

Lecture 2: The Multivariate Normal Distribution

STATS305C: Applied Statistics III

Scott Linderman

March 30, 2022

Last Time...

- ▶ Course Overview
- ▶ Bayes' Rule
- ▶ Normal with unknown mean (normal prior)
- ▶ Normal with unknown precision (χ^2 prior)
- ▶ Normal with unknown variance (χ^{-2} prior)
- ▶ Normal with unknown mean and variance (NIX prior)
- ▶ Posterior marginals (Student's t distribution)

Today...

- ▶ Multivariate normal (MVN) distribution
 - ▶ Generative story
 - ▶ Marginal distributions
 - ▶ Conditional distributions
 - ▶ Linear Gaussian models
- ▶ The Wishart and inverse Wishart distributions
- ▶ Bayesian estimation with a normal-inverse-Wishart (NIW) prior
- ▶ Posterior marginals (multivariate Student's t distribution)

Reading: Bishop, Ch 2.3. See also: Murphy, Ch 2.3 and 3.2.4.

Generative Story

Start with a vector of standard normal random variates, $\mathbf{z} = [z_1, \dots, z_D]^\top$ where $z_d \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.

This is a D -dimensional random variable, but not a very interesting one. All the coordinates are independent! The joint density is,

$$p(\mathbf{z}) = \prod_{d=1}^D \mathcal{N}(z_d \mid 0, 1) \quad (1)$$

$$= \prod_{d=1}^D \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_d^2} \quad (2)$$

$$= (2\pi)^{-\frac{D}{2}} \exp \left\{ -\frac{1}{2} \sum_{d=1}^D z_d^2 \right\} \quad (3)$$

$$= (2\pi)^{-\frac{D}{2}} \exp \left\{ -\frac{1}{2} \mathbf{z}^\top \mathbf{z} \right\} \quad (4)$$

$$\triangleq \mathcal{N}(\mathbf{z} \mid \mathbf{0}, I). \quad (5)$$

Question: What do the contours of this joint density look like in $D = 2$ dimensions?

Spherical Gaussian Density

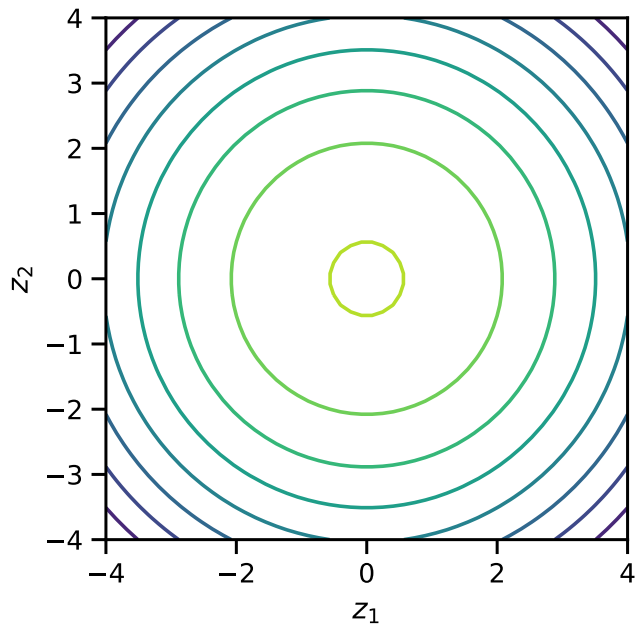


Figure: Contours of the pdf of a **spherical Gaussian** distribution, $\mathcal{N}(\mathbf{z} \mid \mathbf{0}, I)$.

Generative Story II

We can obtain more interesting joint distributions by transforming this random vector.

For example, let \mathbf{U} be an orthogonal $D \times D$ matrix and $\mathbf{\Lambda} = \text{diag}([\lambda_1, \dots, \lambda_D])$ with $\lambda_d > 0$. Define the linearly transformed random variable $\mathbf{x} = \mathbf{U}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{z}$. $\mathbf{\Lambda}^{1/2} = \text{diag}([\sqrt{\lambda_1}, \dots, \sqrt{\lambda_D}])$

Exercise: Compute the mean $\mathbb{E}[\mathbf{x}]$ and covariance $\text{Cov}[\mathbf{x}]$.

