

# The Multivariate Normal and the Choleski and Eigen Decompositions

Rob McCulloch

1. Introduction
2. Change of Variable
3. Orthogonal Matrices and Rotation
4. Multivariate Normal
5. The Choleski Decomposition
6. More on the Multivariate Normal
7. Simulating a Multivariate Normal
8. Likelihood, Sufficiency, and MLE
9. Checking for Normality
10. Weighted Regression

# 1. Introduction

A square matrix  $A = [a_{ij}]$  is *symmetric* if  $a_{ij} = a_{ji}$ .

A square, symmetric matrix is *positive definite* (pd) if

$$x'Ax > 0 \quad \forall x.$$

Our basic example is a covariance matrix.

If  $X = (X_1, X_2, \dots, X_p)'$  is a random (column) vector with  $E(X) = \mu = (\mu_1, \mu_2, \dots, \mu_p)'$ , then the covariance of  $X$  is

$$\Sigma = E((X - \mu)(X - \mu)') = [E((X_i - \mu_i)(X_j - \mu_j))]$$

is symmetric.

Since

$$\text{Var}(a'X) = a'\Sigma a$$

$\Sigma$  is positive definiteness unless some linear combination of the  $X_i$  has 0 variance.



Let's review two basic matrix decompositions for symmetric pd matrices and use them to review basic properties of the multivariate normal distribution.

We'll look at:

(i):

The eigen decomposition.

(i):

The Choleski decomposition.

Later we will also look at the Singular Value Decomposition.

## 2. Change of Variable

To develop the normal distribution based on matrix decompositions, we will need the change of variable formulas, univariate and multivariate.

Let's review these.

Let  $\Theta$  be a random variable with density  $p(\theta)$ .

In Bayesian statistics,  $\theta$  is often used for the parameter of the model so that  $p(\theta)$  is the prior distribution.

The general Bayesian model consists of:

$$f(y | \theta), \quad p(\theta).$$

Rather than think in terms of the parameter  $\theta$  we may want to consider a 1-1 reparametrization

$$\gamma = g(\theta),$$

where  $g$  is 1 to 1.

What is  $p(\gamma)$  ??

## Univariate change of variable

$$\Theta \in \mathbb{R}$$

$$p(\Theta)$$

$$y = g(\Theta) \quad (1 \text{ to } 1)$$

$$\Theta = h(y) \quad (h = g^{-1})$$

$$p(y) = p(h(y)) |h'(y)|$$

or

$$\begin{aligned} p(y) &= p(\Theta(y)) |\Theta'(y)| \\ &= p(\Theta(y)) \left| \frac{d\Theta}{dy} \right| \end{aligned}$$

## Example

Suppose  $Y \sim \text{Bernoulli}(\theta)$  and  $p(\theta) = 1$ .

That is, we have the uniform prior on  $\theta \in (0, 1)$ .

Suppose we want to work with the odds-ratio, instead of the probability.

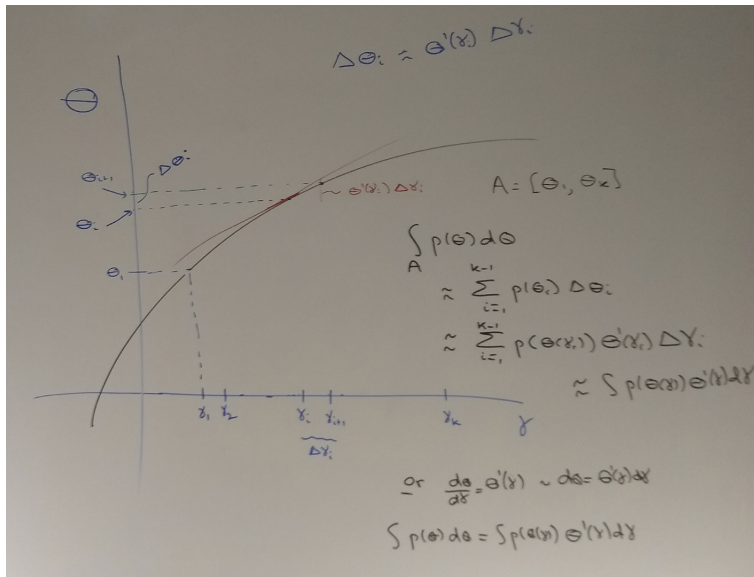
$$\gamma = \frac{1 - \theta}{\theta} \in \mathbb{R}^+$$

$$\theta = \frac{1}{1 + \gamma}$$

$$\frac{d\theta}{d\gamma} = -(1 + \gamma)^{-2}$$

$$p(\gamma) = \frac{1}{(1 + \gamma)^2}$$

# Calculus intuition for univariate change of variable



## Example, linear

Suppose  $X \sim p(x|\alpha)$ , where  $\alpha$  is a “shape” parameter.

