

Multivariate Gaussian

The univariate Gaussian distribution can be written as

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

where μ is the mean and σ^2 is the variance

For a D -dim vector \vec{x} the multivariate Gaussian distribution ~~has the form~~ the following 3 defs are equiv.

Definition 1)

$$f_X(\vec{x}) = N(\vec{x}|\vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1}(\vec{x}-\vec{\mu})\right\}$$

which is equivalent to both

Definition 2) the moment generating function of X is

$$M_X(\vec{t}) \equiv \mathbb{E}[e^{\vec{t}'X}] = \exp\left\{\vec{\mu}'\vec{t} + \frac{1}{2}\vec{t}'\Sigma\vec{t}\right\}$$

Definition 3) X has the same distribution as

$A\vec{z} + \vec{\mu}$
where $\vec{z} = [z_1, \dots, z_k]$ is a sample from $N(0, I)$
and $A_{n \times k}$ satisfied $AA' = \Sigma$

Here $\vec{\mu}$ is a vector, Σ is a PSD matrix
and for \vec{x} all it can be written that

$$\vec{x} \sim N(\vec{x}|\vec{\mu}, \Sigma)$$

Thm: Def's 1, 2, and 3 are equiv. for Σ pos. def. . Def's 2 & 3 are equivalent for Σ pos. semi. def.

Proof that Def 3 \Rightarrow Def 2

For $z_i \sim N(0, 1)$

$$\begin{aligned} M_{z_i}(t_i) &= \mathbb{E}[e^{t_i z_i}] = \int_{-\infty}^{\infty} e^{z_i t_i} \frac{e^{-z_i^2/2}}{\sqrt{2\pi}} dz_i \\ &= \int_{-\infty}^{\infty} e^{z_i t_i} \frac{e^{-(z_i - t_i)^2/2}}{\sqrt{2\pi}} e^{-z_i t_i} e^{t_i^2/2} dz_i \\ &= e^{t_i^2/2} \int_{-\infty}^{\infty} \frac{e^{-(z_i - t_i)^2/2}}{\sqrt{2\pi}} dz_i \\ &= e^{t_i^2/2} \frac{\sqrt{2\pi}}{\sqrt{2\pi}} = e^{t_i^2/2} \end{aligned}$$

If $\vec{z} = [z_1, \dots, z_k]$ is a random sample from $N(0, 1)$ then -

$$\begin{aligned} M_{\vec{z}}(\vec{t}) &= \mathbb{E}[e^{\vec{z}'\vec{t}}] = \mathbb{E}[e^{\sum_i z_i t_i}] \\ &= \mathbb{E}\left[\prod_{i=1}^k e^{z_i t_i}\right] \quad \text{because } z_i \sim \text{iid } N(0, 1) \\ &= \prod_{i=1}^k \mathbb{E}[e^{z_i t_i}] \\ &= \prod_{i=1}^k M_{z_i}(t_i) \\ &= \exp\left\{\sum_{i=1}^k t_i^2/2\right\} \\ &= \exp\left\{\vec{t}'\vec{t}/2\right\} \end{aligned}$$

Which gives us the M.G.F. for $\vec{z} \sim N(0, I)$

If $\vec{X} = A\vec{Z} + \vec{\mu}$ then -

$$\begin{aligned}M_{\vec{X}}(\vec{t}) &= E[\exp\{\vec{X}'\vec{t}\}] \\&= E[\exp\{(A\vec{Z} + \vec{\mu})'\vec{t}\}] \\&= \exp(\vec{\mu}'\vec{t}) \cdot E[\exp\{(\vec{A}\vec{Z})'\vec{t}\}] \\&= \exp(\vec{\mu}'\vec{t}) \cdot E[\exp\{\vec{Z}'A'\vec{t}\}] \\&= \exp(\vec{\mu}'\vec{t}) \cdot M_{\vec{Z}}(A'\vec{t}) \\&= \exp(\vec{\mu}'\vec{t}) \cdot \exp\left\{\frac{1}{2} (A'\vec{t})'(A'\vec{t})\right\} \\&= \exp(\vec{\mu}'\vec{t}) \cdot \exp\left\{\frac{1}{2} \vec{t}' A A' \vec{t}\right\} \\&= \exp\left(\vec{\mu}'\vec{t} + \frac{1}{2} \vec{t}' \Sigma \vec{t}\right)\end{aligned}$$

- Note, if A is invertible then $A'(\vec{X} - \vec{\mu}) \sim N(\vec{0}, I)$ which can be useful for hypothesis testing, data pre-conditioning, etc.
- If two ~~dist~~ distributions have the same MGF then they are identical at all points.