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{z}] \\ &= \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbb{E}[\mathbf{z}] \\ &= \mathbf{0}\end{aligned}$$

$$\begin{aligned}\text{Cov}[\mathbf{x}] &= \text{Cov}[\mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{z}] \\ &= \mathbf{U}\mathbf{\Lambda}^{1/2}\underbrace{\text{Cov}(\mathbf{z})}_{\mathbf{I}}\mathbf{\Lambda}^{1/2}\mathbf{U}^T \\ &= \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T\end{aligned}$$

Generative Story III

Exercise: Explain each step in this derivation of the density:

$$p(\mathbf{x}) = p(\mathbf{z}) \left| \frac{d\mathbf{z}}{d\mathbf{x}} \right| \quad \text{change of variables} \quad (6)$$

$$= p(\Lambda^{-\frac{1}{2}} \mathbf{U}^T \mathbf{x}) |\Lambda^{-\frac{1}{2}} \mathbf{U}^T| \quad \text{Substitute } \mathbf{x} = \mathbf{U} \Lambda^{\frac{1}{2}} \mathbf{z} \quad (7)$$

$$= (2\pi)^{-\frac{D}{2}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \mathbf{U} \Lambda^{-1} \mathbf{U}^T \mathbf{x} \right\} |\Lambda^{-\frac{1}{2}}| |\mathbf{U}^T| \quad \begin{aligned} \mathbf{z} &= \Lambda^{-\frac{1}{2}} \mathbf{U}^T \mathbf{x} \\ &= \Lambda^{-\frac{1}{2}} \mathbf{U}^T \mathbf{x} \end{aligned} \quad (8)$$

$$= (2\pi)^{-\frac{D}{2}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \mathbf{U} \Lambda^{-1} \mathbf{U}^T \mathbf{x} \right\} |\Lambda|^{-\frac{1}{2}} \quad \frac{d\mathbf{z}}{d\mathbf{x}} = \Lambda^{-\frac{1}{2}} \mathbf{U}^T \quad (9)$$

$$= (2\pi)^{-\frac{D}{2}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} \right\} |\Sigma|^{-\frac{1}{2}} \quad (10)$$

where $\Sigma = \mathbf{U} \Lambda \mathbf{U}^T$.

$$= \mathcal{N}(\mathbf{x} | \mu=0, \Sigma)$$

Generative Story IV

Last but not least, add a translation so that $\mathbf{x} = \mathbf{U}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$ for $\boldsymbol{\mu} \in \mathbb{R}^D$.

Question: How does this change the mean and covariance of \mathbf{x} ? $\mathbb{E}[\mathbf{x}] = \mathbb{E}[\mathbf{U}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}] = \boldsymbol{\mu}$

Following the same argument as above,

$$\text{cov}[\mathbf{x}] = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{\top} \triangleq \boldsymbol{\Sigma}$$

$$p(\mathbf{x}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \triangleq \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (11)$$

This is the pdf of the **multivariate normal** distribution with mean vector $\boldsymbol{\mu} \in \mathbb{R}^D$ and positive definite covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$. Again, $\boldsymbol{\Sigma} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{\top}$.

It depends on \mathbf{x} through the **squared Mahalanobis distance** between \mathbf{x} and $\boldsymbol{\mu}$,

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad (12)$$

Question: What do the contours of this density look like?

~~Spherical~~ Gaussian Density

Multivariate

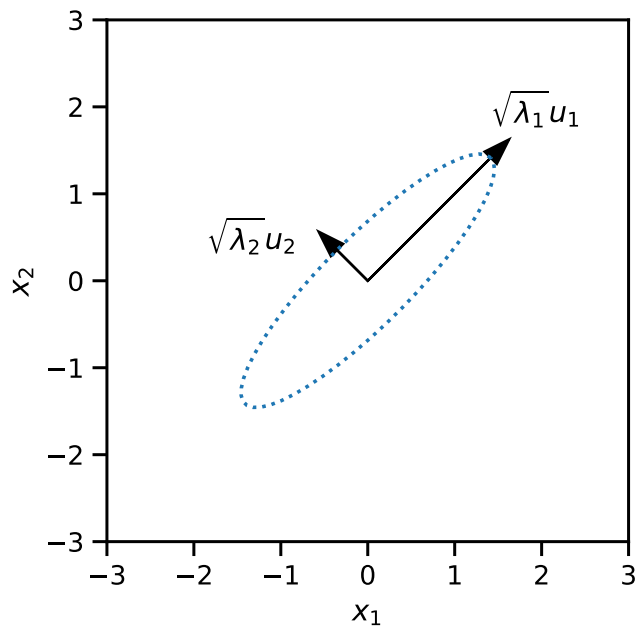


Figure: Contours of the pdf of a **multivariate normal** distribution, $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Generative Story V

We introduced the MVN by scaling \mathbf{z} , applying a change of basis, and adding a translation.

However, we could have applied *any* linear transformation,

$$\mathbf{x} = \mathbf{A}\mathbf{z} + \boldsymbol{\mu}. \quad (13)$$

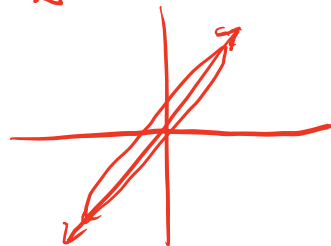
for $\mathbf{A} \in \mathbb{R}^{M \times D}$ and $\boldsymbol{\mu} \in \mathbb{R}^M$. Then,

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (14)$$

with $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^\top$.