Let  $Y = a + bX$ .

$$X = \frac{Y-a}{b}. \quad \frac{dx}{dy} = \frac{1}{b}.$$

$$f(y \mid a, b, \alpha) = p\left(\frac{y-a}{b}\right) \frac{1}{b}$$

## example:

Details:

If scale is omitted, it assumes the default value of 1.

The Gamma distribution with parameters `shape = a` and `scale = s` has density

$$f(x) = 1/(s^a \Gamma(a)) x^{a-1} e^{-(x/s)}$$

for  $x \geq 0$ ,  $a > 0$  and  $s > 0$ . (Here  $\Gamma(a)$  is the function implemented by R's `gamma()` and defined in its help. Note that  $a = 0$  corresponds to the trivial distribution with all mass at point 0.)

$$f(x|a) = \frac{1}{\Gamma(a)} x^{a-1} e^{-x}.$$

$$Y = sX, \quad X = Y/s, \quad dx/dy = 1/s.$$

$$f(y|s, a) = \frac{1}{\Gamma(a)} (y/s)^{a-1} e^{-y/s} (1/s)$$



## Multivariate Change of Variable

$$\Theta \in \mathbb{R}^k$$

$$x = g(\Theta) \quad \Theta = h(x)$$

$$h'(x) = \left[ \frac{\partial \Theta_i}{\partial x_j} \right]_{k \times k}$$

$$p(x) = p(h(x)) |h'(x)| +$$

density of  $x$       density of  $\Theta$       absolute value of the determinant of  $h'(x)$

Example, Linear,  $R^k \Rightarrow R^k$

$$Z = (Z_1, Z_2, \dots, Z_k)'.$$

$\mu \in R^k$ ,  $A$ ,  $k \times k$ , invertible.

$$y = \mu + Az, \quad z = A^{-1}(y - \mu), \quad \frac{dz}{dy} = A^{-1}.$$

$$f(y) = f_z(A^{-1}(y - \mu)) |A^{-1}|.$$

### 3. Orthogonal Matrices and Rotation

A matrix  $p \times p$  matrix  $P$  is orthogonal if

$$P'P = PP' = I$$

where  $I$  is the identity matrix.

This means all the rows and columns have euclidean length 1 and all the rows are orthogonal to each other and all the columns are orthogonal to each other.

$$P = [\phi_1, \phi_2, \dots, \phi_p]$$

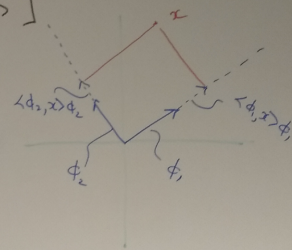
$$I = P^T P = \begin{bmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_p^T \end{bmatrix} [\phi_1, \phi_2, \dots, \phi_p] = [\langle \phi_i, \phi_j \rangle]$$

$$\Rightarrow \langle \phi_i, \phi_j \rangle = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

The columns of  $P$  (or the rows) form an orthonormal basis for  $R^p$ .

$$\begin{aligned} x &\in \mathbb{R}^p \\ x &= PP^T x \\ P^T x &= \begin{bmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_p^T \end{bmatrix} x = \begin{bmatrix} \langle \phi_1, x \rangle \\ \langle \phi_2, x \rangle \\ \vdots \\ \langle \phi_p, x \rangle \end{bmatrix} \end{aligned}$$

$x = \langle \phi_1, x \rangle \phi_1 + \langle \phi_2, x \rangle \phi_2$

$$\begin{aligned} x &= PP^T x = \sum \langle \phi_i, x \rangle \phi_i \\ &= \sum \frac{\langle \phi_i, x \rangle}{\langle \phi_i, \phi_i \rangle} \phi_i \end{aligned}$$


$P$  may be viewed as a rotation.

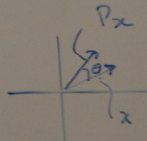
$P$  &  $P^T$  are rotations

$$\begin{aligned}\|Px\|^2 &= (Px)^T(Px) \\ &= x^T P^T P x \\ &= x^T x = \|x\|^2\end{aligned}$$

In  $\mathbb{R}^2$  :  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

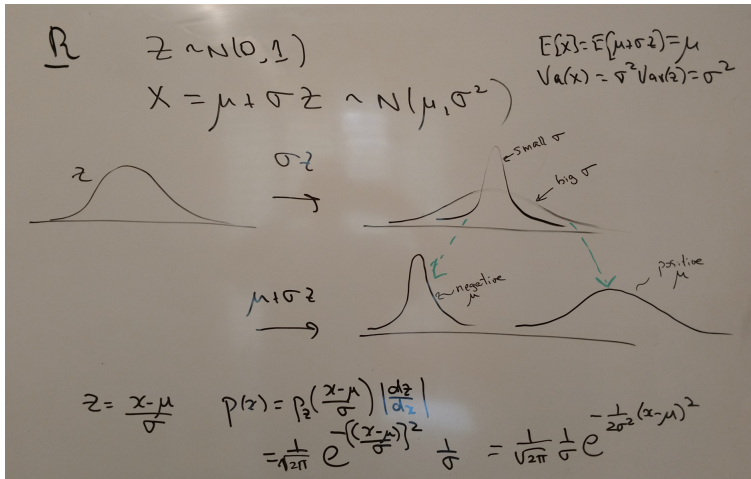
$$P = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} = \begin{array}{l} \text{counter} \\ \text{clockwise} \\ \text{by} \\ \theta \end{array}$$

$$P^T = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} = \begin{array}{l} \text{clockwise} \\ \text{by} \\ \theta \end{array}$$