Proof that Def 2 \Rightarrow Def 3

Since $\Sigma \geq 0$ (and $\Sigma = \Sigma'$) there is an orthogonal matrix $C^{n \times n}$ such that $C' \Sigma C = \Lambda$ where Λ is diagonal with non-negative elements. So

$$\begin{aligned}\Sigma &= C \Lambda C' \\ &= C \Lambda^{1/2} \Lambda^{1/2} C' \\ &= (C \Lambda^{1/2}) (C \Lambda^{1/2})' \\ &= A A'\end{aligned}$$

But we know that the moment generating function of $A\tilde{Z} + \vec{\mu}$ is (by the previous proof)

$$\exp\{\vec{\mu}' \vec{t} + \frac{1}{2} \vec{t}' \Sigma \vec{t}\}$$

i.e. the same as \tilde{X} . Because the m.g.f. uniquely determines the distribution, \tilde{X} has the same distribution as $A\tilde{Z} + \vec{\mu}$.

•

Proof of Def 3 \Rightarrow Def 1: (for pos. det. Σ)

Because Σ is pos. def. there is a non-singular $A_{n \times n}$ s.t. $AA' = \Sigma$. Let ~~\vec{X}~~ $\vec{X} = A\vec{Z} + \mu$ where $\vec{Z} \sim N(0, I)$. The density of \vec{Z} is

$$f_{\vec{Z}}(\vec{z}) = \prod_{i=1}^D (2\pi)^{-1/2} \exp\left\{-\frac{1}{2} z_i^2\right\} = (2\pi)^{-D/2} \exp\left\{-\frac{1}{2} \vec{z}' \vec{z}\right\}$$

By the change of variable rule (transformation of distribution) the density function of \vec{X} is

$$f_{\vec{X}}(\vec{x}) = f_{\vec{Z}}(\vec{z}(\vec{x})) |J| \quad \text{determinant of the jacobian}$$

$$\vec{z}(\vec{x}) = A^{-1}(\vec{x} - \mu), \quad \frac{\partial \vec{z}(\vec{x})}{\partial \vec{x}} = \frac{\partial A^{-1} \vec{x}}{\partial \vec{x}} = A^{-1}$$

$$\Rightarrow J = A^{-1}$$

$$\Rightarrow f_{\vec{X}}(\vec{x}) = f_{\vec{Z}}(\vec{z}(\vec{x})) |A^{-1}| = f_{\vec{Z}}(\vec{z}(\vec{x})) |A|^{-1}$$

Therefore

$$f_{\vec{X}}(\vec{x}) = f_{\vec{Z}}(A^{-1}(\vec{x} - \mu)) |A|^{-1}$$

$$\begin{aligned} \text{Note } |A|^{-1} &= (2\pi)^{-D/2} |A|^{-1} \exp\left\{-\frac{1}{2} [A^{-1}(\vec{x} - \mu)]' [A^{-1}(\vec{x} - \mu)]\right\} \\ &= |A|^{-1/2} |A|^{-1/2} \\ &= |A|^{-1/2} |A'|^{-1/2} \\ &= (|A| |A'|)^{-1/2} \\ &= |AA'|^{-1/2} \\ &= (2\pi)^{-D/2} |AA'|^{-1/2} \exp\left\{-\frac{1}{2} (\vec{x} - \mu)' (A^{-1})' A^{-1} (\vec{x} - \mu)\right\} \\ &= (2\pi)^{-D/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2} (\vec{x} - \mu)' (AA')^{-1} (\vec{x} - \mu)\right\} \\ &= (2\pi)^{-D/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2} (\vec{x} - \mu)' \Sigma^{-1} (\vec{x} - \mu)\right\} \end{aligned}$$

Pf that Def 1 \Rightarrow Def 2 left to reader.