Question: What if $M > D$?

→ degenerate mvn. dist.



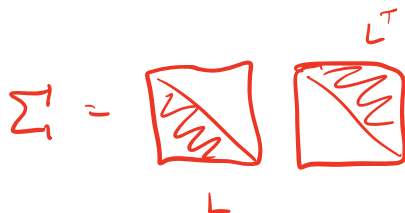
Generative Story VI

Finally, we can go the other direction as well: to sample $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, first compute **any square root** of the covariance matrix $\boldsymbol{\Sigma}^{\frac{1}{2}} \in \mathbb{R}^{D \times D}$ such that $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}^{\frac{1}{2}})(\boldsymbol{\Sigma}^{\frac{1}{2}})^\top$, then set,

$$\mathbf{x} = \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{z} + \boldsymbol{\mu}, \quad (15)$$

where $\mathbf{z} \in \mathbb{R}^D$ is a vector of standard normal random variates.

As before, we could use the eigendecomposition $\boldsymbol{\Sigma}^{\frac{1}{2}} = \mathbf{U}\boldsymbol{\Lambda}^{\frac{1}{2}}$ where $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$. Since the covariance is positive semidefinite, we could also use the **Cholesky decomposition**, $\boldsymbol{\Sigma}^{\frac{1}{2}} = \text{chol}(\boldsymbol{\Sigma})$, in which case the square root is lower triangular.


$$\boldsymbol{\Sigma} = \mathbf{L} \mathbf{L}^\top$$

Linear Transformations of Gaussians

The multivariate normal distribution is closed under linear transformations,

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \Rightarrow \quad \mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top). \quad (16)$$

Marginal Distributions

Let $\mathbf{x} \in \mathbb{R}^D$ be a multivariate normal random vector. Partition the vector and the MVN parameters into two subsets,

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}. \quad (17)$$

The symmetry of $\boldsymbol{\Sigma}$ implies that $\boldsymbol{\Sigma}_{aa}$ and $\boldsymbol{\Sigma}_{bb}$ are symmetric, whereas $\boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^\top$.

Now consider the linear transformation $\mathbf{A} = \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix}$ so that $\mathbf{A}\mathbf{x} = \mathbf{x}_a$.

The linearity property implies that the **marginal distributions** of the multivariate normal are also multivariate normal,

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top) = \mathcal{N}(\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}). \quad (18)$$

Simply extract the corresponding blocks of the mean and covariance matrix to get the marginal.

Conditional Distributions

Let $\mathbf{x} \in \mathbb{R}^D$ be a multivariate normal random vector. Partition the vector and the MVN parameters into two subsets,

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}. \quad (19)$$

The symmetry of $\boldsymbol{\Sigma}$ implies that $\boldsymbol{\Sigma}_{aa}$ and $\boldsymbol{\Sigma}_{bb}$ are symmetric, whereas $\boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^\top$.

We can write the inverse of $\boldsymbol{\Sigma}$ in block form as well,

$$\boldsymbol{\Sigma}^{-1} \triangleq \boldsymbol{\Lambda} = \begin{bmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{bmatrix}. \quad \text{(note: not } \boldsymbol{\Lambda} \text{ from before.)} \quad (20)$$

It too is symmetric, and its blocks involve the **Schur complement**,

$$\boldsymbol{\Lambda}_{aa} = (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba})^{-1} \quad (21)$$

$$\boldsymbol{\Lambda}_{ab} = -(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba})^{-1} \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1}. \quad (22)$$

Conditional Distributions II

What is $p(\mathbf{x}_a | \mathbf{x}_b)$? By Bayes' rule,

$$p(\mathbf{x}_a | \mathbf{x}_b) \propto p(\mathbf{x}_a, \mathbf{x}_b) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (23)$$

$$\propto \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (24)$$

$$\propto \exp \left\{ -\frac{1}{2} (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \boldsymbol{\Lambda}_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) - (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \right\} \quad (25)$$

$$\propto \exp \left\{ -\frac{1}{2} \mathbf{x}_a^\top \mathbf{J}_{a|b} \mathbf{x}_a + \mathbf{x}_a^\top \mathbf{h}_{a|b} \right\} \quad (26)$$

where $\mathbf{J}_{a|b} = \boldsymbol{\Lambda}_{aa}$ and $\mathbf{h}_{a|b} = \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)$.

This is the **information form** of a Gaussian density on \mathbf{x}_a .