## 4. Multivariate Normal

In the univariate normal case it is useful to think a general  $Y \sim N(\mu, \sigma^2)$  as a linear function of a standard normal:



What about the multivariate normal? Can we express it as a linear function of a "standard normal"?

$$\begin{aligned} \mathbb{R}^p \\ \mathbf{z} &= \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_p \end{bmatrix} & z_i &\sim N(0,1) & \mathbf{z} &\sim N(\mathbf{0}, \mathbf{I}) \\ & & i &= 1, 2, \dots, p \\ E(\mathbf{z}) &= \mathbf{0} & \text{Var}(\mathbf{z}) &= \mathbf{I} \\ \mathbf{y} &= \boldsymbol{\mu} + \mathbf{A}\mathbf{z} & \mathbf{A} & p \times p & |\mathbf{A}| \neq 0 \\ E(\mathbf{y}) &= \boldsymbol{\mu} + \mathbf{A}E(\mathbf{z}) = \boldsymbol{\mu} & \mathbf{y} &\sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \text{Var}(\mathbf{y}) &= \mathbf{A}\text{Var}(\mathbf{z})\mathbf{A}^T \\ &= \mathbf{A}\mathbf{I}\mathbf{A}^T = \mathbf{A}\mathbf{A}^T \\ \text{Let } \boldsymbol{\Sigma} &= \mathbf{A}\mathbf{A}^T \end{aligned}$$



The multivariate normal density from the change of variable and  $Y = \mu + AZ$ :

$$f(z) = f(z_1, z_2, \dots, z_p) \\ = \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} z_i^2} = (2\pi)^{-\frac{p}{2}} e^{-\frac{1}{2} \|z\|^2}$$

$$f(y) = (2\pi)^{-\frac{p}{2}} e^{-\frac{1}{2} \|A^{-1}(y-\mu)\|^2} |A^{-1}|$$

$$\begin{aligned} \|A^{-1}(y-\mu)\|^2 &= (y-\mu)^T (A^{-1})^T A^{-1} (y-\mu) \\ &= (y-\mu)^T \{AA^T\}^{-1} (y-\mu) \\ &= (y-\mu)^T \Sigma^{-1} (y-\mu) \end{aligned}$$

$$|\Sigma|^{-\frac{1}{2}} = |AA^T|^{-\frac{1}{2}} = (|A|^2)^{-\frac{1}{2}} = |A^{-1}|$$

$$f(y) = (2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} (y-\mu)^T \Sigma^{-1} (y-\mu)}$$

But, can we choose  $A$  in such a way that it tell us a nice story about how the  $Z_i$  are combined to create a dependent structure embodied in a given  $\Sigma$ ?

Given  $\Sigma$ , there is more than one way to choose  $A$  such that  $\Sigma = AA^T$ !!!!

## Choleski Decomposition

Given symmetric, positive definite  $\Sigma$  we can always write  $\Sigma = AA^T$  where  $A$  is lower triangular.

In  $R^2$  we have:

$$Y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} a_{11} & 0 \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$y_1 = \mu_1 + a_{11} z_1$$

$$y_2 = \mu_2 + a_{21} z_1 + a_{22} z_2$$

$$z_1, z_2 \text{ iid } N(0, 1)$$

## Eigen Decomposition

Also called the spectral decomposition.

We can always write a symmetric positive definite  $\Sigma = PDP^T$ .

$$\Sigma = P D P^T$$

$P$ : orthogonal

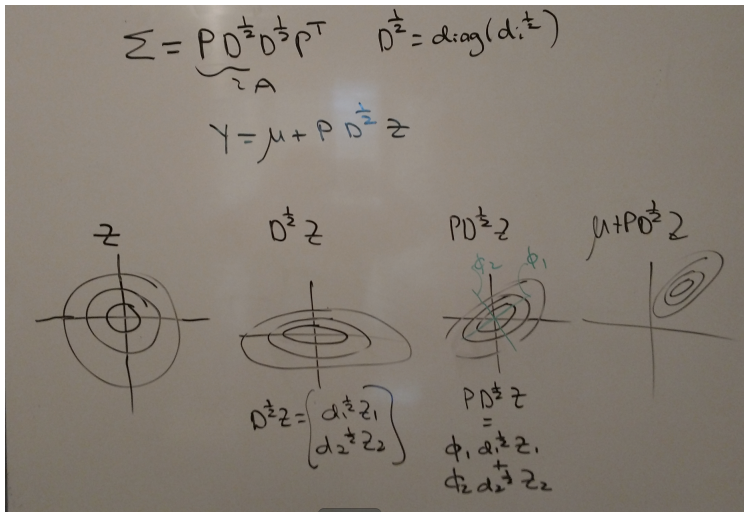
$D$ :  $\text{diag}(d_1, d_2, \dots, d_p)$   $d_i \geq d_{i+1}$   
 $d_1 > 0$

$$P = [\phi_1, \phi_2, \dots, \phi_p]$$

$$\Sigma P = P D \Rightarrow \Sigma \phi_i = d_i \phi_i$$

The columns of  $P$  are the eigen vectors of  $\Sigma$  and the diagonal elements are the corresponding eigen values.