Properties of MVN distributions

1) If $\vec{X} \sim N(\vec{\mu}, \Sigma)$

$$E[\vec{X}] = \vec{\mu}, \quad \text{cov}[\vec{X}] = \Sigma$$

Pf. We know \vec{X} has the same distribution as $A\vec{z} + \vec{\mu}$ so

$$E[\vec{X}] = E[A\vec{z} + \vec{\mu}] = E[A\vec{z}] + \vec{\mu} = \vec{\mu}$$

$$\text{Cov}[\vec{X}] = E[(\vec{X} - E[\vec{X}])(\vec{X} - E[\vec{X}])']$$

$$= E[(\vec{X} - \vec{\mu})(\vec{X} - \vec{\mu})']$$

$$= E[(A\vec{z})(A\vec{z})']$$

$$= E[A\vec{z}\vec{z}'A']$$

$$= A \text{Cov}(\vec{z}) A'$$

$$= A I A' = A A' = \Sigma$$

2) If ~~$\vec{z} \sim N(0, I_{D \times D})$~~ ~~$\vec{z} = \{z_1, \dots, z_D\}$~~
 ~~$\vec{z} = (z_1, \dots, z_D)$~~ , $z_i \sim \text{iid } N(0, \sigma^2)$ then

$$\vec{z} \sim N_D(\vec{0}_D, I_{D \times D} \sigma^2) \quad (\text{matrix normal})$$

$$p(\vec{z}) = \prod_{i=1}^D \frac{1}{\sqrt{2\pi}} \sigma^{-1} \exp \left\{ -\frac{1}{2\sigma^2} (z_i)^2 \right\}$$

$$= (2\pi)^{-D/2} \sigma^{-D} \exp \left\{ -\frac{1}{2} \sum_{i=1}^D z_i^2 / \sigma^2 \right\}$$

$$= (2\pi)^{-D/2} |I \sigma^2|^{-1/2} \exp \left\{ -\frac{1}{2} (\vec{z}' (I \sigma^2) \vec{z}) \right\}$$

Linear transformations of MVN vectors

1. If $\vec{y} \sim N_n(\vec{\mu}, \Sigma)$ and $C_{p \times n}$ is a constant matrix of rank p then $C\vec{y} \sim N_p(C\vec{\mu}, C\Sigma C')$

Proof: By def 3, $\vec{y} = A\vec{z} + \vec{\mu}$ where $AA' = \Sigma$, so

$$C\vec{y} = C(A\vec{z} + \vec{\mu}) \\ = CA\vec{z} + C\vec{\mu}$$

but by def. 3

$$CA\vec{z} + C\vec{\mu} \sim N(C\vec{\mu}, (CA)(CA)')$$

$$= N(C\vec{\mu}, CAA'C')$$

$$= N(C\vec{\mu}, C\Sigma C')$$

2. \vec{y} is MVN iff $\vec{a}'\vec{y}$ is normally distributed for all non-zero constant vectors \vec{a}

Pf. If $\vec{y} \sim N_n(\vec{\mu}, \Sigma)$ then by the previous claim

$$\vec{a}'\vec{y} \sim N(\vec{a}'\vec{\mu}, \vec{a}'\Sigma\vec{a})$$

Conversely, assume $x = \vec{a}'\vec{y}$ is univariate normal for all nonzero \vec{a} . Then since $E[x] = \vec{a}'\vec{\mu}$ and $\text{cov}(x) = \sigma^2\{x\} = \text{var}\{x\} = \text{cov}(\vec{a}'\vec{y}) = \vec{a}'\Sigma\vec{a}$.

The univariate normal M.G.F. is

$$E[\exp(xt)] = M_x(t) = \exp\left\{(\vec{a}'\vec{\mu})t + \frac{1}{2}(\vec{a}'\Sigma\vec{a})t^2\right\}$$

for all t . Setting $t=1$ we have

$$E[\exp(\vec{a}'\vec{y})] = \exp\left\{(\vec{a}'\vec{\mu}) + \frac{1}{2}(\vec{a}'\Sigma\vec{a})\right\} = M(\vec{a})$$

which is the M.G.F. of $N_n(\vec{\mu}, \Sigma)$ and therefore $\vec{y} \sim N_n(\vec{\mu}, \Sigma)$.

Conditional Gaussian Distributions

Important features: if two sets of vars are jointly Gaussian then 1) the cond. dist. of one set given the other is Gaussian and 2) the marginal dist of each set is also Gaussian.

Suppose $x \in \mathbb{R}^D$ and $x \sim \mathcal{N}(x | \mu, \Sigma)$.
let x_a be (w.l.o.g.) the first M components of x and x_b the remaining $D-M$ components.
i.e.

