

## Appendix E

# Multivariate Distributions

### E.1 Review of Definitions

Let's review some definitions from basic probability. When we have a random vector  $\vec{X}$  with  $p$  different components,  $X_1, X_2, \dots, X_p$ , the **joint cumulative distribution function** is

$$F(\vec{a}) = F(a_1, a_2, \dots, a_p) = \Pr(X_1 \leq a_1, X_2 \leq a_2, \dots, X_p \leq a_p) \quad (\text{E.1})$$

Thus

$$F(\vec{b}) - F(\vec{a}) = \Pr(a_1 < X_1 \leq b_1, a_2 < X_2 \leq b_2, \dots, a_p < X_p \leq b_p) \quad (\text{E.2})$$

This is the probability that  $X$  is in a (hyper-)rectangle, rather than just in an interval. The **joint probability density function** is

$$p(\vec{x}) = p(x_1, x_2, \dots, x_p) = \left. \frac{\partial^p F(a_1, \dots, a_p)}{\partial a_1 \dots \partial a_p} \right|_{\vec{a}=\vec{x}} \quad (\text{E.3})$$

Of course,

$$F(\vec{a}) = \int_{-\infty}^{a_1} \int_{-\infty}^{a_2} \dots \int_{-\infty}^{a_p} p(x_1, x_2, \dots, x_p) dx_p \dots dx_2 dx_1 \quad (\text{E.4})$$

(In this case, the order of integration doesn't matter. Why?)

From these, and especially from the joint PDF, we can recover the marginal PDF of any group of variables, say those numbered 1 through  $q$ ,

$$p(x_1, x_2, \dots, x_q) = \int p(x_1, x_2, \dots, x_p) dx_{q+1} dx_{q+2} \dots dx_p \quad (\text{E.5})$$

(What are the limits of integration here?) Then the conditional pdf for some variables given the others — say, use variables 1 through  $q$  to condition those numbered  $q + 1$

through  $p$  — just comes from division:

$$p(x_{q+1}, x_{q+2}, \dots, x_p | X_1 = x_1, \dots, X_q = x_q) = \frac{p(x_1, x_2, \dots, x_p)}{p(x_1, x_2, \dots, x_q)} \quad (\text{E.6})$$

These two tricks can be iterated, so, for instance,

$$p(x_3 | x_1) = \int p(x_3, x_2 | x_1) dx_2 \quad (\text{E.7})$$

## E.2 Multivariate Gaussians

The multivariate Gaussian is just the generalization of the ordinary Gaussian to vectors. Scalar Gaussians are parameterized by a mean  $\mu$  and a variance  $\sigma^2$ , so we write  $X \sim \mathcal{N}(\mu, \sigma^2)$ . Multivariate Gaussians, likewise, are parameterized by a mean vector  $\vec{\mu}$ , and a variance-covariance matrix  $\Sigma$ , written  $\vec{X} \sim \mathcal{MVN}(\vec{\mu}, \Sigma)$ . The components of  $\vec{\mu}$  are the means of the different components of  $\vec{X}$ . The  $i, j^{\text{th}}$  component of  $\Sigma$  is the covariance between  $X_i$  and  $X_j$  (so the diagonal of  $\Sigma$  gives the component variances).

Just as the probability density of scalar Gaussian is

$$p(x) = (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right\} \quad (\text{E.8})$$

the probability density of the multivariate Gaussian is

$$p(\vec{x}) = (2\pi \det \Sigma)^{-p/2} \exp\left\{-\frac{1}{2} (\vec{x} - \vec{\mu}) \cdot \Sigma^{-1} (\vec{x} - \vec{\mu})\right\} \quad (\text{E.9})$$

Finally, remember that the parameters of a Gaussian change along with linear transformations

$$X \sim \mathcal{N}(\mu, \sigma^2) \Leftrightarrow aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2) \quad (\text{E.10})$$

and we can use this to “standardize” any Gaussian to having mean 0 and variance 1 (by looking at  $\frac{X - \mu}{\sigma}$ ). Likewise, if