The geometric picture is:



?worth a thousand words?

Note:

$A$  symmetric, pd.

$$A = PDP'$$

(i)

$$|A| = |P|^2 |D| = |D| = \prod d_{ii}$$

(ii)

$$\text{tr}(A) = \text{tr}(DP'P) = \text{tr}(D) = \sum d_{ii}$$

Note:

$$A = PDP'$$

$$D^{\frac{1}{2}} = \text{diag}(d_{ii}^{\frac{1}{2}}).$$

$$A = PD^{\frac{1}{2}}D^{\frac{1}{2}}P' = PD^{\frac{1}{2}}P'PD^{\frac{1}{2}}P'$$

$$\text{Let } A^{\frac{1}{2}} = PD^{\frac{1}{2}}P'.$$

So,

$$A = A^{\frac{1}{2}}A^{\frac{1}{2}}.$$

$A^{\frac{1}{2}}$  is called the symmetric pd square root of  $A$ .

## 5. The Choleski Decomposition

Not only is the Choleski decomposition very powerful, you can figure out basic things about it very simply!!

Simple and powerful, my favorite!!



## Computing the Choleski:

### Choleski:

$A$  symmetric, positive definite  $\rightarrow \exists$  lower triangular  $L$  such that

$$A = LL'$$

To compute  $L$ , you can recursively solve the system of equations give by

$$LL' = A$$

The simple  $2 \times 2$  case:

$$L = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \quad A = \begin{bmatrix} a_{11} & a_{21} \\ a_{21} & a_{22} \end{bmatrix}$$

$$LL^T = A:$$

$$LL^T = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} L_{11} & L_{21} \\ 0 & L_{22} \end{bmatrix} = \begin{bmatrix} L_{11}^2 & L_{11}L_{21} \\ L_{21}L_{11} & L_{21}^2 + L_{22}^2 \end{bmatrix}$$

$$L_{11} = \sqrt{a_{11}}$$

$$L_{21} = a_{21} / L_{11}$$

$$L_{22} = (a_{22} - L_{21}^2)^{\frac{1}{2}}$$

In general we have:

The image shows two handwritten matrices,  $L$  and  $L^T$ , illustrating the structure of a lower triangular matrix in LU decomposition.

Matrix  $L$  is shown as:

$$L = \begin{bmatrix} L_{11} & 0 & - & - & - & - & 0 \\ L_{21} & L_{22} & 0 & - & - & - & 0 \\ & & & & & & \\ & & & & & & \\ L_{j1} & L_{j2} & - & - & L_{jj} & 0 & - & - & 0 \end{bmatrix}$$

Matrix  $L^T$  is shown as:

$$L^T = \begin{bmatrix} L_{11} & L_{21} & \dots & L_{j1} & \dots & L_{j1} \\ 0 & L_{22} & L_{32} & L_{j2} & \dots & 0 \\ & 0 & L_{33} & L_{j3} & \dots & 0 \\ & & & L_{jj} & \dots & 0 \\ & & & & & & & & 0 \end{bmatrix}$$

Notice that the top  $2 \times 2$  corner is just like the simple  $2 \times 2$  case!

After that we can do solve for  $L$  by iterating over the rows, and doing each row by iterating over the columns.

Assume we know all the rows of  $L$  for rows  $1, 2, \dots, (j-1)$ .

$j^{\text{th}}$  row of  $L$  times first column of  $L'$ :

$$L_{j1} L_{11} = a_{j1} \rightarrow L_{j1} = a_{j1}/L_{11}.$$

$$L = \begin{bmatrix} L_{11} & 0 & \dots & 0 \\ L_{21} & L_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ L_{j1} & L_{j2} & \dots & L_{jj} & 0 & \dots & 0 \end{bmatrix} \quad L^T = \begin{bmatrix} L_{11} & L_{21} & \dots & L_{j1} & 0 & \dots & 0 \\ 0 & L_{22} & \dots & L_{j2} & 0 & \dots & 0 \\ \vdots & 0 & \ddots & L_{jj} & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix}$$

