mathbf{m}_1 - \mathbf{m}_2)\right] \\ \mathbf{m}_c &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}(\Sigma_1^{-1}\mathbf{m}_1 + \Sigma_2^{-1}\mathbf{m}_2) \\ \Sigma_c &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \end{aligned}$$

IV. MOMENTS

A. Mean and covariance of linear forms

First and second moments. Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$\mathbb{E} \langle \mathbf{x} \rangle = \mathbf{m} \quad (1)$$

then

$$\begin{aligned} \text{Cov} \langle \mathbf{x}, \mathbf{x} \rangle &= \text{Var} \langle \mathbf{x} \rangle \\ &= \Sigma \\ &= \mathbb{E} \langle \mathbf{x} \mathbf{x}^\top \rangle - \mathbb{E} \langle \mathbf{x} \rangle \mathbb{E} \langle \mathbf{x}^\top \rangle \\ &= \mathbb{E} \langle \mathbf{x} \mathbf{x}^\top \rangle - \mathbf{m} \mathbf{m}^\top \end{aligned}$$

As for any other distribution is holds for gaussians that

$$\begin{aligned} \mathbb{E} \langle \mathbf{Ax} \rangle &= \mathbf{A} \mathbb{E} \langle \mathbf{x} \rangle \\ \text{Var} \langle \mathbf{Ax} \rangle &= \mathbf{A} \text{Var} \langle \mathbf{x} \rangle \mathbf{A}^\top \\ \text{Cov} \langle \mathbf{Ax}, \mathbf{By} \rangle &= \mathbf{A} \text{Cov} \langle \mathbf{x}, \mathbf{y} \rangle \mathbf{B}^\top \end{aligned}$$

B. Mean and variance of square forms

Mean and variance of square forms. Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$\begin{aligned} \mathbb{E} \langle \mathbf{x} \mathbf{x}^\top \rangle &= \Sigma + \mathbf{m} \mathbf{m}^\top \\ \mathbb{E} \langle \mathbf{x}^\top \mathbf{Ax} \rangle &= \text{Tr}(\mathbf{A}\Sigma) + \mathbf{m}^\top \mathbf{A} \mathbf{m} \\ \text{Var} \langle \mathbf{x}^\top \mathbf{Ax} \rangle &= 2\sigma^4 \text{Tr}(\mathbf{A}^2) + 4\sigma^2 \mathbf{m}^\top \mathbf{A}^2 \mathbf{m} \\ \text{Cov} \langle (\mathbf{x} - \mathbf{m}')^\top \mathbf{A}(\mathbf{x} - \mathbf{m}') \rangle &= (\mathbf{x} - \mathbf{m}')^\top \mathbf{A}(\mathbf{x} - \mathbf{m}') + \text{Tr}(\mathbf{A}\Sigma) \end{aligned}$$

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ and \mathbf{A} and \mathbf{B} to be symmetric, then

$$\text{Cov} \langle \mathbf{x}^\top \mathbf{Ax}, \mathbf{x}^\top \mathbf{Bx} \rangle = 2\sigma^4 \text{Tr}(\mathbf{AB})$$

C. Cubic forms

$$\begin{aligned} \mathbb{E} \langle \mathbf{x} \mathbf{b}^\top \mathbf{x} \mathbf{x}^\top \rangle &= \mathbf{m} \mathbf{b}^\top (\mathbf{M} + \mathbf{m} \mathbf{m}^\top) + (\mathbf{M} + \mathbf{m} \mathbf{m}^\top) \mathbf{b} \mathbf{m}^\top \\ &\quad + \mathbf{b}^\top \mathbf{m} (\mathbf{M} - \mathbf{m} \mathbf{m}^\top) \end{aligned}$$

D. Moments

$$\begin{aligned} \mathbb{E} \langle \mathbf{x} \rangle &= \sum_k \rho_k \mathbf{m}_k \\ \text{Cov} \langle \mathbf{x} \rangle &= \sum_k \sum_{k'} \rho_k \rho_{k'} (\Sigma_k + \mathbf{m}_k \mathbf{m}_k^\top + \mathbf{m}_{k'} \mathbf{m}_{k'}^\top) \end{aligned}$$

V. WHITENING

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{\Sigma})$, then

$$\mathbf{z} = \mathbf{\Sigma}^{-1/2}(\mathbf{x} - \mathbf{m}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Conversely having $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ one can generate data $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{\Sigma})$ by setting

$$\mathbf{x} = \mathbf{\Sigma}^{1/2}\mathbf{z} + \mathbf{m} \sim \mathcal{N}(\mathbf{m}, \mathbf{\Sigma})$$

Note that $\mathbf{\Sigma}^{1/2}$ means the matrix which fulfills $\mathbf{\Sigma}^{1/2}\mathbf{\Sigma}^{1/2} = \mathbf{\Sigma}$, and that it exists and is unique since $\mathbf{\Sigma}$ is positive definite.

VI. CHI-SQUARE

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{\Sigma})$ and \mathbf{x} to be n dimensional, then

$$z = (\mathbf{x} - \mathbf{m})^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{m}) \sim \chi_n^2$$

where χ_n^2 denotes the Chi square distribution with n degrees of freedom.

VII. MIXTURE OF GAUSSIANS

A. Density

The variable \mathbf{x} is distributed as a mixture of gaussians if it has the density

$$p(\mathbf{x}) = \sum_{k=1}^K \rho_k \frac{1}{\sqrt{\det(2\pi\mathbf{\Sigma}_k)}} \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{m}_k)^\top \mathbf{\Sigma}_k^{-1}(\mathbf{x} - \mathbf{m}_k) \right]$$

where ρ_k sum to 1 and the $\mathbf{\Sigma}_k$ all are positive definite.

B. Derivatives

Defining $p(s) = \sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)$, then

$$\begin{aligned} \frac{\partial \ln p(s)}{\partial \rho_j} &= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)} \frac{\partial}{\partial \rho_j} \ln[\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)] \\ &= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)} \frac{1}{\rho_j} \\ \frac{\partial \ln p(s)}{\partial \boldsymbol{\mu}_j} &= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_j} \ln[\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)] \\ &= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)} [-\mathbf{\Sigma}_k^{-1}(s - \boldsymbol{\mu}_k)] \\ \frac{\partial \ln p(s)}{\partial \mathbf{\Sigma}_j} &= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)} \frac{\partial}{\partial \mathbf{\Sigma}_j} \ln[\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)] \\ &= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \mathbf{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \mathbf{\Sigma}_k)} \\ &\quad \frac{1}{2} [\mathbf{\Sigma}_j^{-\top} + \mathbf{\Sigma}_j^{-\top}(s - \boldsymbol{\mu}_j)(s - \boldsymbol{\mu}_j)^\top \mathbf{\Sigma}_j^{-\top}] \end{aligned}$$

Note, ρ_k and $\mathbf{\Sigma}_k$ must be constrained.