\mathbf{J}, \mathbf{h} are "natural" params

Conditional Distributions III

Completing the square (recall Lecture 1), we find that

$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}), \quad (27)$$

where

$$\boldsymbol{\Sigma}_{a|b} \preceq \boldsymbol{\Sigma}_{aa}$$

$$\boldsymbol{\Sigma}_{a|b} = \mathbf{J}_{a|b}^{-1} \quad (28)$$

$$= \boldsymbol{\Lambda}_{aa}^{-1} \quad (29)$$

$$= \boldsymbol{\Sigma}_{aa} - \underbrace{\boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}}_{\text{schur complement}} \quad (30)$$

and

$$\boldsymbol{\mu}_{a|b} = \mathbf{J}_{a|b}^{-1} \mathbf{h}_{a|b} \quad (31)$$

$$= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (32)$$

$$= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b). \quad (33)$$

Linear Gaussian Models

Now suppose $\mathbf{x} \sim \mathcal{N}(\mathbf{b}, \mathbf{Q})$ and $\mathbf{y} \mid \mathbf{x} \sim \mathcal{N}(\mathbf{C}\mathbf{x} + \mathbf{d}, \mathbf{R})$. What is the joint distribution $p(\mathbf{x}, \mathbf{y})$?

Reparameterize the joint model as,

$$\mathbf{x} = \mathbf{b} + \mathbf{Q}^{\frac{1}{2}} \mathbf{z}_x \quad (34)$$

$$\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{d} + \mathbf{R}^{\frac{1}{2}} \mathbf{z}_y \quad (35)$$

$$= \mathbf{C}\mathbf{b} + \mathbf{C}\mathbf{Q}^{\frac{1}{2}} \mathbf{z}_x + \mathbf{d} + \mathbf{R}^{\frac{1}{2}} \mathbf{z}_y \quad (36)$$

where \mathbf{z}_x and \mathbf{z}_y are independent standard normal vectors.

Now combine them into a single vector and rearrange terms,

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{C}\mathbf{b} + \mathbf{d} \end{bmatrix} + \begin{bmatrix} \mathbf{Q}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{C}\mathbf{Q}^{\frac{1}{2}} & \mathbf{R}^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \mathbf{z}_x \\ \mathbf{z}_y \end{bmatrix} \quad (37)$$

Linear Gaussian Models II

Thus, \mathbf{x} and \mathbf{y} are jointly Gaussian,

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (38)$$

where

$$\boldsymbol{\mu} = \begin{bmatrix} \mathbf{b} \\ \mathbf{C}\mathbf{b} + \mathbf{d} \end{bmatrix} \quad (39)$$

and

$$\boldsymbol{\Sigma} = \begin{bmatrix} \mathbf{Q}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{C}\mathbf{Q}^{\frac{1}{2}} & \mathbf{R}^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \mathbf{Q}^{\frac{1}{2}} & \mathbf{Q}^{\frac{1}{2}}\mathbf{C}^{\top} \\ \mathbf{0} & \mathbf{R}^{\frac{1}{2}} \end{bmatrix} = \begin{bmatrix} \mathbf{Q} & \mathbf{Q}\mathbf{C}^{\top} \\ \mathbf{C}\mathbf{Q} & \mathbf{C}\mathbf{Q}\mathbf{C}^{\top} + \mathbf{R} \end{bmatrix}. \quad (40)$$

Question: What is the marginal distribution of \mathbf{y} ?

$$\mathbf{y} \sim \mathcal{N}(\mathbf{C}\mathbf{b} + \mathbf{d}, \mathbf{C}\mathbf{Q}\mathbf{C}^{\top} + \mathbf{R})$$

Maximum Likelihood Estimation

Exercise: The log likelihood is,

$$x_n \overset{\text{iid}}{\sim} \mathcal{N}(\mu, \Sigma) \quad \mathcal{L}(\mu, \Sigma) = \sum_{n=1}^N \log p(\mathbf{x}_n \mid \mu, \Sigma). \quad (41)$$

Take gradients and set them to zero to obtain the the maximum likelihood estimates,

$$\mu_{\text{ML}}, \Sigma_{\text{ML}} = \arg \max \mathcal{L}(\mu, \Sigma). \quad (42)$$

Show that,

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (43)$$

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{\text{ML}})(\mathbf{x}_n - \mu_{\text{ML}})^{\top}. \quad (44)$$

Bayesian Estimation: Multivariate Normal with Unknown Mean

Like last time, let's start with a simple Bayesian model of Gaussian data with a known covariance but an unknown mean:

$$\boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (45)$$

$$\mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (46)$$

Goal: Infer the posterior distribution $p(\boldsymbol{\mu} \mid \mathbf{X}, \boldsymbol{\eta})$ where $\boldsymbol{\eta} = (\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0, \boldsymbol{\Sigma})$.

$$p(\boldsymbol{\mu} \mid \mathbf{X}, \boldsymbol{\eta}) \propto \mathcal{N}(\boldsymbol{\mu} \mid \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (47)$$

Bayes' rule

$$\propto \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) \right\} \prod_{n=1}^N \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \quad (48)$$

$$\propto \exp \left\{ -\frac{1}{2} \boldsymbol{\mu}^\top \mathbf{J}_N \boldsymbol{\mu} + \boldsymbol{\mu}^\top \mathbf{h}_N \right\} \quad (49)$$

where $\mathbf{J}_N = \boldsymbol{\Sigma}_0^{-1} + N\boldsymbol{\Sigma}^{-1}$ and $\mathbf{h}_N = \boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 + \sum_{n=1}^N \boldsymbol{\Sigma}^{-1}\mathbf{x}_n$.