$j^{th}$  row of  $L$  times second column of  $L'$ :

$$L_{j1}L_{21} + L_{j2}L_{22} = a_{j2} \rightarrow L_{j2} = (a_{j2} - L_{j1}L_{21})/L_{22}$$

$$\begin{array}{c} L \qquad \qquad \qquad L^T \\
 \left[ \begin{array}{cccccccc} L_{11} & 0 & - & - & - & - & - & 0 \\ L_{21} & L_{22} & 0 & - & - & - & - & 0 \\ & & & & & & & \\ & & & & & & & \\ L_{j1} & L_{j2} & - & - & - & L_{ji} & 0 & - & - & - & 0 \end{array} \right] & \left[ \begin{array}{cccc} L_{11} & L_{21} & \dots & L_{j1} & L_{j1} \\ 0 & L_{22} & & L_{j2} & L_{j2} \\ & 0 & & L_{ji} & L_{ji} \\ & & & 0 & L_{jj} \\ & & & & 0 \\ & & & & & & & & & & \\ & & & & & & & & & & \\ 0 & 0 & & 0 & 0 & - & - & - & 0 \end{array} \right]
 \end{array}$$

$j^{th}$  row of  $L$  times  $i^{th}$  column of  $L'$ , ( $j > i$ ):

$$\sum_{k=1}^i L_{jk} L_{ik} = a_{ji} \rightarrow L_{ji} = (a_{ji} - \sum_{k=1}^{(i-1)} L_{jk} L_{ik}) / L_{ii}$$

and, finally,

$$\begin{array}{c}
 L \qquad \qquad \qquad L^T \\
 \left[ \begin{array}{cccccccc}
 L_{11} & 0 & - & - & - & - & - & 0 \\
 L_{21} & L_{22} & 0 & - & - & - & - & 0 \\
 & & & & & & & \\
 & & & & & & & \\
 & & & & & & & \\
 & & & & & & & \\
 & & & & & & & \\
 L_{j1} & L_{j2} & - & - & - & L_{jj} & 0 & - & - & - & 0
 \end{array} \right]
 \left[ \begin{array}{cccccccc}
 L_{11} & L_{21} & \dots & L_{j1} & \dots & L_{j1} \\
 0 & L_{22} & & L_{2j} & & L_{2j} \\
 & 0 & & L_{3j} & & L_{3j} \\
 & & & \vdots & & \vdots \\
 & & & L_{ji} & & 0 \\
 & & & 0 & & L_{jj} \\
 & & & \vdots & & \vdots \\
 & & & \vdots & & \vdots \\
 & & & \vdots & & \vdots \\
 0 & 0 & & 0 & & 0 & - & - & - & 0
 \end{array} \right]
 \end{array}$$

$j^{\text{th}}$  row of  $L$  times  $j^{\text{th}}$  column of  $L'$ , ( $j > i$ ):

$$\sum_{k=1}^j L_{jk}^2 = a_{jj} \rightarrow L_{jj} = (a_{jj} - \sum_{k=1}^{(j-1)} L_{jk}^2)^{1/2}$$

## Other basic properties:

(i)

For  $L$  Lt (lower triangular),  $L^{-1}$  is Lt and fast to compute.

(ii) The system

$$Lx = b$$

is quickly recursively solved.

(iii)

If  $A$  is symmetric, pd, then the system

$$Ax = b$$

can be solved by

$$A = LL' \rightarrow LL'x = b \rightarrow L'x = L^{-1}b$$

Let  $y = L^{-1}b$  and solve for  $y$  using  $Ly = b$ .

Then solve for  $x$  using  $L'x = y$ .



As previously noted:

Solve  $y \approx Xb$  for  $b$ .

solve:  $X^T X b = X^T y$

$X = QR$ :  $X^T X = R^T R$   $X^T y = R^T Q^T y$

solve:  $R^T R b = R^T Q^T y$

solve:  $Rb = Q^T y$

See Murphy, page section 7.5.2.

$QR$  is  $O(np^2)$ .

## 6. More on the Multivariate Normal

We'll use the Choleski decomposition to derive fundamental properties of the multivariate normal distribution.

- ▶ (a) The marginal from a multivariate normal.
- ▶ (b) For normals, uncorrelated  $\Rightarrow$  independent.
- ▶ (c) The conditional from a multivariate normal.
- ▶ (d) Linear of normal is normal.

We partition a normal vector into  $X$  and  $Y$ .

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

$$\Sigma_{yx} = \Sigma_{xy}^T$$

$$\Sigma_{xy} = E \left( (X - \mu_x)(Y - \mu_y)^T \right)$$

We take the choleski root of  $\Sigma$ .

$$\Sigma = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

$$\Sigma = L L^T \quad L \text{ lower triangular}$$

$$L = \begin{bmatrix} L_1 & 0 \\ A & L_2 \end{bmatrix} \quad L_i \text{ lower tri}$$

We have  $(X, Y)'$  in terms of the choleski.

We have  $\Sigma$  in terms of the choleski.

