

Lecture 13: Transformations of Random Variables

In this lecture we will discuss transformations of random vectors and review some additional properties of random vectors. These concepts are essential for understanding the state estimation approaches (e.g., Kalman filtering) to be presented later.

Linear Transformation of a Gaussian Random Variable

Consider a Gaussian random variable $x \sim \mathcal{N}(\mu_x, \sigma_x^2)$ and the linear scalar-valued transformation

$$y = ax + b. \quad (1)$$

The r.v. is transformed through the above function to yield a new r.v. $y \sim \mathcal{N}(\mu_y, \sigma_y^2)$. To find the mean μ_y we take the expected operator of both sides:

$$\mu_y = E[y] = E[ax + b] \quad (2)$$

$$= aE[x] + E[b] \quad (3)$$

$$\implies \mu_y = a\mu_x + b \quad (4)$$

and the variance σ_y^2 is similarly obtained by applying the definition of variance to the new r.v. y with (1) and (4):

$$\begin{aligned} \sigma_y^2 &= E[(y - \mu_y)^2] \\ &= E[(ax + b - (a\mu_x + b))^2] \\ &= E[(ax - a\mu_x)^2] \\ &= E[a^2(x - \mu_x)^2] \\ &= a^2E[(x - \mu_x)^2] \\ &= a^2\sigma_x^2 \\ \implies \sigma_y &= a\sigma_x \end{aligned}$$

Thus, the transformed r.v. y has a mean that is obtained by propagating $x = \mu_y$ through the function $y = ax + b$ and the variance scales by a factor of a^2 .

Aside: In MATLAB the function `randn()` generates a zero-mean and unit variance random number (i.e., from $\mathcal{N}(0, 1)$). We can use the property described above to draw a number from $\mathcal{N}(\mu_y, \sigma_y^2)$ by transforming `randn()` appropriately, for example:

```
mu_y = 2; % desired mean
var_y = 3; % desired variance
r = randn([100 1]); % a hundred unit variance random numbers
y = sqrt(var_y)*r + mu_y % random numbers with mean mu_y and variance var_y
```

Linear Transformation of a Gaussian Random Vector

Now we revisit the previous example for the case of a random vector instead of a random variable. Consider a Gaussian random vector $x \sim \mathcal{N}(\mu_x, P_x)$ with mean $\mu \in \mathbb{R}^n$ and covariance $P_x \in \mathbb{R}^{n \times n}$ that is transformed through the linear matrix equation:

$$y = Ax + b, \quad (5)$$

where $y \sim \mathcal{N}(\mu_y, P_y)$ is a new Gaussian random vector and $A \in \mathbb{R}^{n \times m}$ and $b \in \mathbb{R}^m$. To find the mean μ_y we take the expected operator of both sides of (5):

$$\mu_y = E[y] = E[Ax + b] \quad (6)$$

$$= AE[x] + E[b] \quad (7)$$

$$\implies \mu_y = A\mu_x + b \quad (8)$$

Similarly, we can obtain the covariance by using the definition of covariance of a random vector and substituting (5) and (8)

$$P_y = E[yy^T] - \mu_y\mu_y^T \quad (9)$$

$$= E[(Ax + b)(Ax + b)^T] - (A\mu_x + b)(A\mu_x + b)^T \quad (10)$$

$$= E[(Ax + b)(x^T A^T + b^T)] - (A\mu_x + b)(\mu_x^T A^T + b^T) \quad (11)$$

$$= E[(Ax x^T A^T + Ax b^T + b x^T A^T + b b^T)] - A\mu_x \mu_x^T A^T - A\mu_x b^T - b \mu_x^T A^T - b b^T \quad (12)$$

$$= AE[xx^T]A^T + AE[x]b^T + bE[x^T]A^T + b b^T - A\mu_x \mu_x^T A^T - A\mu_x b^T - b \mu_x^T A^T - b b^T \quad (13)$$

$$= AE[xx^T]A^T + A\mu_x b^T + b \mu_x^T A^T - A\mu_x \mu_x^T A^T - A\mu_x b^T - b \mu_x^T A^T \quad (14)$$

$$= AE[xx^T]A^T - A\mu_x \mu_x^T A^T \quad (15)$$

$$= A \left[E[xx^T] - \mu_x \mu_x^T \right] A^T \quad (16)$$

$$\implies P_y = AP_x A^T \quad (17)$$

Aside: We can use the above property to generate a random vector with a desired mean and covariance by transforming a zero-mean and unit-diagonal covariance random vector $x \sim \mathcal{N}(0, P_x)$ by the above linear equation where $P_x = \mathbf{1}_{n \times n}$ is the identity matrix. If the target distribution has a covariance P_y then we can obtain the desired transformation by finding a matrix A such that:

$$\begin{aligned} P_y &= AP_x A^T \\ &= A \mathbf{1} A^T \\ &= AA^T \end{aligned}$$

The A matrix can be found using Cholsky decomposition (i.e., chol in MATLAB).