Bayesian Estimation: Multivariate Normal with Unknown Mean II

Completing the square,

$$p(\boldsymbol{\mu} \mid \mathbf{X}, \boldsymbol{\eta}) = \mathcal{N}(\boldsymbol{\mu} \mid \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) \quad (50)$$

where

$$\boldsymbol{\Sigma}_N = \mathbf{J}_N^{-1} = (\boldsymbol{\Sigma}_0^{-1} + N\boldsymbol{\Sigma}^{-1})^{-1} \quad (51)$$

$$\boldsymbol{\mu}_N = \mathbf{J}_N^{-1} \mathbf{h}_N = \boldsymbol{\Sigma}_N \left(\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \sum_{n=1}^N \boldsymbol{\Sigma}^{-1} \mathbf{x}_n \right) \quad (52)$$

Question: What does the posterior converge to in the uninformative limit $\boldsymbol{\Sigma}_0 \rightarrow \infty$?

$$\boldsymbol{\Sigma}_N \rightarrow \frac{1}{N} \boldsymbol{\Sigma}$$

$$\boldsymbol{\mu}_N \rightarrow \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \boldsymbol{\mu}_{ML}$$

Multivariate Normal with Unknown *Precision*

Now imagine the mean is known but not the covariance. Like last time, we will actually start by parameterizing the model in terms of the **precision**, i.e. inverse covariance, $\Lambda = \Sigma^{-1}$.

In the univariate case, we used a scaled χ^2 distribution for the prior, which was defined as the average of squared standard normal variates. We will take an analogous approach here.

Let $\Lambda = \sum_{i=1}^{\nu_0} \mathbf{z}_i \mathbf{z}_i^\top$ where $\mathbf{z}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \Lambda_0)$ with covariance Λ_0 .

Then, Λ follows a **Wishart distribution*** with ν_0 degrees of freedom and scale Λ_0 , written,

$$\Lambda \sim W(\nu_0, \Lambda_0). \quad (53)$$

Its expected value is $\mathbb{E}[\Lambda] = \nu_0 \Lambda_0$.

* There is an unfortunate asymmetry between this definition and the definition of the scaled χ^2 distribution from the previous lecture. Whereas we defined the scaled χ^2 as the *average* of squared Gaussian random variables, here we've defined the Wishart to be the sum of "squared" (really the outer product of) multivariate normal random vectors. This is to be consistent with the textbook definitions of the Wishart distribution, even though I find it more convenient to parameterize the distribution in terms of its mean.

The Wishart Distribution

Let \mathcal{S}_D denote the set of $D \times D$ positive definite matrices, and let $\mathbf{\Lambda} \in \mathcal{S}_D$ be a random variable.

The Wishart distribution is a distribution on the set of positive definite matrices. Its pdf is,

$$W(\mathbf{\Lambda} \mid \nu_0, \mathbf{\Lambda}_0) = \frac{1}{2^{\frac{\nu_0 D}{2}} |\mathbf{\Lambda}_0|^{\frac{\nu_0}{2}} \Gamma_D\left(\frac{\nu_0}{2}\right)} |\mathbf{\Lambda}|^{\frac{\nu_0 - D - 1}{2}} e^{-\frac{1}{2} \text{Tr}(\mathbf{\Lambda}_0^{-1} \mathbf{\Lambda})} \quad (54)$$

where $\nu_0 > 0$ specifies the **degrees of freedom** and $\mathbf{\Lambda}_0 \in \mathcal{S}_D$ is the **scale**. Γ_D is the multivariate gamma function.

The mean of the Wishart distribution is $\nu_0 \mathbf{\Lambda}_0$ and the mode is $(\nu_0 - D - 1) \mathbf{\Lambda}_0$ for $\nu_0 \geq (D + 1)$.

In the univariate case where $D = 1$, we have $W(\lambda \mid \nu_0, \nu_0^{-1} \lambda_0) = \chi^2(\nu_0, \lambda_0)$.

The Wishart distribution plays a key role in **frequentist statistics** as the distribution of the sample covariance matrix of mean-zero multivariate normal r.v.'s. In **Bayesian statistics**, it arises as the conjugate prior for the precision of a multivariate normal distribution.