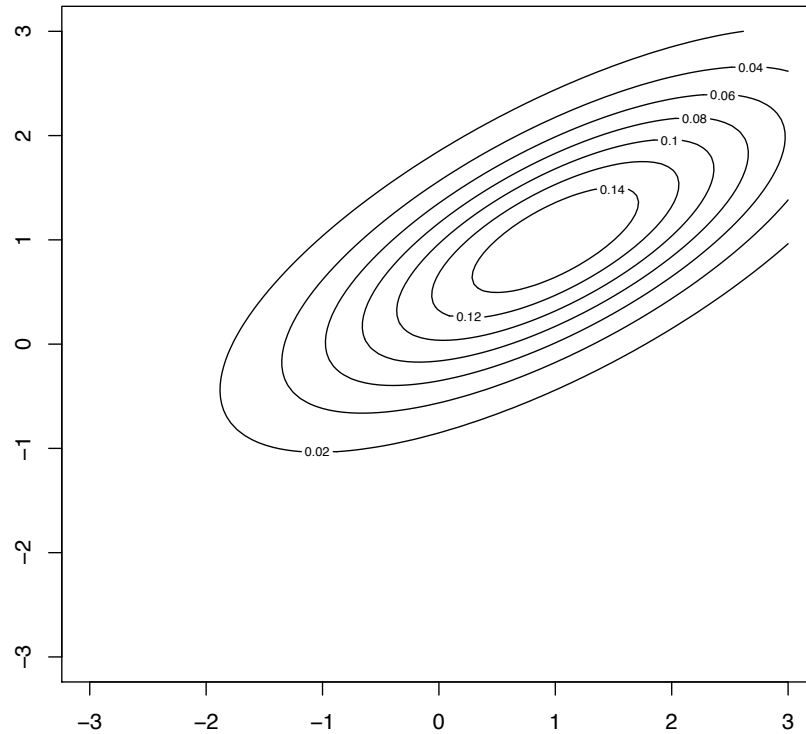
$$\vec{X} \sim \mathcal{MVN}(\vec{\mu}, \Sigma) \quad (\text{E.11})$$

then

$$\mathbf{a}\vec{X} + \vec{b} \sim \mathcal{MVN}(\mathbf{a}\vec{\mu} + \vec{b}, \mathbf{a}\Sigma\mathbf{a}^T) \quad (\text{E.12})$$

In fact, the analogy between the ordinary and the multivariate Gaussian is so complete that it is very common to not really distinguish the two, and write  $\mathcal{N}$  for both.

The multivariate Gaussian density is most easily visualized when  $p = 2$ , as in Figure E.1. The probability contours are ellipses. The density changes comparatively slowly along the major axis, and quickly along the minor axis. The two points marked + in the figure have equal geometric distance from  $\vec{\mu}$ , but the one to its right lies on a higher probability contour than the one above it, because of the directions of their displacements from the mean.



```

library(mvtnorm)
x.points <- seq(-3,3,length.out=100)
y.points <- x.points
z <- matrix(0,nrow=100,ncol=100)
mu <- c(1,1)
sigma <- matrix(c(2,1,1,1),nrow=2)
for (i in 1:100) {
  for (j in 1:100) {
    z[i,j] <- dmvnorm(c(x.points[i],y.points[j]),mean=mu,sigma=sigma)
  }
}
contour(x.points,y.points,z)

```

FIGURE E.1: Probability density contours for a two-dimensional multivariate Gaussian, with mean  $\vec{\mu} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  (solid dot), and variance matrix  $\Sigma = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$ . Using `expand.grid`, as in Chapter 4, would be more elegant coding than this double `for` loop.