Transformation of a Random Vector through a Nonlinear Function [1, Sec. 14.1]

Consider a random vector \mathbf{x} with mean $\boldsymbol{\mu}_x$ and a symmetric p.d.f. around the mean (e.g., a Gaussian p.d.f.). Now, suppose this random vector is transformed by a nonlinear function

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) \quad (18)$$

to give a new random vector \mathbf{y} ; we wish to determine its mean and covariance. Begin by expanding (18) as a Taylor series around a nominal value $\boldsymbol{\mu}_x$:

$$\mathbf{y} = \mathbf{h}(\boldsymbol{\mu}_x) + D_{\tilde{\mathbf{x}}} \mathbf{h} + \frac{1}{2!} D_{\tilde{\mathbf{x}}}^2 \mathbf{h} + \frac{1}{3!} D_{\tilde{\mathbf{x}}}^3 \mathbf{h} + \dots \quad (19)$$

where $\tilde{\mathbf{x}} = \mathbf{x} - \boldsymbol{\mu}_x$ is a new r.v. describing the deviation from the nominal value and the operator $D_{\tilde{\mathbf{x}}}^k \mathbf{f}$:

$$\begin{aligned} D_{\tilde{\mathbf{x}}}^k \mathbf{f} &= \left(\tilde{x}_1 \frac{\partial}{\partial x_1} + \tilde{x}_2 \frac{\partial}{\partial x_2} + \dots + \tilde{x}_n \frac{\partial}{\partial x_n} \right)^k \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \\ &= \left(\sum_{i=1}^n \tilde{x}_i \frac{\partial}{\partial x_i} \right)^k \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \end{aligned}$$

where the partial derivative applies element-wise to the column vector of functions $\mathbf{f}(\mathbf{x})$.

Note: In our earlier lecture on linearization of ODEs we saw this operator written in a slightly different form using the Jacobian, however both forms are equivalent:

$$(J_{\mathbf{x}} \mathbf{h}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \tilde{\mathbf{x}} = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \dots & \frac{\partial h_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_n}{\partial x_1} & \dots & \frac{\partial h_n}{\partial x_n} \end{bmatrix}_{\mathbf{x}=\boldsymbol{\mu}_x} \begin{bmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_n \end{bmatrix} \quad (20)$$

$$= \begin{bmatrix} \frac{\partial h_1}{\partial x_1} \tilde{x}_1 + \dots + \frac{\partial h_1}{\partial x_n} \tilde{x}_n \\ \vdots \\ \frac{\partial h_n}{\partial x_1} \tilde{x}_1 + \dots + \frac{\partial h_n}{\partial x_n} \tilde{x}_n \end{bmatrix}_{\mathbf{x}=\boldsymbol{\mu}_x} \quad (21)$$

$$= \tilde{x}_1 \begin{bmatrix} \frac{\partial h_1}{\partial x_1} \\ \vdots \\ \frac{\partial h_n}{\partial x_1} \end{bmatrix}_{\mathbf{x}=\boldsymbol{\mu}_x} + \dots + \tilde{x}_n \begin{bmatrix} \frac{\partial h_1}{\partial x_n} \\ \vdots \\ \frac{\partial h_n}{\partial x_n} \end{bmatrix}_{\mathbf{x}=\boldsymbol{\mu}_x} \quad (22)$$

$$= \tilde{x}_1 \left(\frac{\partial}{\partial x_1} \right) \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} + \dots + \tilde{x}_n \left(\frac{\partial}{\partial x_n} \right) \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \quad (23)$$

$$= \left(\tilde{x}_1 \frac{\partial}{\partial x_1} + \dots + \tilde{x}_n \frac{\partial}{\partial x_n} \right) \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \quad (24)$$

$$= D_{\tilde{\mathbf{x}}} \quad (25)$$

By construction, this new r.v. is zero-mean, $E[\tilde{x}] = 0$. To find the mean of the transformed r.v. \mathbf{y} take the expected value of both sides:

$$\boldsymbol{\mu}_y \triangleq E[\mathbf{y}] = E[\mathbf{h}(\boldsymbol{\mu}_x) + D_{\tilde{x}}\mathbf{h} + \frac{1}{2!}D_{\tilde{x}}^2\mathbf{h} + \frac{1}{3!}D_{\tilde{x}}^3\mathbf{h} + \dots] \quad (26)$$

$$= E[\mathbf{h}(\boldsymbol{\mu}_x)] + E[D_{\tilde{x}}\mathbf{h}] + E[\frac{1}{2!}D_{\tilde{x}}^2\mathbf{h}] + E[\frac{1}{3!}D_{\tilde{x}}^3\mathbf{h}] + E[\dots] \quad (27)$$