$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} + \begin{bmatrix} L_1 & 0 \\ A & L_2 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$X = \mu_x + L_1 z_1$$

$$Y = \mu_y + A z_1 + L_2 z_2$$

$$\begin{aligned} L L^T &= \begin{bmatrix} L_1 & 0 \\ A & L_2 \end{bmatrix} \begin{bmatrix} L_1^T & A^T \\ 0 & L_2^T \end{bmatrix} = \begin{bmatrix} L_1 L_1^T & L_1 A^T \\ A L_1^T & A A^T + L_2 L_2^T \end{bmatrix} \\ &= \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \end{aligned}$$

(a) Marginal of  $X$ , (b) uncorrelated implies independence.

$$(a) \quad X = \mu_X + L_1 Z$$

$$\Rightarrow X \sim N(\mu_X, L_1 L_1^T) = N(\mu_X, \Sigma_{XX})$$

$$(b) \quad \Sigma_{XY} = 0 \Rightarrow A = 0$$

$$X = \mu_X + L_1 Z_1$$

$$Y = \mu_Y + L_2 Z_2$$

$$Z_1 \perp Z_2 \Rightarrow X \perp Y$$

$$(c) Y|X=x.$$

$$(c) \quad Y - \mu_Y = A Z_1 + L_2 Z_2$$

$Z_1$  and  $X$  are 1-1

$$Y - \mu_Y = A L_1^{-1} (X - \mu_X) + L_2 Z_2$$

$$Z_2 \perp\!\!\!\perp X$$

$$Y|X=x \sim N(\mu_Y + A L_1^{-1} (x - \mu_X), L_2 L_2^T)$$

Solve for  $x$  coefficients in terms of  $\Sigma$ .

$$A L_1^T = \sum y x$$

$$A L_1^{-1} L_1 L_1^T = \sum y x$$

$$A L_1^{-1} \Sigma_{xx} = \sum y x$$

$$A L_1^{-1} = \sum y x \Sigma_{xx}^{-1}$$



$$Y|X=x.$$

$$Y = \mu_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (X - \mu_X) + E$$

$$E = I_2 Z_2 \perp\!\!\!\perp X$$

$$\Sigma_{YY} = \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XX} \Sigma_{XX}^{-1} \Sigma_{XY} + \text{Var}(E)$$

$$\text{Var}(E) = \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$$

$$Y|X=x \sim N\left(\mu_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (x - \mu_X), \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}\right)$$

An important special case:

$$\begin{bmatrix} Y \\ X \end{bmatrix} \quad \begin{array}{l} Y \in \mathbb{R} \\ X \in \mathbb{R}^p \end{array}$$

$$\begin{bmatrix} Y \\ X \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_Y \\ \mu_X \end{bmatrix}, \begin{bmatrix} \sigma_Y^2 & \sigma_{YX} \\ \sigma_{XY} & \Sigma_{XX} \end{bmatrix} \right)$$

$1 \times p$   
 $\downarrow$   
 $p \times 1$   
 $\uparrow$   
 $p \times p$

$$\sigma_{YX} = E((Y - \mu_Y)(X - \mu_X)^T); \quad \sigma_{XY} = \sigma_{YX}^T$$

$$Y | X=x \sim \mathcal{N}(\mu_Y + \sigma_{YX} \Sigma_{XX}^{-1} (x - \mu_X), \sigma_Y^2 - \sigma_{YX} \Sigma_{XX}^{-1} \sigma_{XY})$$

So, if  $(X, Y)$  are multivariate normal,  $X \in R^p$ ,  $Y \in R$ , then,

$$Y | X=x \sim N(\mu_Y + \sigma_{YX} \Sigma_{XX}^{-1} (x - \mu_X), \\ \sigma_Y^2 - \sigma_{YX} \Sigma_{XX}^{-1} \sigma_{XY})$$

$$\beta = \Sigma_{XX}^{-1} \sigma_{XY}$$

$$\sigma^2 = \sigma_Y^2 - \sigma_{YX} \Sigma_{XX}^{-1} \sigma_{XY}$$

$$Y | X=x \sim N(\mu_Y + (x - \mu_X)^T \beta, \sigma^2)$$

So that the conditional distribution of  $Y | X = x$  has the form of the standard multiple regression model with iid homoscedastic normal errors.

(4)

Example

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

$$x, y \sim N(0, 1)$$

$$y = \rho x + E \quad E \sim N(0, 1 - \rho^2)$$

Note  $\mu_x, \mu_y = 0$ 

$$E((y - Bx)x^T) = 0$$

$$\Rightarrow \Sigma_{yx} = B \Sigma_{xx}$$

$$B = \Sigma_{yx} \Sigma_{xx}^{-1}$$

$$\Rightarrow y = Bx + E \quad E \perp x.$$

(d)

$$Y \sim N(\mu, \Sigma)$$

$$\Rightarrow a + BY \sim N(a + B\mu, B\Sigma B')$$

$$X = a + BY$$

$$= a + B[\mu + AZ] ; AA^T = \Sigma$$

$$= a + B\mu + BAZ ; (BA)(BA^T) = BAA^TB^T \\ = B\Sigma B^T$$

$$\sim N(a + B\mu, B\Sigma B^T)$$

## 7. Simulating a Multivariate Normal

Suppose we wish to draw  $Y \sim N(\mu, \Sigma)$ .