### E.2.1 Linear Algebra and the Covariance Matrix

We can use some facts from linear algebra to understand the general pattern here, for arbitrary multivariate Gaussians in an arbitrary number of dimensions. The covariance matrix  $\Sigma$  is symmetric and positive-definite, so we know from linear algebra that it can be written in terms of its eigenvalues and eigenvectors:

$$\Sigma = \mathbf{v}^T \mathbf{d} \mathbf{v} \quad (\text{E.13})$$

where  $\mathbf{d}$  is the diagonal matrix of the eigenvalues of  $\Sigma$ , and  $\mathbf{v}$  is the matrix whose columns are the eigenvectors of  $\Sigma$ . Because the eigenvectors are all of length 1, and they are all perpendicular to each other, it is easy to check that  $\mathbf{v}^T \mathbf{v} = \mathbf{I}$ , so  $\mathbf{v}^{-1} = \mathbf{v}^T$  and  $\mathbf{v}$  is an orthogonal matrix. What actually shows up in the equation for the multivariate Gaussian density is  $\Sigma^{-1}$ , which is

$$(\mathbf{v}^T \mathbf{d} \mathbf{v})^{-1} = \mathbf{v}^{-1} \mathbf{d}^{-1} (\mathbf{v}^T)^{-1} = \mathbf{v}^T \mathbf{d}^{-1} \mathbf{v} \quad (\text{E.14})$$

Geometrically, orthogonal matrices represent rotations. Multiplying by  $\mathbf{v}$  rotates the coordinate axes so that they are parallel to the eigenvectors of  $\Sigma$ . Probabilistically, this tells us that the axes of the probability-contour ellipse are parallel to those eigenvectors. The radii of those axes are proportional to the square roots of the eigenvalues. To see *that*, look carefully at the math. Fix a level for the probability density whose contour we want, say  $f_0$ . Then we have

$$f_0 = (2\pi \det \Sigma)^{-p/2} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu}) \cdot \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\} \quad (\text{E.15})$$

$$c = (\vec{x} - \vec{\mu}) \cdot \Sigma^{-1} (\vec{x} - \vec{\mu}) \quad (\text{E.16})$$

$$= (\vec{x} - \vec{\mu})^T \mathbf{v}^T \mathbf{d}^{-1} \mathbf{v} (\vec{x} - \vec{\mu}) \quad (\text{E.17})$$

$$= (\vec{x} - \vec{\mu})^T \mathbf{v}^T \mathbf{d}^{-1/2} \mathbf{d}^{-1/2} \mathbf{v} (\vec{x} - \vec{\mu}) \quad (\text{E.18})$$

$$= (\mathbf{d}^{-1/2} \mathbf{v} (\vec{x} - \vec{\mu}))^T (\mathbf{d}^{-1/2} \mathbf{v} (\vec{x} - \vec{\mu})) \quad (\text{E.19})$$

$$= \left\| \mathbf{d}^{-1/2} \mathbf{v} (\vec{x} - \vec{\mu}) \right\|^2 \quad (\text{E.20})$$

where  $c$  combines  $f_0$  and all the other constant factors, and  $\mathbf{d}^{-1/2}$  is the diagonal matrix whose entries are one over the square roots of the eigenvalues of  $\Sigma$ . The  $\mathbf{v}(\vec{x} - \vec{\mu})$  term takes the displacement of  $\vec{x}$  from the mean,  $\vec{\mu}$ , and replaces the components of that vector with its projection on to the eigenvectors. Multiplying by  $\mathbf{d}^{-1/2}$  then scales those projections, and so the radii have to be proportional to the square roots of the eigenvalues.<sup>1</sup>

<sup>1</sup>If you know about principal components analysis and think that all this manipulation of eigenvectors and eigenvalues of the covariance matrix seems familiar, you're right; this was one of the ways in which PCA was originally discovered. But PCA does not require any distributional assumptions. If you do not know about PCA, read Chapter 16.

### E.2.2 Conditional Distributions and Least Squares

Suppose that  $\vec{X}$  is bivariate, so  $p = 2$ , with mean vector  $\vec{m} = (\mu_1, \mu_2)$ , and variance matrix  $\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ . One can show (exercise!) that the conditional distribution of  $X_2$  given  $X_1$  is Gaussian, and in fact

$$X_2|X_1 = x_1 \sim \mathcal{N}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}) \quad (\text{E.21})$$

To understand what is going on here, remember from Chapter 1 that the optimal slope for linearly regressing  $X_2$  on  $X_1$  would be  $\text{Cov}[X_2, X_1] / \mathbb{V}[X_1]$ . This is *precisely* the same as  $\Sigma_{21}\Sigma_{11}^{-1}$ . So in the bivariate Gaussian case, the best linear regression and the optimal regression are exactly the same — there is no need to consider nonlinear regressions. Moreover, we get the same conditional variance for each value of  $x_1$ , so the regression of  $X_2$  on  $X_1$  is homoskedastic, with independent Gaussian noise. This is, in short, exactly the situation which all the standard regression formulas aim at.