Since $\mathbf{h}(\boldsymbol{\mu}_x)$ is not a random variable then $E[\mathbf{h}(\boldsymbol{\mu}_x)] = \mathbf{h}(\boldsymbol{\mu}_x)$ and the first term in (27) can be simplified. To simplify the second term in (27), use the definition of $D_{\tilde{x}}\mathbf{h}$ and take its expected value:

$$E[D_{\tilde{x}}\mathbf{h}] = E\left[\sum_{i=1}^n \tilde{x}_i \frac{\partial}{\partial x_i} \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x}\right] \quad (28)$$

$$= \sum_{i=1}^n E[\tilde{x}_i] \frac{\partial}{\partial x_i} \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \quad (29)$$

$$= 0 \quad (30)$$

In the above expression we distributed the expected value to only the first term since $\frac{\partial}{\partial x_i} \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x}$ is a constant when evaluated at the known mean value and we used the fact that $E[\tilde{x}] = 0$. Now, consider the fourth (cubic) term in (27). Again, use the definition of $D_{\tilde{x}}\mathbf{h}$ and take the expected value:

$$E[D_{\tilde{x}}^3\mathbf{h}] = E\left[\sum_{i=1}^n \left[\tilde{x}_i \frac{\partial}{\partial x_i}\right]^3 \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x}\right] \quad (31)$$

$$= \sum_{i=1}^n E[\tilde{x}_i^3] \left[\frac{\partial}{\partial x_i}\right]^3 \mathbf{h}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \quad (32)$$

$$= 0 \quad (33)$$

This follows from the fact that all odd moments are zero (e.g., $E[\tilde{x}_i^3] = 0$ for all $i = 1, \dots, n$) for a pdf that is zero-mean and symmetrically distributed (i.e., $f_x(-x) = f_x(x)$). For a proof of this fact see the Appendix and [1, p.58-59]. Thus, the Taylor series (27) only contains even terms:

$$\boldsymbol{\mu}_y \triangleq E[\mathbf{y}] = E[\mathbf{h}(\boldsymbol{\mu}_x)] + E[\frac{1}{2!}D_{\tilde{x}}^2\mathbf{h}] + E[\frac{1}{4!}D_{\tilde{x}}^4\mathbf{h}] + E[\dots] \quad (34)$$

We have shown that if we simply propagate the mean of \mathbf{x} (i.e., $\boldsymbol{\mu}_x$) through the nonlinear function to give an approximate mean $\boldsymbol{\mu}_y = \mathbf{h}(\boldsymbol{\mu}_x)$ then we will be accurate only to first order. To improve accuracy we can evaluate the second order term in the Taylor series, and so on.

First-Order Nonlinear Transformation of a Gaussian Random Vector

We have just shown above that the Taylor series expansion of the expected value of a random vector contains only even terms (34). Here we go back to the nonlinear transformation we started with (19) but only consider the first two terms (i.e., a linear approximation):

$$\begin{aligned} \mathbf{y} &= \mathbf{h}(\boldsymbol{\mu}_x) + D_{\tilde{x}}\mathbf{h} \\ &= \mathbf{h}(\boldsymbol{\mu}_x) + (\mathbf{J}_x\mathbf{h})\Big|_{\mathbf{x}=\boldsymbol{\mu}_x} \tilde{\mathbf{x}} \end{aligned}$$

where $\tilde{x} = x - \mu_x$. This expression is in the form (5) with $A = (J_x h)|_x$ and $b = h(\mu_x)$. Thus, for the first-order linear approximation we can use our previous results for transformations of linear variables (8) and (17). That is, to first-order, the mean and covariance of the transformed random vector are:

$$\mu_y = h(\mu_x) \quad (35)$$

$$P_y = \left[(J_x h)|_{x=\mu_x} \right] P_x \left[(J_x h)|_{x=\mu_x} \right]^T \quad (36)$$

Example: First-order Error Propagation (1D case). Suppose we have a random variable $X \sim \mathcal{N}(\mu_x = 9, \sigma_x^2 = 0.2)$ that is transformed through the nonlinear function $z = f(x) = \sqrt{x}$ into a new random variable $Z \sim \mathcal{N}(\mu_z, \sigma_z^2)$. We wish to determine μ_z and σ_z^2 .

Solution. The mean μ_z is found by propagating the mean μ_x through the transformation.

$$\mu_z = f(\mu_x) = \sqrt{9} = 3$$

The first-order error propagation law (36) reduces in the 1D case to:

$$\sigma_z^2 = \left[J_f|_{x=\mu_x} \right]^2 \sigma_x^2$$

where J_f is the derivative of $f(x)$ with respect to x . Thus

$$J = \frac{d}{dx} \sqrt{x} = \frac{1}{2\sqrt{x}}$$

Evaluating with the mean μ_x

$$J_f|_{x=\mu_x} = \frac{1}{2\sqrt{\mu_x}} = \frac{1}{2\sqrt{9}} = \frac{1}{6}$$

Finally, applying the 1D error-propagation law:

$$\sigma_z^2 = \left[J_f|_{x=\mu_x} \right]^2 \sigma_x^2 = \left[\frac{1}{6} \right]^2 (0.2) = 0.00556$$

Thus, the random variable Y is characterized as $Z \sim \mathcal{N}(3, 0.00556)$.

The following example illustrates the same principle for a two-dimensional Gaussian random vector.

Example: First-order Error Propagation (2D case).