The Wishart Distribution

$$\Lambda \sim W(\nu_0 = 4, \Lambda_0 = \frac{I}{4})$$

$$\Sigma^{-1} = \Lambda$$

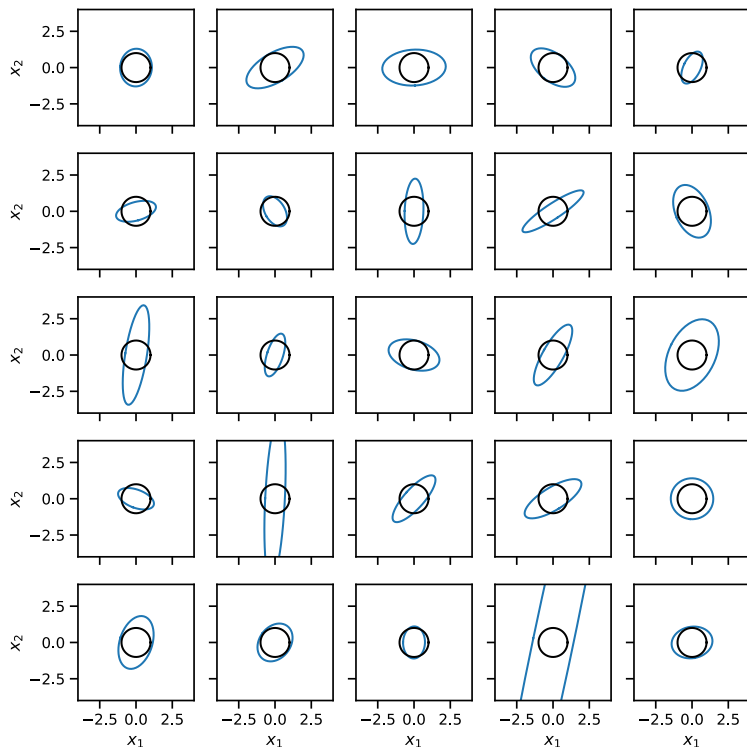


Figure: Visualizing Λ^{-1} where $\Lambda \sim W(\nu_0, \Lambda_0)$.

Multivariate Normal with Unknown Precision II

The Wishart distribution is a conjugate prior for the precision of a multivariate normal distribution.

$$\Lambda \sim W(\nu_0, \Lambda_0), \quad (55)$$

$$\mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \Lambda^{-1}). \quad (56)$$

Then, letting $\boldsymbol{\eta} = (\boldsymbol{\mu}, \nu_0, \Lambda_0)$,

$$p(\Lambda \mid \mathbf{X}, \boldsymbol{\eta}) \propto W(\Lambda \mid \nu_0, \Lambda_0) \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}, \Lambda^{-1})$$

"Trace trick"
→ $e^{-1/2 (\mathbf{x}_n - \boldsymbol{\mu})^\top \Lambda (\mathbf{x}_n - \boldsymbol{\mu})}$
= $e^{-1/2 \text{Tr}((\mathbf{x}_n - \boldsymbol{\mu})^\top \Lambda (\mathbf{x}_n - \boldsymbol{\mu}))}$
= $e^{-1/2 \text{Tr}(\underbrace{\Lambda (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top}_{\mathbb{R}^{D \times D}})}$

$$= e^{-1/2 \text{Tr}(\underbrace{\Lambda (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top}_{\mathbb{R}^{D \times D}})} \quad (57)$$

$$\propto |\Lambda|^{\frac{\nu_0 - D - 1}{2}} e^{-\frac{1}{2} \text{Tr}(\Lambda_0^{-1} \Lambda)} \prod_{n=1}^N |\Lambda|^{\frac{1}{2}} e^{-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^\top \Lambda (\mathbf{x}_n - \boldsymbol{\mu})} \quad (58)$$

$$\propto |\Lambda|^{\frac{\nu_0 + N - D - 1}{2}} \exp \left\{ -\frac{1}{2} \text{Tr} \left(\left[\Lambda_0^{-1} + \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top \right] \Lambda \right) \right\} \quad (59)$$

Multivariate Normal with Unknown Precision III

We recognize this as yet another Wishart distribution,

$$p(\Lambda \mid \mathbf{X}, \boldsymbol{\eta}) \propto W(\Lambda \mid \nu_N, \Lambda_N), \quad (60)$$

where

$$\nu_N = \nu_0 + N \quad (61)$$

$$\Lambda_N = \left[\Lambda_0^{-1} + \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top \right]^{-1} \quad (62)$$

$\Lambda_0 \rightarrow \infty$

Question: What is the posterior mean under the uninformative prior where $\nu_0 \rightarrow 0$ or in the large data limit where $N \rightarrow \infty$?

$$\mathbb{E}[\Lambda \mid \mathbf{X}, \boldsymbol{\eta}] = \nu_N^{-1} \Lambda_N \longrightarrow N \left[\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top \right]^{-1} = \Sigma_{ML}^{-1}$$

