

Random Vectors

\mathbf{x} is a $p \times 1$ random vector with a pdf probability density function $f(\mathbf{x}): \mathbf{R}^p \rightarrow \mathbf{R}$. Many books write \mathbf{X} for the random vector and $\mathbf{X}=\mathbf{x}$ for the realization of its value.

$$E[\mathbf{X}] = \int \mathbf{x} f(\mathbf{x}) d\mathbf{x} = \boldsymbol{\mu}.$$

Theorem: $E[A\mathbf{x}+\mathbf{b}] = AE[\mathbf{x}]+\mathbf{b}$

Covariance Matrix $E[(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})'] = \text{var}(\mathbf{x}) = \Sigma$ (note the location of transpose)

Theorem: $\Sigma = E[\mathbf{x}\mathbf{x}'] - \boldsymbol{\mu}\boldsymbol{\mu}'$

If \mathbf{y} is a random variable: covariance $C(\mathbf{x}, \mathbf{y}) = E[(\mathbf{x}-\boldsymbol{\mu})(\mathbf{y}-\mathbf{v})']$

Theorem: For constants \mathbf{a} , A , $\text{var}(\mathbf{a}'\mathbf{x}) = \mathbf{a}'\Sigma\mathbf{a}$, $\text{var}(A\mathbf{x}+\mathbf{b}) = A\Sigma A'$, $C(\mathbf{x}, \mathbf{x}) = \Sigma$, $C(\mathbf{x}, \mathbf{y}) = C(\mathbf{y}, \mathbf{x})'$

Theorem: If \mathbf{x} , \mathbf{y} are independent RVs, then $C(\mathbf{x}, \mathbf{y}) = 0$, but not conversely.

Theorem: Let \mathbf{x}, \mathbf{y} have same dimension, then $\text{var}(\mathbf{x}+\mathbf{y}) = \text{var}(\mathbf{x}) + \text{var}(\mathbf{y}) + C(\mathbf{x}, \mathbf{y}) + C(\mathbf{y}, \mathbf{x})$

Normal Random Vectors

The Central Limit Theorem says that if a focal random variable x consists of the sum of many other independent random variables, then the focal random variable will asymptotically have a distribution that is basically of the form e^{-x^2} , which we call "normal" because it is so common.

Normal random variable has pdf $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\left(\frac{x-\mu}{\sigma}\right)^2 / 2} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu) \frac{1}{\sigma^2} (x-\mu) / 2}$

Denote \mathbf{x} $p \times 1$ normal random variable with pdf

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-(\mathbf{x}-\boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

where $\boldsymbol{\mu}$ is the mean vector and Σ is the covariance matrix: $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$.

Bivariate Normal $f(x_1, x_2) = \frac{1}{2\pi\sqrt{\sigma_{11}\sigma_{22} - \sigma_{12}^2}} e^{-\begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix}^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} / 2}$ Note

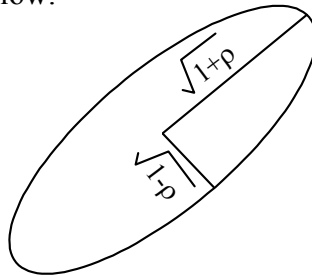
Recall variance σ_{11} is also sometimes written σ_1^2 and by symmetry $\sigma_{12} = \sigma_{21}$. The correlation is $\rho_{12} = \sigma_{12} / \sqrt{\sigma_{11}\sigma_{22}} = \sigma_{12} / (\sigma_1\sigma_2)$.

Theorem: Eigenvalue of Σ^{-1} is reciprocal of eigenvalue of Σ and eigenvectors are identical.

Proof: Let $\Sigma^{-1}\mathbf{x} = \lambda\mathbf{x}$. Then $\mathbf{x} = \lambda\Sigma\mathbf{x}$ or $\Sigma\mathbf{x} = (1/\lambda)\mathbf{x}$.

Contour of constant probability is ellipsoid $(\mathbf{x}-\boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})=c^2$ for some c . This is an ellipse centered at $\boldsymbol{\mu}$ and with axis that point in the directions of the eigenvectors of $\boldsymbol{\Sigma}$ with length $c\sqrt{\lambda_i}$, that is the axes are $\pm c\sqrt{\lambda_i} \mathbf{e}_i$ where λ_i and \mathbf{e}_i are the eigenvalues and eigenvectors of the covariance matrix $\boldsymbol{\Sigma}$.

Suppose that $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$, then eigenvalues are defined by $(1-\lambda)^2-\rho^2=0$ or $\lambda=1\pm\rho$. The eigenvectors are values of (x_1, x_2) such that $-\rho x_1 + \rho x_2 = 0$ and $x_1^2 + x_2^2 = 1$; these are $(1/\sqrt{2}, \pm 1/\sqrt{2})$. If the correlation ρ is positive, then the eigenvector $(1/\sqrt{2}, 1/\sqrt{2})$ is stretched to a length greater than 1, $\sqrt{1+\rho}$, while the eigenvector $(1/\sqrt{2}, -1/\sqrt{2})$ is shrunk to a length less than 1, $\sqrt{1-\rho}$. See the figure below.