Given the random vector $x \sim \mathcal{N}(\mu_x, P_x)$ where

$$\mu_x = \begin{bmatrix} 16 \\ 30 \end{bmatrix} \quad P_x = \begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix}$$

find the mean μ_z and covariance P_z which characterizes the random variable $z \sim \mathcal{N}(\mu_z, P_z)$ (to first-order) that results from the transformation

$$z = f(x)$$

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix} = \begin{bmatrix} x_1^{1/4} + x_2 \\ x_1 x_2 / 10 \end{bmatrix}$$

Solution. The mean μ_z is obtained by propagating the mean μ_x through the system:

$$\mu_z = f(\mu_x) = \begin{bmatrix} 16^{1/4} + 30 \\ 30(16)/10 \end{bmatrix} = \begin{bmatrix} 32 \\ 48 \end{bmatrix}$$

Next, compute the Jacobian

$$J_f = \begin{bmatrix} \partial f_1 / \partial x_1 & \partial f_1 / \partial x_2 \\ \partial f_2 / \partial x_1 & \partial f_2 / \partial x_2 \end{bmatrix} = \begin{bmatrix} \left[\frac{1}{4}\right] x_1^{-3/4} & 1 \\ x_2 / 10 & x_1 / 10 \end{bmatrix}$$

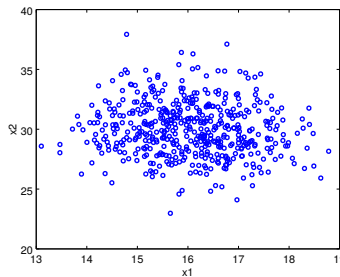
and evaluate the Jacobian at the mean μ_x

$$(J_x f)|_{x=\mu_x} = \begin{bmatrix} \left[\frac{1}{4}\right] (16)^{-3/4} & 1 \\ 30/10 & 16/10 \end{bmatrix} = \begin{bmatrix} 0.03125 & 1 \\ 3 & 1.6 \end{bmatrix}$$

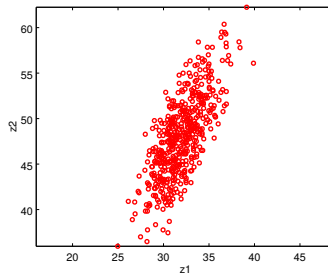
We have all the necessary elements to compute the covariance P_z using the error-propagation law:

$$\begin{aligned} P_z &= \left[(J_x f)|_{x=\mu_x} \right] P_x \left[(J_x f)|_{x=\mu_x} \right]^T \\ &= \begin{bmatrix} 0.03125 & 1 \\ 3 & 1.6 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 0.03125 & 1 \\ 3 & 1.6 \end{bmatrix}^T \\ &= \begin{bmatrix} 5.001 & 8.094 \\ 8.094 & 21.800 \end{bmatrix} \end{aligned}$$

The solution can be visualized numerically by generating samples from $x \sim \mathcal{N}(\mu_x, P_x)$



and propagating them through the system



Joint Gaussian Vector Distributions

Suppose that $\mathbf{x} \sim (\boldsymbol{\mu}_x, \mathbf{Q})$ is a Gaussian random vector taking values in \mathbb{R}^n and $\mathbf{y} \sim (\boldsymbol{\mu}_y, \mathbf{R})$ is a Gaussian random vector taking values in \mathbb{R}^m . Then the stacked random vector $\mathbf{z} = [\mathbf{x}^T, \mathbf{y}^T]^T$ jointly describes both random vectors (i.e., $p(\mathbf{x}, \mathbf{y})$). The mean is, of course,

$$E[\mathbf{z}] = \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \quad (37)$$

and the covariance is

$$\mathbf{P}_z = E \left[\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \right) \left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \right)^T \right] \quad (38)$$

$$= E \left[\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \right) \left(\begin{bmatrix} \mathbf{x}^T & \mathbf{y}^T \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_x^T & \boldsymbol{\mu}_y^T \end{bmatrix} \right) \right] \quad (39)$$

$$= E \left[\left(\begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix} \right) \left(\begin{bmatrix} (\mathbf{x} - \boldsymbol{\mu}_x)^T & (\mathbf{y} - \boldsymbol{\mu}_y)^T \end{bmatrix} \right) \right] \quad (40)$$

$$= E \left[\left(\begin{bmatrix} (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T & (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^T \\ (\mathbf{y} - \boldsymbol{\mu}_y)(\mathbf{x} - \boldsymbol{\mu}_x)^T & (\mathbf{y} - \boldsymbol{\mu}_y)(\mathbf{y} - \boldsymbol{\mu}_y)^T \end{bmatrix} \right) \right] \quad (41)$$

$$= \begin{bmatrix} \text{Cov}(\mathbf{x}, \mathbf{x}) & \text{Cov}(\mathbf{x}, \mathbf{y}) \\ \text{Cov}(\mathbf{y}, \mathbf{x}) & \text{Cov}(\mathbf{y}, \mathbf{y}) \end{bmatrix} \quad (42)$$

$$= \begin{bmatrix} \mathbf{Q} & \text{Cov}(\mathbf{x}, \mathbf{y}) \\ \text{Cov}(\mathbf{y}, \mathbf{x}) & \mathbf{R} \end{bmatrix} \quad (43)$$

The terms $\text{Cov}(\mathbf{x}, \mathbf{y})$ and $\text{Cov}(\mathbf{y}, \mathbf{x})$ describe the covariance of one random vector with the other. If they are entirely independent then these are zero matrices.