Visualizing the Posterior

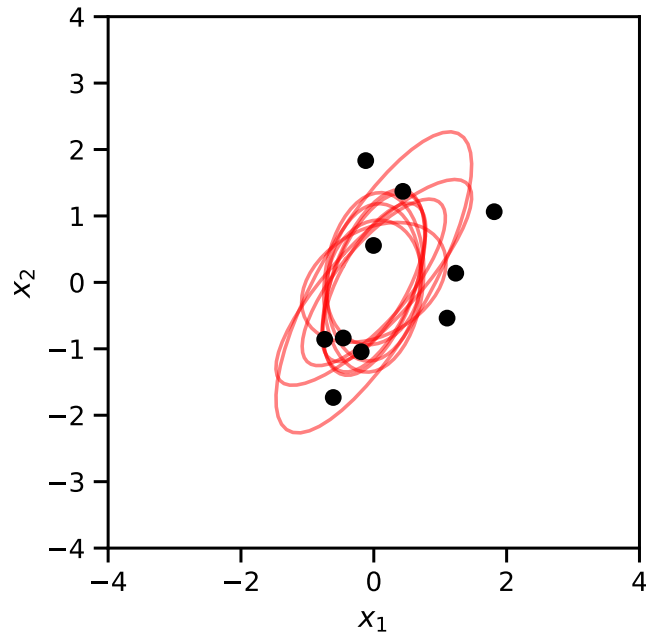


Figure: Visualizing Λ^{-1} under the posterior $\Lambda \sim W(\nu_N, \Lambda_N)$.

Multivariate Normal with Unknown *Covariance*

As with the scaled inverse chi-squared distribution in the univariate case, we define the **inverse Wishart** distribution as,

$$\Lambda \sim W(\nu_0, \Lambda_0) \iff \Sigma = \Lambda^{-1} \sim IW(\nu_0, \Sigma_0). \quad (63)$$

where $\Sigma_0 = \Lambda_0^{-1}$.

It too is a distribution on the set of positive definite matrices. Its pdf is,

$$IW(\Sigma \mid \nu_0, \Sigma_0) = \frac{|\Sigma_0|^{\frac{\nu_0}{2}}}{2^{\frac{\nu_0 D}{2}} \Gamma_D\left(\frac{\nu_0}{2}\right)} |\Sigma|^{-\frac{\nu_0 + D + 1}{2}} e^{-\frac{1}{2} \text{Tr}(\Sigma_0 \Sigma^{-1})} \quad (64)$$

where $\nu_0 > 0$ specifies the **degrees of freedom** and $\Sigma_0 \in \mathcal{S}_D$ is the **scale**.

Its mean is $\frac{\Sigma_0}{\nu - D - 1}$ for $\nu_0 \geq (D + 1)$ and its mode is $\frac{\Sigma_0}{\nu + D + 1}$.

In the univariate case where $D = 1$, we have $IW(\sigma^2 \mid \nu_0, \nu_0 \sigma_0^2) = \chi^2(\nu_0, \sigma_0^2)$.

Multivariate Normal with Unknown *Covariance*

Exercise: Consider the model,

$$\Sigma \sim \text{IW}(\nu_0, \Sigma_0), \quad (65)$$

$$\mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \Sigma). \quad (66)$$

Show that,

$$p(\Sigma \mid \mathbf{X}, \boldsymbol{\eta}) = \text{IW}(\nu_N, \Sigma_N) \quad (67)$$

where

$$\nu_N = \nu_0 + N \quad (68)$$

$$\Sigma_N = \Sigma_0 + \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top. \quad (69)$$

Multivariate Normal with Unknown Mean and Covariance

As you can probably guess, the normal-inverse-chi-squared distribution for the univariate case generalizes to the multivariate case too. Here, the conjugate prior for the mean and covariance is the **normal inverse Wishart** (NIW) distribution,

$$\text{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \boldsymbol{\mu}_0, \kappa_0, \nu_0, \boldsymbol{\Sigma}_0) = \text{IW}(\boldsymbol{\Sigma} \mid \nu_0, \boldsymbol{\Sigma}_0) \mathcal{N}(\boldsymbol{\mu} \mid \boldsymbol{\mu}_0, \boldsymbol{\Sigma}/\kappa_0). \quad (70)$$

Exercise: Show that under this prior the posterior is,

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{X}, \boldsymbol{\eta}) = \text{NIW}(\boldsymbol{\mu}_N, \kappa_N, \nu_N, \boldsymbol{\Sigma}_N) \quad (71)$$