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$

let
$$\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}$$

and let
$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \leftarrow \begin{array}{l} \text{note } \Sigma_{aa}, \Sigma_{bb} \\ \text{symmetric} \\ \text{and } \Sigma_{ab} = \Sigma_{ba}^T \end{array}$$

be a covariance matrix

The precision matrix $\Lambda \equiv \Sigma^{-1}$ is sometimes easier to work with.

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

since Λ is also symmetric the same stuff holds

To start let's find $p(x_a | x_b)$

By the product rule we know that we can get $p(x_a | x_b)$ simply by fixing x_b in the joint $p(x_a, x_b)$ and renormalizing by inspection, i.e.

$$\begin{aligned} p(x_a | x_b) &\propto p(x_a, x_b) \\ &= \frac{p(x_a, x_b)}{\int p(x_a, x_b) dx_a} = \frac{p(x_a, x_b)}{p(x_b)} \end{aligned}$$

But we will proceed by inspection

$$\begin{aligned} \ln p(x_a | x_b, \mu, \Sigma) &\propto -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) + \text{const} \\ &= -\frac{1}{2} (x_a - \mu_a)^T \Lambda_{aa} (x_a - \mu_a) - \frac{1}{2} (x_a - \mu_a)^T \Lambda_{ab} (x_b - \mu_b) \\ &\quad - \frac{1}{2} (x_b - \mu_b)^T \Lambda_{ba} (x_a - \mu_a) - \frac{1}{2} (x_b - \mu_b)^T \Lambda_{bb} (x_b - \mu_b) + \text{const} \end{aligned}$$

but by inspection we can see that, as a function of ~~μ_a~~ x_a this is again a quadratic form

Reminder: Completing the square

$$-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) = -\frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu + \text{const}$$

If we collect the terms quadratic in x_a we can immediately identify the ~~mean~~ $\text{cov} \Sigma$ of $x_a | x_b$ is,

$$\begin{aligned} -\frac{1}{2} x^T \Sigma^{-1} x &= -\frac{1}{2} x_a^T \Lambda_{aa} x_a \\ \Rightarrow \Sigma^{-1} &= \Lambda_{aa} \end{aligned}$$

to identify the mean we have to collect the linear terms in x_a

$$\begin{aligned} &x_a^T \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b) \} \\ \Rightarrow x_a^T \Lambda_{aa} \mu &= x_a^T \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b) \} \\ \Rightarrow \mu &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b) \end{aligned}$$

yielding

$$x_a | x_b \sim N(x_a | \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b), \Lambda_{aa}^{-1})$$

If we wish to have everything in terms of the covariance matrix instead of the precision matrix, we need the following matrix inversion identity ~~fact~~.

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{pmatrix}$$

where $M = (A - BD^{-1}C)^{-1}$ (called the Schur complement)

$$\text{Since } \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

and using the inversion identity we have

$$\begin{aligned} \Lambda_{aa} &= (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1} \\ \Lambda_{ab} &= -(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1} \Sigma_{ab} \Sigma_{bb}^{-1} \end{aligned}$$

which gives us (after some algebra)

$$\mu_{ab} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$$

$$\Sigma_{ab} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

Note that μ_{ab} is a linear function of x_b and the covariance is independent of x_a . This is an example of a linear gaussian model.

Marginal Gaussian Distributions

We know now that if $p(x_a, x_b)$ is Gaussian then $p(x_a | x_b)$ is also Gaussian, what about $p(x_a)$?

$$p(x_a) = \int p(x_a, x_b) dx_b$$

Using the same tricks (completing squares, recognizing quadratic forms, etc.) it can be arrived at that

$$\begin{aligned} \mathbb{E}[x_a] &= \mu_a \\ \text{cov}[x_a] &= \Sigma_{aa} \end{aligned}$$

i.e. the marginal distribution of x_a arises from simply ignoring the parts of the distribution corresponding to x_b .