Conditional Gaussian Vector Distributions

Having described the joint distribution $p(\mathbf{x}, \mathbf{y})$ above we now consider the conditional distributions (e.g., $p(\mathbf{x}|\mathbf{y})$ which is the p.d.f. of \mathbf{x} given a particular value of \mathbf{y}). Consider again the case above of a stacked random vector $\mathbf{z} = [\mathbf{x}^T, \mathbf{y}^T]^T$ with the mean partitioned as

$$E[\mathbf{z}] = \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}. \quad (44)$$

Denote the partitioned covariance matrix (45) as

$$\mathbf{P}_z = \begin{bmatrix} \mathbf{P}_x & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_y \end{bmatrix} \quad (45)$$

The symmetry of \mathbf{P}_z implies that $\mathbf{P}_{xy} = \mathbf{P}_{yx}^T$. In many situations, it is convenient to work with the inverse of the covariance matrix, called the *precision matrix*, which we denote as:

$$\boldsymbol{\Lambda}_z := \mathbf{P}_z^{-1} = \begin{bmatrix} \boldsymbol{\Lambda}_{xx} & \boldsymbol{\Lambda}_{xy} \\ \boldsymbol{\Lambda}_{yx} & \boldsymbol{\Lambda}_{yy} \end{bmatrix} \quad (46)$$

Note however, that $\boldsymbol{\Lambda}_{xx}$ is a subblock of the precision matrix but $\boldsymbol{\Lambda}_{xx} \neq \mathbf{P}_{xx}^{-1}$. To obtain the mean and covariance as a function of the original covariance matrix (45) we can use the matrix identity associated with Schur's complement.

Aside (Schur's complement): Consider a matrix \mathbf{P} that can be partitioned into four sub-blocks as

$$\mathbf{P} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}. \quad (47)$$

The inverse of this matrix is given by the following identity:

$$\mathbf{P}^{-1} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{B} & (\mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1}) \end{bmatrix} \quad (48)$$

where the new matrix is introduced

$$\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} \quad (49)$$

Applying the above inverse identity to (45) and (46) it follows that

$$\Lambda_{xx} = (\mathbf{P}_{xx} - \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}\mathbf{P}_{yx})^{-1} \quad (50)$$

$$\Lambda_{xy} = -(\mathbf{P}_{xx} - \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}\mathbf{P}_{yx})^{-1}\mathbf{P}_{xy}\mathbf{P}_{yy}^{-1} \quad (51)$$

The quantity $\Lambda_{xx}^{-1}\Lambda_{xy}$ that appears in the mean then simplifies as

$$\Lambda_{xx}^{-1}\Lambda_{xy} = -(\mathbf{P}_{xx} - \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}\mathbf{P}_{yx})[(\mathbf{P}_{xx} - \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}\mathbf{P}_{yx})^{-1}\mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}] \quad (52)$$

$$= -\mathbf{P}_{xy}\mathbf{P}_{yy}^{-1} \quad (53)$$

so an equivalent expression for the mean and covariance of the condition probability is

$$\implies \mu_{x|y} = \mu_x + \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}(\mathbf{y} - \mu_y) \quad (54)$$

$$\implies \mathbf{P}_{x|y} = \mathbf{P}_{xx} - \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}\mathbf{P}_{yx} \quad (55)$$

Marginal Gaussian Vector Distributions

For the partitioned random vector $\mathbf{z} = [\mathbf{x}^T, \mathbf{y}^T]^T$ we have thusfar defined the the joint distribution $p(\mathbf{x}, \mathbf{y})$ as well as the condition distributions $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y}|\mathbf{x})$. The marginal distributions for each component are defined by the integrals

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \quad (56)$$

$$p(\mathbf{y}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{x} \quad (57)$$

where we “integrate out” the undesired variables. It can be shown (see [2, Sec 2.3.2]) that the marginal distribution is also Gaussian, and is intuitively given by

$$\mathbf{x} \sim \mathcal{N}(\mu_x, \mathbf{P}_{xx}) \quad (58)$$

$$\mathbf{y} \sim \mathcal{N}(\mu_y, \mathbf{P}_{yy}) \quad (59)$$

Product of Two Gaussians

Consider two independent distributions $\mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\mathcal{N}(\mathbf{b}, \mathbf{B})$ defined for random vectors \mathbf{x} and \mathbf{y} , respectively, of the same size. For a given \mathbf{x} or \mathbf{y} each of the p.d.f.s evaluate to a scalar value and the product of these two density values yields a new density that is also Gaussian. This new distribution is $Z^{-1}\mathcal{N}(\mathbf{c}, \mathbf{C})$ where