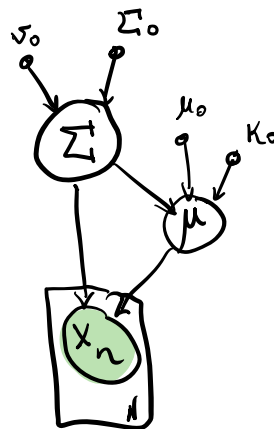
where

$$\nu_N = \nu_0 + N \quad (72)$$

$$\kappa_N = \kappa_0 + N \quad (73)$$

$$\boldsymbol{\mu}_N = \frac{1}{\kappa_N} \left(\kappa_0 \boldsymbol{\mu}_0 + \sum_{n=1}^N \mathbf{x}_n \right) \quad (74)$$

$$\boldsymbol{\Sigma}_N = \boldsymbol{\Sigma}_0 + \kappa_0 \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^\top + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top - \kappa_N \boldsymbol{\mu}_N \boldsymbol{\mu}_N^\top \quad (75)$$



Visualizing the Posterior

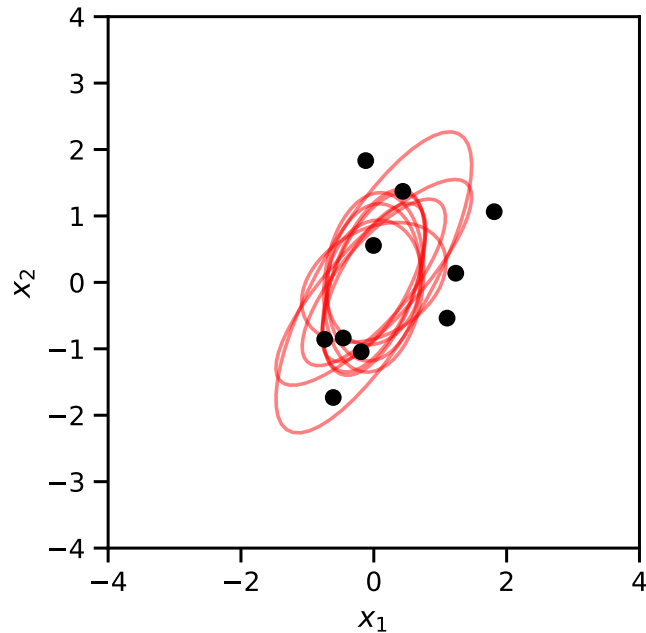
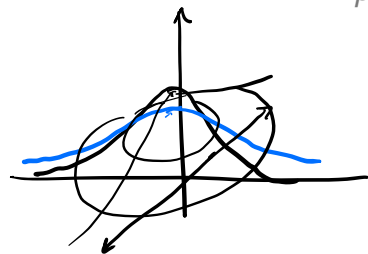


Figure: Visualizing μ and Λ^{-1} under the posterior $\mu, \Lambda \sim \text{NW}(\mu_N, \kappa_N, \nu_N, \Lambda_N)$.

Posterior Marginals

And again, just like in the univariate case, we find that the posterior marginal of $\boldsymbol{\mu} \in \mathbb{R}^D$ is

$$p(\boldsymbol{\mu} \mid \mathbf{X}, \boldsymbol{\eta}) = \int p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{X}, \boldsymbol{\eta}) d\boldsymbol{\Sigma} \quad (76)$$



$$= \int \mathcal{N}(\boldsymbol{\mu} \mid \boldsymbol{\mu}_N, \boldsymbol{\Sigma}/\kappa_N) \text{IW}(\boldsymbol{\Sigma} \mid \nu_N, \boldsymbol{\Sigma}_N) d\boldsymbol{\Sigma} \quad (77)$$

$$= \text{St} \left(\nu_N, \boldsymbol{\mu}_N, \frac{1}{\kappa_N(\nu_N - D + 1)} \boldsymbol{\Sigma}_N \right) \quad (78)$$

where St denotes the **multivariate Student's t** distribution with density,

$$\text{St}(\mathbf{x} \mid \nu, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\Gamma(\frac{\nu+D}{2})}{\Gamma(\frac{\nu}{2})} (\nu\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \left[1 + \frac{\Delta^2}{\nu} \right]^{-\frac{\nu+D}{2}} \quad (79)$$

where $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ is the squared Mahalanobis distance.

The mean of the multivariate Student's t is $\boldsymbol{\mu}$ and the covariance is $\frac{\nu}{\nu+2} \boldsymbol{\Sigma}$.

$$p(\Sigma|x,\gamma) = \int p(\Sigma, \mu|x, \gamma) d\mu$$

$$= \int NIW(\Sigma, \mu | \nu_N, \Sigma_N, \mu_N, \kappa_N) d\mu$$

$$= \int IW(\Sigma | \nu_N, \Sigma_N) \cdot N(\mu | \mu_N, \Sigma / \kappa_N) d\mu$$

$$= IW(\Sigma | \nu_N, \Sigma_N) \underbrace{\int N(\mu | \mu_N, \Sigma / \kappa_N) d\mu}_{=1}$$

$$= IW(\Sigma | \nu_N, \Sigma_N)$$