More generally, if  $X_1, X_2, \dots, X_p$  are multivariate Gaussian, then conditioning on  $X_1, \dots, X_q$  gives the remaining variables  $X_{q+1}, \dots, X_p$  a Gaussian distribution as well.

If we say that  $\vec{\mu} = (\vec{\mu}_A, \vec{\mu}_B)$  and  $\Sigma = \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}$ , where  $A$  stands for the conditioning variables and  $B$  for the conditioned, then

$$\vec{X}_B | \vec{X}_A = \vec{x}_a \sim \mathcal{MVN}(\vec{\mu}_B + \Sigma_{BA}\Sigma_{AA}^{-1}(\vec{x}_a - \vec{\mu}_A), \Sigma_{BB} - \Sigma_{BA}\Sigma_{AA}^{-1}\Sigma_{AB}) \quad (\text{E.22})$$

(Remember that here  $\Sigma_{BA} = \Sigma_{AB}^T$  [Why?].) This, too, is just doing a linear regression of  $\vec{X}_B$  on  $\vec{X}_A$ .

### E.2.3 Projections of Multivariate Gaussians

A useful fact about multivariate Gaussians is that all their univariate projections are also Gaussian. That is, if  $\vec{X} \sim \mathcal{MVN}(\vec{\mu}, \Sigma)$ , and we fix any unit vector  $\vec{w}$ , then  $\vec{w} \cdot \vec{X}$  has a Gaussian distribution. This is easy to see if  $\Sigma$  is diagonal: then  $\vec{w} \cdot \vec{X}$  reduces to a sum of independent Gaussians, which we know from basic probability is also Gaussian. But we can use the eigen-decomposition of  $\Sigma$  to check that this holds more generally.

One can also show that the converse is true: if  $\vec{w} \cdot \vec{X}$  is a univariate Gaussian for *every* choice of  $\vec{w}$ , then  $\vec{X}$  must be multivariate Gaussian. This fact is more useful for probability theory than for data analysis<sup>2</sup>, but it's still worth knowing.

### E.2.4 Computing with Multivariate Gaussians

Computationally, it is not hard to write functions to calculate the multivariate Gaussian density, or to generate multivariate Gaussian random vectors. Unfortunately, no

<sup>2</sup>It's a special case of a result called the **Cramér-Wold theorem**, or the **Cramér-Wold device**, which asserts that two random vectors  $\vec{X}$  and  $\vec{Y}$  have the same distribution if and only if  $\vec{w} \cdot \vec{X}$  and  $\vec{w} \cdot \vec{Y}$  have the same distribution for every  $\vec{w}$ .

one seems to have thought to put a standard set of such functions in the basic set of R packages, so you have to use a different library. The MASS library contains a function, `mvrnorm`, for generating multivariate Gaussian random vectors. The `mvtnorm` contains functions for calculating the density, cumulative distribution and quantiles of the multivariate Gaussian, as well as generating random vectors<sup>3</sup>. The package `mixtools`, which will use in Chapter 19 for mixture models, includes functions for the multivariate Gaussian density and for random-vector generation.

## E.3 Inference with Multivariate Distributions

As with univariate distributions, there are several ways of doing statistical inference for multivariate distributions. Here I will focus on parametric inference, since non-parametric inference is covered in Chapter 14.

Parameter estimation by maximum likelihood, the sampling distribution of the MLE, and the resulting hypothesis tests and confidence sets work exactly as they do for one-dimensional distributions. That is to say, they are special cases of general results about estimation by minimizing a loss function, described in App. H.4.

### E.3.1 Model Comparison

Out of sample, models can be compared on log-likelihood. When a strict out-of-sample comparison is not possible, we can use cross-validation.