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \quad (60)$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}) \quad (61)$$

and the normalizing constant is

$$Z^{-1} = (2\pi)^{-n/2} |\mathbf{A} + \mathbf{B}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^T (\mathbf{A} + \mathbf{B})^{-1} (\mathbf{a} - \mathbf{b})\right) \quad (62)$$

Appendix: Linearization Review

Consider the vector-valued nonlinear function $\mathbf{f}(\cdot)$ with a scalar argument x . We can approximate the function $\mathbf{f}(x)$ around some nominal value $x = \bar{x}$ by defining a Taylor series. Defining the deviation from the nominal point as $\tilde{x} = x - \bar{x}$, the Taylor series expansion of $\mathbf{f}(x)$ is:

$$\begin{aligned}\mathbf{f}(x) &= \mathbf{f}(\bar{x}) + \left. \frac{d\mathbf{f}}{dx} \right|_{\bar{x}} \tilde{x} + \frac{1}{2!} \left. \frac{d^2\mathbf{f}}{dx^2} \right|_{\bar{x}} \tilde{x}^2 + \frac{1}{3!} \left. \frac{d^3\mathbf{f}}{dx^3} \right|_{\bar{x}} \tilde{x}^3 + \dots \\ &= \sum_{k=0}^{\infty} \frac{1}{k!} \left. \frac{d^k\mathbf{f}}{dx^k} \right|_{x=\bar{x}} (\tilde{x})^k\end{aligned}$$

Now suppose the function argument \mathbf{x} is a vector. For example, if \mathbf{x} is a 2x1 vector, then the Taylor series expansion is:

$$\begin{aligned}\mathbf{f}(\mathbf{x}) &= \mathbf{f}(\bar{\mathbf{x}}) + \left(\left. \frac{\partial \mathbf{f}}{\partial x_1} \right|_{\bar{\mathbf{x}}} \tilde{x}_1 + \left. \frac{\partial \mathbf{f}}{\partial x_2} \right|_{\bar{\mathbf{x}}} \tilde{x}_2 \right) + \frac{1}{2!} \left(\left. \frac{\partial^2 \mathbf{f}}{\partial x_1^2} \right|_{\bar{\mathbf{x}}} \tilde{x}_1^2 + \left. \frac{\partial^2 \mathbf{f}}{\partial x_2^2} \right|_{\bar{\mathbf{x}}} \tilde{x}_2^2 + 2 \left. \frac{\partial^2 \mathbf{f}}{\partial x_1 \partial x_2} \right|_{\bar{\mathbf{x}}} \tilde{x}_1 \tilde{x}_2 \right) \\ &\quad + \frac{1}{3!} \left(\left. \frac{\partial^3 \mathbf{f}}{\partial x_1^3} \right|_{\bar{\mathbf{x}}} \tilde{x}_1^3 + \left. \frac{\partial^3 \mathbf{f}}{\partial x_2^3} \right|_{\bar{\mathbf{x}}} \tilde{x}_2^3 + 3 \left. \frac{\partial^3 \mathbf{f}}{\partial x_1^2 \partial x_2} \right|_{\bar{\mathbf{x}}} \tilde{x}_1^2 \tilde{x}_2 + 3 \left. \frac{\partial^3 \mathbf{f}}{\partial x_1 \partial x_2^2} \right|_{\bar{\mathbf{x}}} \tilde{x}_1 \tilde{x}_2^2 \right) + \dots \\ &= \mathbf{f}(\bar{\mathbf{x}}) + \left(\tilde{x}_1 \frac{\partial}{\partial x_1} + \tilde{x}_2 \frac{\partial}{\partial x_2} \right) \mathbf{f} \Big|_{\bar{\mathbf{x}}} + \frac{1}{2!} \left(\tilde{x}_1^2 \frac{\partial^2}{\partial x_1^2} + \tilde{x}_2^2 \frac{\partial^2}{\partial x_2^2} + 2 \tilde{x}_1 \tilde{x}_2 \frac{\partial^2}{\partial x_1 \partial x_2} \right) \mathbf{f} \Big|_{\bar{\mathbf{x}}} \\ &\quad + \frac{1}{3!} \left(\tilde{x}_1^3 \frac{\partial^3}{\partial x_1^3} + \tilde{x}_2^3 \frac{\partial^3}{\partial x_2^3} + 3 \tilde{x}_1^2 \tilde{x}_2 \frac{\partial^3}{\partial x_1^2 \partial x_2} + 3 \tilde{x}_1 \tilde{x}_2^2 \frac{\partial^3}{\partial x_1 \partial x_2^2} \right) \mathbf{f} \Big|_{\bar{\mathbf{x}}} + \dots\end{aligned}$$

Define the operator $D_{\bar{\mathbf{x}}}^k \mathbf{f}$

$$\begin{aligned}D_{\bar{\mathbf{x}}}^k \mathbf{f} &= \left(\tilde{x}_1 \frac{\partial}{\partial x_1} + \tilde{x}_2 \frac{\partial}{\partial x_2} + \dots + \tilde{x}_n \frac{\partial}{\partial x_n} \right)^k \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\bar{\mathbf{x}}} \\ &= \left(\sum_{i=1}^n \tilde{x}_i \frac{\partial}{\partial x_i} \right)^k \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\bar{\mathbf{x}}}\end{aligned}$$