Theorems: The moment generating function (mgf) for multivariate normal is

$$\phi_{\mathbf{x}}(\mathbf{t}) = E[e^{\mathbf{t}'\mathbf{x}}] = e^{\mathbf{t}'\boldsymbol{\mu} + \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}}$$

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \Rightarrow x_i \sim N(\mu_i, \sigma_{ii})$$

$$y = \mathbf{a}'\mathbf{x} \Rightarrow y \sim N(\mathbf{a}'\boldsymbol{\mu}, \mathbf{a}'\boldsymbol{\Sigma}\mathbf{a}) \text{ and mgf } \phi_y(\tau) = e^{\tau\mathbf{a}'\boldsymbol{\mu} + \frac{1}{2}\tau^2\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a}}$$

$$x_1|x_2 \sim N(\mu_1 + \sigma_{12}/\sigma_{22}(x_2 - \mu_2), \sigma_{11} - \sigma_{12}^2/\sigma_{22})$$

Theorem $(\mathbf{x}-\boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}) \sim \chi^2_p$

Proof: (note: chi-sq is the sum of the squares of independent normals). Using spectral decomposition, $(\mathbf{x}-\boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}) = (\mathbf{x}-\boldsymbol{\mu})'(\mathbf{P}\boldsymbol{\Lambda}\mathbf{P}')^{-1}(\mathbf{x}-\boldsymbol{\mu}) = (\mathbf{x}-\boldsymbol{\mu})'\mathbf{P}\boldsymbol{\Lambda}^{-1/2}\boldsymbol{\Lambda}^{-1/2}\mathbf{P}'(\mathbf{x}-\boldsymbol{\mu})$. From above $\boldsymbol{\Lambda}^{-1/2}\mathbf{P}'(\mathbf{x}-\boldsymbol{\mu}) \sim N(0, \boldsymbol{\Lambda}^{-1/2}\mathbf{P}'\boldsymbol{\Sigma}\mathbf{P}\boldsymbol{\Lambda}^{-1/2}) = N(0, \boldsymbol{\Lambda}^{-1/2}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^{-1/2}) = N(0, \mathbf{I})$, so quadratic form is the sum of independent squared normal random variables. QED

Normal data matrix $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1' \\ \mathbf{M} \\ \mathbf{x}_n' \end{bmatrix}$ where \mathbf{x}_i is iid $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. This is a $n \times p$ matrix of random variables.

Each row is independent of other rows and identically distributed.

$$P(\mathbf{X}) = (2\pi)^{-np/2} |\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})/2\right)$$

$$= (2\pi)^{-np/2} |\boldsymbol{\Sigma}|^{-n/2} \exp\left[-\text{tr}\left(\boldsymbol{\Sigma}^{-1}\left(\sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})'\right)/2\right)\right]$$

Note: $\mathbf{x}'\mathbf{A}\mathbf{x} = \text{tr}(\mathbf{x}'\mathbf{A}\mathbf{x}) = \text{tr}(\mathbf{A}\mathbf{x}\mathbf{x}')$

Aside: On Union Intersection Tests

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{I}) \Rightarrow y = \mathbf{a}'\mathbf{x} \sim N(\mathbf{a}'\boldsymbol{\mu}, \mathbf{a}'\mathbf{a})$$

$$H_0: \boldsymbol{\mu} = \mathbf{0} \Leftrightarrow y_a = \mathbf{a}'\mathbf{x} \sim N(0, \mathbf{a}'\mathbf{a}) \text{ for all } \mathbf{a}$$

$$H_{0a}: \mathbf{a}'\boldsymbol{\mu} = 0$$

$H_0 = \cap H_{0a}$ note: if you find one \mathbf{a} that violates H_{0a} then H_0 cannot be true.

Let's test H_{0a} using $z_a = y_a / \sqrt{a'a}$. The rejection region is $R_a = \{z_a | z_a^2 > c^2\}$. What about H_0 ?

$R = \cup R_a$. So H_0 is accepted if and only if $z_a^2 < c^2$ for all a . The worst case scenario is $\max_a z_a^2$.

So, if $\max_a z_a^2 < c^2$ then this will be true for all a .

Suppose that we have independent draws of a random vector. Let \mathbf{x}_j be the j th draw. Define

$\mathbf{y} = \mathbf{a}'\mathbf{x}_j$. Then we know that $\bar{y} = \mathbf{a}'\bar{\mathbf{x}}$ and $s_y^2 = \mathbf{a}'\mathbf{S}\mathbf{a}$. Compute $t^2 = \frac{n(\mathbf{a}'(\bar{\mathbf{x}} - \boldsymbol{\mu}))^2}{\mathbf{a}'\mathbf{S}\mathbf{a}}$. Following the

Union Intersection test procedure we would like to find the value of \mathbf{a} that is the worst case scenario. The Cauchy-Schwartz inequality helps here: $(\mathbf{x}'\mathbf{y})^2 \leq (\mathbf{x}'\mathbf{x})(\mathbf{y}'\mathbf{y})$ (this is a consequence of $\mathbf{x}'\mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos \theta$): $(\mathbf{a}'(\bar{\mathbf{x}} - \boldsymbol{\mu}))^2 \leq (\mathbf{a}'\mathbf{S}\mathbf{a}) ((\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}))$ and can only "=" if $\mathbf{a} = \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu})$.