Let  $Z = (Z_1, Z_2, \dots, Z_p)'$ ,  $Z_j \sim N(0, 1)$ , *iid*.

Then let,

$$Y = \mu + AZ$$

where,

$$\Sigma = AA'$$

If  $A$  is cholesky, multiplication  $AZ$  is fast.

## 8. Likelihood, Sufficiency, and MLE

Let's use our spectral decomposition to learn about the multivariate normal likelihood.

Let,

$$X_i \sim N_p(\mu, \Sigma), \text{ iid, } i = 1, 2, \dots, n.$$

$$X_i = (X_{i1}, X_{i2}, \dots, X_{ip})'$$

Recall that for a parametric model,

$$f(y \mid \theta), \quad \theta \in \Theta,$$

given data,  $y$ , the maximum likelihood estimator is obtained by finding the  $\theta$  that makes what you have seen most likely:

$$\hat{\theta} = \operatorname{argmax}_{\theta} f(y \mid \theta)$$

In practice we often maximize the log of the likelihood or minimize the negative of the log likelihood.



## Example:

### Bernoulli: MLE

$$Y_i \sim \text{Bern}(p) \quad Y_i \in \{0, 1\}$$

$$\begin{aligned} p(y_1, y_2, \dots, y_n | p) &= \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i} \\ &= p^k (1-p)^{n-k} \quad k = \#(Y_i = 1) \end{aligned}$$

$$\log p = k \log p + (n-k) \log(1-p)$$

$$\begin{aligned} \text{FOC: } \frac{k}{p} - \frac{(n-k)}{1-p} &= 0 \Rightarrow (n-k)p = k(1-p) \\ &\Rightarrow p = \frac{k}{n} \end{aligned}$$

FOC: "first order condition",  $f' = 0$ .

So, the observed sample frequency is the MLE!

In our problem we will observe  $X_i = x_i$  for  
 $X_i \sim N_p(\mu, \Sigma)$ , *iid*,  $i = 1, 2, \dots, n$ .

note:

$x$  a  $p$  dimensional column vector.  $A$   $p \times p$ .

$$x'Ax = tr(x'Ax) = tr(Axx'),$$

where  $tr$  is the trace.

$$\begin{aligned}
 p(x_1, x_2, \dots, x_n) &= \\
 &\prod_{i=1}^n (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left[-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)\right] \\
 &= (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left[-\frac{1}{2} \sum_i (x_i - \mu)^T \Sigma^{-1}(x_i - \mu)\right] \\
 &= (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left[-\frac{1}{2} \sum_i \text{tr}(\Sigma^{-1}(x_i - \mu)(x_i - \mu)^T)\right] \\
 &= (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left[-\frac{1}{2} \text{tr}(\Sigma^{-1} \sum_i (x_i - \mu)(x_i - \mu)^T)\right]
 \end{aligned}$$

Note: 
$$\bar{x} = \frac{1}{n} \sum x_i ; \begin{aligned} \sum (x_i - \bar{x}) &= \sum x_i - n\bar{x} \\ &= n\bar{x} - n\bar{x} = 0. \end{aligned}$$

$$\begin{aligned} & \sum (\boldsymbol{x}_i - \mu)(\boldsymbol{x}_i - \mu)^T \\ = & \sum ((\boldsymbol{x}_i - \bar{\boldsymbol{x}}) - (\mu - \bar{\boldsymbol{x}}))(\quad)^T \\ = & \sum (\boldsymbol{x}_i - \bar{\boldsymbol{x}})(\boldsymbol{x}_i - \bar{\boldsymbol{x}})^T + n(\mu - \bar{\boldsymbol{x}})(\mu - \bar{\boldsymbol{x}})^T \\ & \text{+ cross terms } \end{aligned}$$

" $\circ$

Cross terms:

$$\sum (\mu - \bar{x})(x_i - \bar{x})^T = (\mu - \bar{x}) \sum (x_i - \bar{x})^T = 0$$

$$A = \sum_i (x_i - \bar{x})(x_i - \bar{x})'$$

$$\begin{aligned} & \text{tr} \left( \Sigma^{-1} \sum_i (x_i - \mu)(x_i - \mu)' \right) \\ &= \text{tr} \left( \Sigma^{-1} (A + n(\bar{x} - \mu)(\bar{x} - \mu)') \right) \\ &= \text{tr} \Sigma^{-1} A + n \text{tr} \Sigma^{-1} (\bar{x} - \mu)(\bar{x} - \mu)' \\ &= \text{tr} \Sigma^{-1} A + n (\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu) \end{aligned}$$

## Sufficiency:

Given data, functions of the data are *sufficient* if they are all we need to compute the likelihood.

Clearly, for iid MVN data,

$$\bar{x} \text{ and } A$$

are sufficient.