In Summary

Given a joint Gaussian dist $\mathcal{N}(\vec{x} | \vec{\mu}, \Sigma)$ with $\Lambda \equiv \Sigma^{-1}$ and

$$\vec{x} = \begin{pmatrix} \vec{x}_a \\ \vec{x}_b \end{pmatrix}, \quad \vec{\mu} = \begin{pmatrix} \vec{\mu}_a \\ \vec{\mu}_b \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

Condition

$$p(\vec{x}_a | \vec{x}_b) = \mathcal{N}(\vec{x}_a | \vec{\mu}_{a|b}, \Lambda_{aa}^{-1})$$

$$\vec{\mu}_{a|b} = \vec{\mu}_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\vec{x}_b - \vec{\mu}_b)$$

Marginal

$$p(\vec{x}_a) = \mathcal{N}(\vec{x}_a | \vec{\mu}_a, \Sigma_{aa})$$

Bayes Thm for Gaussian Variables

Take

$$p(\vec{x}) = \mathcal{N}(\vec{x} | \vec{\mu}, \Lambda^{-1})$$

and

$$p(\vec{y} | \vec{x}) = \mathcal{N}(\vec{y} | A\vec{x} + \vec{b}, L^{-1})$$

Λ, L - precision matrices

μ, A, b - params governing means

We want $p(\vec{y})$ and $p(\vec{x} | \vec{y})$

To do this we define a R.V.

$$\vec{z} = \begin{pmatrix} \vec{x} \\ \vec{y} \end{pmatrix} \leftarrow \text{stacking vectors}$$

and construct the joint distribution of this R.V.
by considering the log of the joint

$$\begin{aligned} \ln p(\vec{z}) &= \ln p(\vec{x}) + \ln p(\vec{y} | \vec{x}) \\ &= -\frac{1}{2} (\vec{x} - \vec{\mu})^T \Lambda (\vec{x} - \vec{\mu}) \\ &\quad - \frac{1}{2} (\vec{y} - A\vec{x} - \vec{b})^T L (\vec{y} - A\vec{x} - \vec{b}) + \text{const} \end{aligned}$$

↑
terms indep
of \vec{x} & \vec{y}

In this we complete the square by
first collecting quadratic terms in \vec{x} & \vec{y}

$$-\frac{1}{2} \vec{x}^T (\Lambda + A^T L A) \vec{x} - \frac{1}{2} \vec{y}^T L \vec{y} + \frac{1}{2} \vec{x}^T A^T L \vec{y} + \frac{1}{2} \vec{y}^T L A \vec{x}$$

and note that this equals

$$-\frac{1}{2} \begin{pmatrix} \vec{x} \\ \vec{y} \end{pmatrix}^T \begin{pmatrix} \Lambda + A^T L A & -A^T L \\ -L A & L \end{pmatrix} \begin{pmatrix} \vec{x} \\ \vec{y} \end{pmatrix} \equiv -\frac{1}{2} \vec{z}^T R \vec{z}$$

where

$$R = \begin{pmatrix} \Lambda + A^T L A & -A^T L \\ -L A & L \end{pmatrix}$$

which remembering $-\frac{1}{2}(\mathbf{z} - \mu_z)^T \text{cov}[\mathbf{z}]^{-1}(\mathbf{z} - \mu_z)$

$$= -\frac{1}{2} \mathbf{z}^T \text{cov}[\mathbf{z}]^{-1} \mathbf{z} + \mathbf{z}^T \text{cov}[\mathbf{z}]^{-1} \mu_z$$

$$\Rightarrow \text{cov}[\mathbf{z}]^{-1} = R, \quad \text{cov}[\mathbf{z}] = R^{-1}$$

$$\Rightarrow \text{cov}[\mathbf{z}] = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A \Lambda^{-1} & L^{-1} + A \Lambda^{-1} A^T \end{pmatrix}$$

verify? HW

from partitioned inverse eqn.

Now that we have the covariance matrix for \mathbf{z} we can find the mean of the distribution (completing square). For this we need the linear terms from the log-post

$$\mathbf{x}^T \Lambda \mu + \mathbf{y}^T L \mathbf{b} + \mathbf{x}^T A^T L \mathbf{b} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \Lambda \mu + A^T L \mathbf{b} \\ L \mathbf{b} \end{pmatrix}$$

$$\Rightarrow \mathbb{E}[\mathbf{z}] = \mu_z = R^{-1} \begin{pmatrix} \Lambda \mu + A^T L \mathbf{b} \\ L \mathbf{b} \end{pmatrix}$$

which after some algebra yields

$$\mathbf{z} \sim \mathcal{N}\left(\begin{pmatrix} \mu \\ A\mu + \mathbf{b} \end{pmatrix}, \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A \Lambda^{-1} & L^{-1} + A \Lambda^{-1} A^T \end{pmatrix}\right)$$

but, since we have analytic forms for the conditionals and marginals of Gaussian distributions we essentially have our answers.