In sample, a likelihood ratio test can be used. This has two forms, depending on the relationship between the models. Suppose that there is a large or wide model, with parameter  $\Theta$ , and a narrow or small model, with parameter  $\theta$ , which we get by fixing some of the components of  $\Theta$ . Thus the dimension of  $\Theta$  is  $q$  and that of  $\theta$  is  $r < q$ . Since every distribution we can get from the narrow model we can also get from the wide model, in-sample the likelihood of the wide model must always be larger. Thus

$$\ell(\hat{\Theta}) - \ell(\hat{\theta}) \geq 0 \quad (\text{E.23})$$

Here we have a clear null hypothesis, which is that the data comes from the narrower, smaller model. Under this null hypothesis, as  $n \rightarrow \infty$ ,

$$2[\ell(\hat{\Theta}) - \ell(\hat{\theta})] \rightsquigarrow \chi_{q-r}^2 \quad (\text{E.24})$$

provided that the restriction imposed by the small model doesn't place it on the boundary of the parameter space of  $\Theta$ . (See Appendix I.)

For instance, suppose that  $\vec{X}$  is bivariate, and the larger model is an unrestricted Gaussian, so  $\Theta = \left\{ (\mu_1, \mu_2), \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12} & \Sigma_{22} \end{bmatrix} \right\}$ . A possible narrow model might impose the assumption that the components of  $\vec{X}$  are uncorrelated, so  $\theta = \left\{ (\mu_1, \mu_2), \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{bmatrix} \right\}$ .

<sup>3</sup>It also has such functions for multivariate  $t$  distributions, which are to multivariate Gaussians exactly as ordinary  $t$  distributions are to univariate Gaussians.

This is a restriction on the broader model, but not one which is on the boundary of the parameter space, so the large-sample  $\chi^2$  distribution should apply. A restriction which *would* be on the boundary would be to insist that  $X_2$  was constant, so  $\Sigma_{22} = 0$ . (This would also force  $\Sigma_{12} = 0$ .)

If, on the other hand, that we have two models, with parameters  $\theta$  and  $\psi$ , and they are completely non-nested, meaning there are no parameter combinations where

$$p(\cdot; \theta) = p(\cdot; \psi) \quad (\text{E.25})$$

then in many ways things become easier. For *fixed* parameter values  $\theta_0, \psi_0$ , the mean log-likelihood ratio is just an average of IID terms:

$$\frac{1}{n} [\ell(\theta_0) - \ell(\psi_0)] \equiv \frac{1}{n} \sum_{i=1}^n \Lambda_i \quad (\text{E.26})$$

$$= \frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i; \theta_0)}{p(x_i; \psi_0)} \quad (\text{E.27})$$

By the law of large numbers, then, the mean log-likelihood ratio converges to an expected value  $\mathbb{E}[\Lambda]$ . This is positive if  $\theta_0$  has a higher expected log-likelihood than  $\psi_0$ , and negative the other way around. Furthermore, by the central limit theorem, as  $n$  grows, the fluctuations around this expected value are Gaussian, with variance  $\sigma_\Lambda^2/n$ . We can estimate  $\sigma_\Lambda^2$  by the sample variance of  $\log \frac{p(x_i; \theta_0)}{p(x_i; \psi_0)}$ .

Ordinarily, we don't have just a single parameter value for each model, but also ordinarily,  $\hat{\theta}_{MLE}$  and  $\hat{\psi}_{MLE}$  both converge to limits, which we can call  $\theta_0$  and  $\psi_0$ . At the cost of some fancy probability theory, one can show that, in the non-nested case,

$$\frac{\sqrt{n} \ell(\hat{\theta}) - \ell(\hat{\psi})}{n \sigma_\Lambda^2} \rightsquigarrow \mathcal{N}(\mathbb{E}[\Lambda], 1) \quad (\text{E.28})$$

and that we can consistently estimate  $\mathbb{E}[\Lambda]$  and  $\sigma_\Lambda^2$  by “plugging in”  $\hat{\theta}$  and  $\hat{\psi}$  in place of  $\theta_0$  and  $\psi_0$ . This gives the **Vuong test** for comparing the two models Vuong (1989). The null hypothesis in the Vuong test is that the two models are equally good (and neither is exactly true). In this case,

$$V = \frac{1}{\sqrt{n}} \frac{\ell(\hat{\theta}) - \ell(\hat{\psi})}{\hat{\sigma}_\Lambda} \rightsquigarrow \mathcal{N}(0, 1) \quad (\text{E.29})$$

If  $V$  is significantly positive, we have evidence in favor of the  $\theta$  model being better (though not necessarily *true*), while if it is significantly negative we have evidence in favor of the  $\psi$  model being better.