Then using this definition we may write the Taylor expansion for any sized vector argument \mathbf{x} as

$$\begin{aligned}\mathbf{f}(\mathbf{x}) &= \mathbf{f}(\bar{\mathbf{x}}) + D_{\bar{\mathbf{x}}} \mathbf{f} + \frac{1}{2!} D_{\bar{\mathbf{x}}}^2 \mathbf{f} + \frac{1}{3!} D_{\bar{\mathbf{x}}}^3 \mathbf{f} + \dots \\ &= \sum_{k=0}^{\infty} \frac{1}{k!} D_{\bar{\mathbf{x}}}^k \mathbf{f}\end{aligned}\tag{63}$$

Appendix: Odd moments of a symmetric probability distribution are zero

Suppose there is a random variable X with mean of zero $E[X] = 0$, and with a symmetric probability density function (pdf): $p(x) = p(-x)$. The i -th moment of X can be written

$$\begin{aligned}m_i &= E[X^i] \\ &= \int_{-\infty}^{\infty} x^i p(x) dx \\ &= \int_{-\infty}^0 x^i p(x) dx + \int_0^{\infty} x^i p(x) dx\end{aligned}$$

Changing the limits of the first term:

$$\int_{-\infty}^0 x^i p(x) dx = - \int_0^{-\infty} x^i p(x) dx$$

Then substituting $u = -x$, the right hand side becomes:

$$\begin{aligned} - \int_0^{-\infty} x^i p(x) dx &= \int_0^{\infty} (-u)^i p(-u) du \\ &= (-1)^i \int_0^{\infty} u^i p(u) du \end{aligned}$$

Thus

$$m_i = (-1)^i \underbrace{\int_0^{\infty} u^i p(u) du}_A + \underbrace{\int_0^{\infty} x^i p(x) dx}_A$$

and it is clear that the moment is zero for odd i .

Appendix: Conditional Gaussian Distribution

Recall that the multivariate p.d.f. is given by

$$f_z(z) = \frac{\exp\left(-\frac{1}{2}(z - \mu_z)^T P_z^{-1}(z - \mu_z)\right)}{(2\pi)^{n/2} \sqrt{|P_z|}} \quad (64)$$

First we establish the fact that, for a general Gaussian random vector $z \sim \mathcal{N}(\mu_z, P_z)$ the argument of the exponent in (64) has components that are quadratic in z , linear in z , and independent of z :

$$-\frac{1}{2}\Delta = -\frac{1}{2}(z - \mu_z)^T P_z^{-1}(z - \mu_z) \quad (65)$$

$$= -\frac{1}{2}(z^T P_z^{-1} z - z^T P_z^{-1} \mu_z - \mu_z^T P_z^{-1} z + \mu_z^T P_z^{-1} \mu_z) \quad (66)$$

$$= \underbrace{-\frac{1}{2}z^T P_z^{-1} z}_{\text{quadratic in } z} + \underbrace{z^T P_z^{-1} \mu_z}_{\text{linear in } z} - \frac{1}{2} \underbrace{\mu_z^T P_z^{-1} \mu_z}_{\text{independent of } z} \quad (67)$$

Recall that Δ is the Mahalanobis distance between z and μ . Our strategy to find the conditional probability is to expand this expression with the components of the stacked vector z while treating y as a fixed quantity and equating terms quadratic, linear, and independent of x to (67).

Expand the Mahalanobis distance first:

$$-\frac{1}{2}\Delta = (\mathbf{z} - \boldsymbol{\mu}_z)^T \mathbf{P}_z^{-1} (\mathbf{z} - \boldsymbol{\mu}_z) \quad (68)$$

$$= -\frac{1}{2} \left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \right)^T \begin{bmatrix} \boldsymbol{\Lambda}_{xx} & \boldsymbol{\Lambda}_{xy} \\ \boldsymbol{\Lambda}_{yx} & \boldsymbol{\Lambda}_{yy} \end{bmatrix} \left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \right) \quad (69)$$

$$= -\frac{1}{2} \begin{bmatrix} (\mathbf{x} - \boldsymbol{\mu}_x)^T & (\mathbf{y} - \boldsymbol{\mu}_y)^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\Lambda}_{xx} & \boldsymbol{\Lambda}_{xy} \\ \boldsymbol{\Lambda}_{yx} & \boldsymbol{\Lambda}_{yy} \end{bmatrix} \begin{bmatrix} (\mathbf{x} - \boldsymbol{\mu}_x) \\ (\mathbf{y} - \boldsymbol{\mu}_y) \end{bmatrix} \quad (70)$$