By inspection

$$E[y] = A\mu + b$$

$$\text{cov}[y] = L^{-1} + A\mathcal{L}^{-1}A^T$$

And now to get $p(\vec{x}|\vec{y})$ we use conditional

And, remembering

$$p(x_a|x_b) = \mathcal{N}(x_a | \mu_{a|b}, \mathcal{L}_{aa}^{-1})$$

$$\mu_{a|b} = \mu_a - \mathcal{L}_{aa}^{-1}\mathcal{L}_{ab}(x_b - \mu_b)$$

~~$$\Rightarrow p(\vec{y}|\vec{x}) = \mathcal{N}(\vec{y} | \dots)$$~~

$$\Rightarrow p(\vec{x}|\vec{y}) = \mathcal{N}(\vec{x} | \mu_{x|y}, \mathcal{L}_{xx}^{-1})$$
$$\mu_{x|y} = \mu_x - \mathcal{L}_{xx}^{-1}\mathcal{L}_{xy}(y - \mu_y)$$

$$\text{where } \mathcal{L}_{xx}^{-1} = (\mathcal{L} + A^T L A)^{-1}$$

$$\mathcal{L}_{xy} = A^T L$$

$$\mu_y = A\mu + b$$

$$\mu_x = \mu$$

$$\mu_{x|y} = \mu + (\mathcal{L} + A^T L A)^{-1} A^T L (y - A\mu - b)$$

$$\Rightarrow p(\vec{x}|\vec{y}) = \mathcal{N}(\vec{x} | E[x|y], \text{cov}[x|y])$$

where

$$E[x|y] = (\mathcal{L} + A^T L A)^{-1} (A^T L (y - b) + \mathcal{L}\mu)$$

$$\text{cov}[x|y] = (\mathcal{L} + A^T L A)^{-1}$$

Summarizing it

$$p(x) = \mathcal{N}(x | \mu, \mathcal{L}^{-1}) \quad \text{and}$$

$$p(y|x) = \mathcal{N}(y | Ax + b, L^{-1})$$

then

$$p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\mathcal{L}^{-1}A^T)$$

$$p(x|y) = \mathcal{N}(x | \Sigma \{A^T L (y - b) + \mathcal{L}\mu\}, \Sigma)$$

where

$$\Sigma = (\mathcal{L} + A^T L A)^{-1}$$

Maximum likelihood for the Gaussian

Given a dataset $X = (x_1, \dots, x_n)^T$ where $x_i \sim \mathcal{N}(\mu, \Sigma)$ by assumption we can estimate μ and Σ from the data using maximum likelihood

$$\ln p(X|\mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

to find an ML estimator for μ we take the partial w.r.t. μ

$$\begin{aligned} \frac{\partial}{\partial \mu} \ln p(X|\mu, \Sigma) &= \cancel{\frac{\partial}{\partial \mu}} - \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial \mu} \left[(x_n - \mu)^T \Sigma^{-1} (x_n - \mu) \right] \\ &= -\frac{1}{2} \sum_{n=1}^N \cancel{2} \Sigma^{-1} (x_n - \mu) \end{aligned}$$

now we set this equal to zero and solve

$$\sum_{n=1}^N (x_n - \mu) = 0 \Rightarrow \boxed{\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n}$$

To solve for the ML estimator for Σ we need the following ~~two~~ facts

$$\boxed{\frac{\partial \ln|X|}{\partial X} = (X^{-1})^T = (X^T)^{-1}}$$

$$\text{and } \boxed{\frac{\partial (a^T X b)}{\partial X} = b a^T}$$

$$\text{and } \boxed{\frac{\partial a^T X^{-1} b}{\partial X} = -X^{-T} a b^T X^{-T}}$$

$$\frac{\partial}{\partial \Sigma} \ln P(X|\mu, \Sigma) = -\frac{N}{2} \frac{\partial}{\partial \Sigma} \ln |\Sigma|$$

$$-\frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial \Sigma} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

$$= -\frac{N}{2} (\Sigma^{-1})^T + \frac{1}{2} \sum_{n=1}^N \Sigma (x_n - \mu)(x_n - \mu)^T \Sigma$$

Now we set this expression equal to zero and solve

$$\frac{N}{2} \Sigma^{-1} = \frac{1}{2} \sum_{n=1}^N \Sigma^{-1} (x_n - \mu)(x_n - \mu)^T \Sigma^{-1}$$

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T$$

But we don't know μ . We do know that μ_{ML} optimizes the same.