The cases where two models *partially* overlap is complicated; see Vuong (1989) for the gory details<sup>4</sup>.

<sup>4</sup>If you are curious about why this central-limit-theorem argument doesn't work in the nested case, notice that when we have nested models, and the null hypothesis is true, then  $\hat{\theta} \rightarrow \hat{\psi}$ , so the numerator in the Vuong test statistic,  $[\ell(\hat{\theta}) - \ell(\hat{\psi})]/n$ , is converging to zero, but so is the denominator  $\sigma_\Lambda^2$ . Since  $0/0$  is undefined, we need to use a stochastic version of L'Hopital's rule, which gives us back Eq. E.24. See, yet again, Vuong (1989).

### E.3.2 Goodness-of-Fit

For univariate distributions, we often assess goodness-of-fit through the Kolmogorov-Smirnov (KS) test<sup>5</sup>, where the test statistic is

$$d_{KS} = \max_a |\hat{F}_n(a) - F(a)| \quad (\text{E.30})$$

with  $\hat{F}_n$  being the empirical CDF, and  $F$  its theoretical counterpart. The null hypothesis here is that the data were drawn IID from  $F$ , and what Kolmogorov and Smirnov did was to work out the distribution of  $d_{KS}$  under this null hypothesis, and show it was the same for all  $F$  (at least for large  $n$ ). This lets us actually calculate  $p$  values.

We could use such a test statistic for multivariate data, where we'd just take the maximum over vectors  $a$ , rather than scalars. But the problem is that we do not know its sampling distribution under the null hypothesis in the multivariate case — Kolmogorov and Smirnov's arguments don't work there — so we don't know whether a given value of  $d_{KS}$  is large or small or what.

There is however a fairly simple approximate way of turning univariate tests into multivariate ones. Suppose our data consists of vectors  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ . Pick a unit vector  $\vec{w}$ , and set  $z_i = \vec{w} \cdot \vec{x}_i$ . Geometrically, this is just the projection of the data along the direction  $\vec{w}$ , but these projections are *univariate* random variables. If the  $\vec{x}_i$  were drawn from  $F$ , then the  $z_i$  must have been drawn from the corresponding projection of  $F$ , call it  $F_{\vec{w}}$ . If we can work out the latter distribution, then we can apply our favorite univariate test to the  $z_i$ . If the fit is bad, then we know that the  $\vec{x}_i$  can't have come from  $F$ . If the fit is good for the  $z_i$ , then the fit is also good for the  $\vec{x}_i$  — at least along the direction  $\vec{w}$ . Now, we can either carefully pick  $\vec{w}$  to be a direction which we care about for some reason, or we can choose it *randomly*. If the projection of the  $\vec{x}_i$  along several random directions matches that of  $F$ , it becomes rather unlikely that they fail to match overall<sup>6</sup>.

To summarize:

1. Choose a random unit vector  $\vec{W}$ . (For instance, let  $\vec{U} \sim \mathcal{MVN}(0, \mathbf{I}_p)$ , and  $\vec{W} = \vec{U} / \|\vec{U}\|$ .)
2. Calculate  $Z_i = \vec{W} \cdot \vec{x}_i$ .
3. Calculate the corresponding projection of the theoretical distribution  $F$ , call it  $F_{\vec{W}}$ .
4. Apply your favorite univariate goodness-of-fit test to  $Z_i$  and  $F_{\vec{W}}$ .

<sup>5</sup>I discuss the KS test here for concreteness. Much the same ideas apply to the Anderson-Darling test, the Cramér-von Mises test, and others which, not being such good ideas, were only invented by one person.