$$= -\frac{1}{2} \begin{bmatrix} (\mathbf{x} - \boldsymbol{\mu}_x)^T & (\mathbf{y} - \boldsymbol{\mu}_y)^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\Lambda}_{xx}(\mathbf{x} - \boldsymbol{\mu}_x) + \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y) \\ \boldsymbol{\Lambda}_{yx}(\mathbf{x} - \boldsymbol{\mu}_x) + \boldsymbol{\Lambda}_{yy}(\mathbf{y} - \boldsymbol{\mu}_y) \end{bmatrix} \quad (71)$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_x)^T (\boldsymbol{\Lambda}_{xx}(\mathbf{x} - \boldsymbol{\mu}_x) + \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y)) - \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_y)^T (\boldsymbol{\Lambda}_{yx}(\mathbf{x} - \boldsymbol{\mu}_x) + \boldsymbol{\Lambda}_{yy}(\mathbf{y} - \boldsymbol{\mu}_y)) \quad (72)$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\Lambda}_{xx}(\mathbf{x} - \boldsymbol{\mu}_x) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y) - \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_y)^T \boldsymbol{\Lambda}_{yx}(\mathbf{x} - \boldsymbol{\mu}_x) \quad (73)$$

$$- \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_y)^T \boldsymbol{\Lambda}_{yy}(\mathbf{y} - \boldsymbol{\mu}_y) \quad (74)$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\Lambda}_{xx}(\mathbf{x} - \boldsymbol{\mu}_x) - (\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y) - \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_y)^T \boldsymbol{\Lambda}_{yy}(\mathbf{y} - \boldsymbol{\mu}_y) \quad (75)$$

$$= -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Lambda}_{xx} \mathbf{x} + \mathbf{x}^T \boldsymbol{\Lambda}_{xx} \boldsymbol{\mu}_x - \frac{1}{2}\boldsymbol{\mu}_x^T \boldsymbol{\Lambda}_{xx} \boldsymbol{\mu}_x - \mathbf{x}^T \boldsymbol{\Lambda}_{xy} \mathbf{y} + \mathbf{x}^T \boldsymbol{\Lambda}_{xy} \boldsymbol{\mu}_y + \boldsymbol{\mu}_x^T \boldsymbol{\Lambda}_{xy} \mathbf{y} - \boldsymbol{\mu}_x^T \boldsymbol{\Lambda}_{xy} \boldsymbol{\mu}_y \quad (76)$$

$$- \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_y)^T \boldsymbol{\Lambda}_{yy}(\mathbf{y} - \boldsymbol{\mu}_y) \quad (77)$$

$$= \underbrace{-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Lambda}_{xx} \mathbf{x} + \mathbf{x}^T \boldsymbol{\Lambda}_{xx} \boldsymbol{\mu}_x}_{\text{quadratic in } \mathbf{x}} + \underbrace{-\mathbf{x}^T \boldsymbol{\Lambda}_{xy} \mathbf{y} + \mathbf{x}^T \boldsymbol{\Lambda}_{xy} \boldsymbol{\mu}_y}_{\text{linear in } \mathbf{x}} \quad (78)$$

$$- \underbrace{\frac{1}{2}\boldsymbol{\mu}_x^T \boldsymbol{\Lambda}_{xx} \boldsymbol{\mu}_x + \boldsymbol{\mu}_x^T \boldsymbol{\Lambda}_{xy} \mathbf{y} - \boldsymbol{\mu}_x^T \boldsymbol{\Lambda}_{xy} \boldsymbol{\mu}_y - \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_y)^T \boldsymbol{\Lambda}_{yy}(\mathbf{y} - \boldsymbol{\mu}_y)}_{\text{independent of } \mathbf{x}} \quad (79)$$

Comparing the quadratic term of (67) with (79) it is evident that the covariance of the conditional distribution $p(\mathbf{x}|\mathbf{y})$ is given by

$$\implies \mathbf{P}_{x|y} = \boldsymbol{\Lambda}_{xx}^{-1}. \quad (80)$$

Also, the linear term in (79) can be written as

$$\begin{aligned} \mathbf{x}^T \boldsymbol{\Lambda}_{xx} \boldsymbol{\mu}_x - \mathbf{x}^T \boldsymbol{\Lambda}_{xy} \mathbf{y} + \mathbf{x}^T \boldsymbol{\Lambda}_{xy} \boldsymbol{\mu}_y &= \mathbf{x}^T (\boldsymbol{\Lambda}_{xx} \boldsymbol{\mu}_x - \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y)) \\ &= \mathbf{x}^T \boldsymbol{\Lambda}_{xx} (\boldsymbol{\mu}_x - \boldsymbol{\Lambda}_{xx}^{-1} \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y)) \end{aligned}$$

and comparing with (79) it follows that

$$\implies \boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x - \boldsymbol{\Lambda}_{xx}^{-1} \boldsymbol{\Lambda}_{xy}(\mathbf{y} - \boldsymbol{\mu}_y). \quad (81)$$

Thus, we can conclude that the condition probability distribution is characterized by

$$\mathbf{x}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_{x|y}, \mathbf{P}_{x|y}) \quad (82)$$

References

- [1] Dan Simon. *Optimal State Estimation: Kalman, H infinity, and Nonlinear Approaches*. John Wiley & Sons, 2006.
- [2] Christopher M Bishop and Nasser M Nasrabadi. *Pattern Recognition and Machine Learning*. Springer, 2006.