$p + \frac{p(p+1)}{2}$  quantities instead of the  $np$  data.

## What is $A$ ?

The  $k, j$  element of  $A$  is:

$$A_{jk} = \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

The sample covariance between  $X_j$  and  $X_k$  is

$$s_{jk} = \frac{A_{jk}}{(n-1)}$$

The sample variance of  $X_j$  is

$$s_{jj} = \frac{A_{jj}}{(n-1)}$$

MLE:

$$L \propto |\Sigma|^{-n/2} \exp(\text{tr}(-\frac{1}{2}\Sigma^{-1}A)) \exp(-\frac{n}{2}(\bar{x} - \mu)' \Sigma^{-1}(\bar{x} - \mu))$$

Clearly, for any  $\Sigma$ , maximum over  $\mu$  is attained at

$$\hat{\mu} = \bar{x}$$

Notation:  $\text{etr}(A) = \exp(\text{tr}(A))$ .



$$L(\hat{\mu}, \Sigma) \propto |\Sigma|^{-\frac{n}{2}} \text{etr}\left[-\frac{1}{2} \Sigma^{-1} A\right]$$

$$A = T T' \quad ; \quad \Gamma = T' \Sigma^{-1} T$$

$$|\Gamma| = \frac{|T|^2}{|\Sigma|}$$

$$L(\hat{\mu}, \Gamma) \propto |\Gamma|^{\frac{n}{2}} \text{etr}\left[-\frac{n}{2} \frac{\Gamma}{n}\right]$$

$\gamma_i$  : roots of  $\frac{f}{n}$

$$L \propto \prod_i \gamma_i^{\frac{n}{2}} e^{-\frac{n}{2} \gamma_i}$$

$$\text{Let } a = \frac{n}{2}$$

$$\max_x x^a e^{-ax}$$

$$\max_x a \log(x) - ax$$

Foc: first order condition  
— set derivative equal to 0

$$\frac{a}{x} = a \Rightarrow x^* = 1.$$

$$\text{So } \frac{\hat{\Gamma}}{n} = P I P^T = I$$

$$\hat{\Gamma} = n I$$

Now we simply solve for  $\hat{\Sigma}$

$$\hat{\Gamma} = T' \hat{\Sigma}^{-1} T = n I$$

$$\begin{aligned} \hat{\Sigma}^{-1} &= n (T')^{-1} T^{-1} = n (T T')^{-1} \\ &= n A^{-1} \end{aligned}$$

$$\hat{\Sigma} = \frac{A}{n}$$

## 9. Checking for Normality

Suppose  $Y \sim N(\mu, \Sigma)$ .

$$\Sigma = PD^{\frac{1}{2}}D^{\frac{1}{2}}P'.$$

$$\Sigma^{-1} = PD^{-\frac{1}{2}}D^{-\frac{1}{2}}P'.$$

Then

$$Y = \mu + PD^{\frac{1}{2}}Z, \quad Z \sim N(0, I).$$

So,

$$Z = D^{-\frac{1}{2}}P'(Y - \mu).$$

$$Z'Z = (Y - \mu)'PD^{-\frac{1}{2}}D^{-\frac{1}{2}}P'(Y - \mu) = (Y - \mu)'\Sigma^{-1}(Y - \mu).$$

So,

$$(Y - \mu)'\Sigma^{-1}(Y - \mu) = Z'Z = \sum Z_i^2 \sim \chi_p^2.$$

So if  $Y_i \sim N(\mu, \Sigma)$ , *iid*, then

$$D_i = (Y_i - \hat{\mu})' \hat{\Sigma}^{-1} (Y_i - \hat{\mu}) \approx \chi_p^2, \text{ iid}$$

So you can check to see if the  $D_i$  look right.

I usually use a qqplot.

## 10. Weighted Regression

### Review Linear Case

Usual:  $Y = X\beta + \varepsilon$   $\varepsilon \sim N(0, \sigma^2 I)$

Consider:  $Y = X\beta + \varepsilon$   $\varepsilon \sim N(0, \Sigma)$

$\Sigma = LL^T$  choleski:  $L$  is  $\Delta$  lower triangular.

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

$$\begin{aligned}\tilde{\varepsilon} &= L^{-1}\varepsilon; \quad E(\tilde{\varepsilon}) = 0 \quad \text{Var}(\tilde{\varepsilon}) = L^{-1}[LL^T](L^{-1})^T \\ &= (L^{-1}L)(L^T(L^{-1})^T) = I.\end{aligned}$$

$$L^{-1}Y = L^{-1}X\beta + L^{-1}\varepsilon$$

$$\tilde{Y} = \tilde{X}\beta + \tilde{\varepsilon} \quad \tilde{\varepsilon} \sim N(0, I)$$

$$\begin{aligned}\hat{\beta} &= (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y} = (X^T (L^{-1})^T L^{-1} X)^{-1} X^T (L^{-1})^T L^{-1} Y \\ &= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y\end{aligned}$$