<sup>6</sup>Theoretically, we appeal to the Cramér-Wold device again: the random vectors  $\vec{X}$  and  $\vec{Y}$  have the same distribution if and only if  $\vec{w} \cdot \vec{X}$  and  $\vec{w} \cdot \vec{Y}$  have the same distribution for every  $\vec{w}$ . Failing to match for any  $\vec{w}$  implies that  $\vec{X}$  and  $\vec{Y}$  have different distributions. Conversely, if  $\vec{X}$  and  $\vec{Y}$  differ in distribution at all,  $\vec{w} \cdot \vec{X}$  must differ in distribution from  $\vec{w} \cdot \vec{Y}$  for *some* choice of  $\vec{w}$ . Randomizing the choice of  $\vec{w}$  gives us power to detect a lot of differences in distribution.



[[ATTN: Multiple comparisons as an appendix topic?]]

5. Repeat (1)-(4) multiple times, with correction for multiple testing.

## E.4 Uncorrelated $\neq$ Independent

As you know, two random variables  $X$  and  $Y$  are **uncorrelated** when their correlation coefficient is zero:

$$\rho(X, Y) = 0 \quad (\text{E.31})$$

Since

$$\rho(X, Y) = \frac{\text{Cov}[X, Y]}{\sqrt{\mathbb{V}[X]}\sqrt{\mathbb{V}[Y]}} \quad (\text{E.32})$$

being uncorrelated is the same as having zero covariance. Since

$$\text{Cov}[X, Y] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \quad (\text{E.33})$$

having zero covariance, and so being uncorrelated, is the same as

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \quad (\text{E.34})$$

One says that “the expectation of the product factors”. If  $\rho(X, Y) \neq 0$ , then  $X$  and  $Y$  are **correlated**.

As you also know, two random variables are **independent** when their joint probability distribution is the product of their marginal probability distributions: for all  $x$  and  $y$ ,

$$p_{X,Y}(x, y) = p_X(x)p_Y(y) \quad (\text{E.35})$$

Equivalently<sup>7</sup>, the conditional distribution is the same as the marginal distribution:

$$p_{Y|X}(y|x) = p_Y(y) \quad (\text{E.36})$$

If  $X$  and  $Y$  are not independent, then they are **dependent**. If, in particular,  $Y$  is a function of  $X$ , then they are always dependent<sup>8</sup>

If  $X$  and  $Y$  are independent, then they are also uncorrelated. To see this, write the expectation of their product:

$$\mathbb{E}[XY] = \int \int xy p_{X,Y}(x, y) dx dy \quad (\text{E.37})$$

$$= \int \int xy p_X(x) p_Y(y) dx dy \quad (\text{E.38})$$

$$= \int x p_X(x) \left( \int y p_Y(y) dy \right) dx \quad (\text{E.39})$$

$$= \left( \int x p_X(x) dx \right) \left( \int y p_Y(y) dy \right) \quad (\text{E.40})$$

$$= \mathbb{E}[X]\mathbb{E}[Y] \quad (\text{E.41})$$

<sup>7</sup>Why is this equivalent?

<sup>8</sup>For the sake of mathematical quibblers: a *non-constant* function of  $X$ .

However, if  $X$  and  $Y$  are uncorrelated, then they can *still* be dependent. To see an extreme example of this, let  $X$  be uniformly distributed on the interval  $[-1, 1]$ . If  $X \leq 0$ , then  $Y = -X$ , while if  $X$  is positive, then  $Y = X$ . You can easily check for yourself that:

- $Y$  is uniformly distributed on  $[0, 1]$
- $\mathbb{E}[XY|X \leq 0] = \int_{-1}^0 -x^2 dx = -1/3$
- $\mathbb{E}[XY|X > 0] = \int_0^1 x^2 dx = +1/3$
- $\mathbb{E}[XY] = 0$  (*hint*: law of total expectation).
- The joint distribution of  $X$  and  $Y$  is not uniform on the rectangle  $[-1, 1] \times [0, 1]$ , as it would be if  $X$  and  $Y$  were independent (Figure E.2).

The only general case when lack of correlation implies independence is when the joint distribution of  $X$  and  $Y$  is a multivariate Gaussian.

## E.5 Exercises

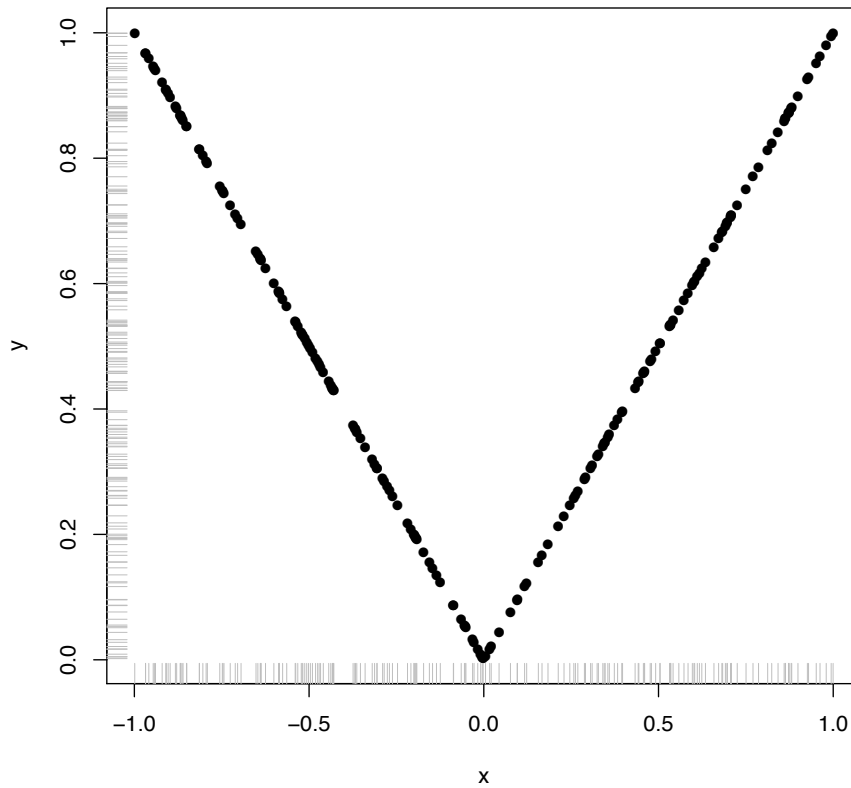
1. Write a function to calculate the density of a multivariate Gaussian with a given mean vector and covariance matrix. Check it against an existing function from one of the packages mentioned in §E.2.4.
2. Write a function to generate multivariate Gaussian random vectors, using `rmnorm`.
3. If  $\vec{X}$  has mean  $\vec{\mu}$  and variance-covariance matrix  $\Sigma$ , and  $\vec{w}$  is a fixed, non-random vector, find the mean and variance of  $w \cdot X$ .
4. If  $\vec{X} \sim \mathcal{MVN}(\vec{\mu}, \Sigma)$ , and  $\mathbf{b}$  and  $\mathbf{c}$  are two non-random matrices, find the covariance matrix of  $\mathbf{b}\vec{X}$  and  $\mathbf{c}\vec{X}$ .
5. One multivariate generalization of the Pareto distribution is defined by

$$\Pr(X_1 \geq x_1, X_2 \geq x_2, \dots, X_p \geq x_p) = \left( \sum_{j=1}^p (x_j/s_j) - p + 1 \right)^{-a} \quad (\text{E.42})$$

when all  $x_j \geq s_j$ .

- (a) Find the joint pdf of all the variables.
- (b) Show that the marginal distribution of each  $X_j$  is a univariate Pareto distribution, and find its parameters.
- (c) Show that the conditional distribution of any  $X_j$  given the other variables is a univariate Pareto distribution, and find its parameters. *Hint*: It is not the same as the marginal distribution.

[[TODO: Move model-comparison, goodness-of-fit stuff to main text]]



```
x <- runif(200,min=-1,max=1)
y <- ifelse(x>0,x,-x)
plot(x,y,pch=16)
rug(x,side=1,col="grey")
rug(y,side=2,col="grey")
```

FIGURE E.2: An example of two random variables which are uncorrelated but strongly dependent. The grey “rug plots” on the axes show the marginal distributions of the samples from  $X$  and  $Y$ .