Taking this worst case scenario, the $\max_a t^2 = n (\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu})$.

Theorem: The interval $\mathbf{a}'\bar{\mathbf{x}} \pm \sqrt{\frac{p(n-1)}{n-p} \frac{\mathbf{a}'\mathbf{S}\mathbf{a}}{n}} F_{p,n-p}(\alpha)$ will contain $\mathbf{a}'\boldsymbol{\mu}$ a fraction $1-\alpha\%$ of the time, simultaneously for all possible \mathbf{a} .

Comparison of Traditional and Simultaneous Confidence Intervals

Suppose that you had $H_{0i}: \mu_i = 0$ for $i=1,2,\dots,p$. If you ignore the fact that there are several simultaneous test, you would do this one variable at a time, computing confidence intervals:

$$\bar{x}_i \pm \sqrt{s_{ii} / n} \cdot t_{n-1}(\alpha/2).$$

As we have seen before, the confidence for these as a whole is not $1-\alpha\%$, but rather $(1-\alpha)^p$: for 6 variables $0.95^6 = 0.75$. Hence the rectangular region sketched out by these intervals is really only a 75% confidence region. If you had 13 variables, then this region will capture the truth in all dimensions only 50% the time. We have a false sense of high accuracy.

The above Union Intersection test would sequentially set $\mathbf{a}' = (0, \dots, 0, 1, 0, \dots, 0)$ where the 1 is in the i th entry and then calculate the intervals

$$\bar{x}_i \pm \sqrt{s_{ii} / n} \cdot \sqrt{\frac{p(n-1)}{n-p}} F_{p,n-p}(\alpha).$$

These simultaneous intervals will be much wider than the above, but we can then say that there is a 95% confidence that all variables will be covered by the combined rectangle. These intervals are the "shadows of the 95% confidence ellipse in a p -dimensional space.

How much wider are these simultaneous intervals? It depends on n and p . As you can see the intervals are much wider, making it very difficult to say with high confidence, "All elements of my theory are true."

$$\sqrt{\frac{p(n-1)}{n-p}} F_{p,n-p}(\alpha)$$

n	$t_{n-1}(0.025)$	$p=4$	$p=10$
15	2.145	4.14	11.52
25	2.064	3.60	6.39
50	2.010	3.31	5.05
100	1.970	3.19	4.61
∞	1.960	3.08	4.28

Generalization of t-test to T²-test

Neither the traditional nor the simultaneous interval tests take into account that the variables may be correlated with one another. The 95% confidence region should not be a rectangle, but rather an ellipse. How should we handle this? This is not that complicated.

In the single normal variable case, we test using $t \equiv \frac{\bar{x} - \mu}{s / \sqrt{n}}$. When we have several variables that we want to combined without having ‘s canceling +’s, we use the Hotelling T²-distribution of the variable $t^2 \equiv \frac{(\bar{x} - \mu)^2}{s^2 / n}$. For p-variate normal vector case, the equivalent statistic is

$$T_{p,n-1}^2 \equiv \frac{(\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu})}{1/n} = n(\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}).$$

This statistic has normals^{squared} on the top (the x’x terms) and normals^{squared} in the bottom (since S is composed from squared normals). That is, it is the ratio of χ^2 s and hence has the F-distribution:

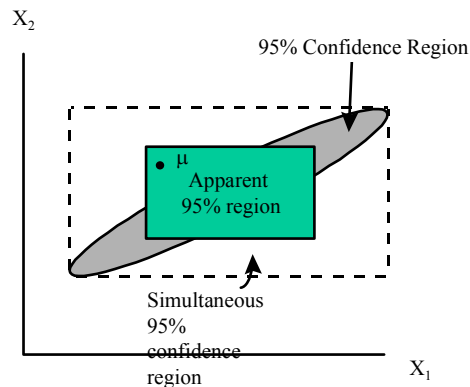
$$T_{p,n-1}^2 \sim \frac{(n-1)p}{n-p} F_{p,n-p}.$$

Thus if we had a null hypothesis that the p-variate variable \mathbf{x} had mean $\boldsymbol{\mu}$, then we would construct the above T² statistic and see if it exceeded the critical value found in an F-distribution table. This will tell us whether our theoretical value $\boldsymbol{\mu}$ is covered by the confidence ellipsoid 1- α % of the time in repeated samples. The 1- α % confidence ellipsoid has axes determined by

$n(\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) \leq c^2 = \frac{(n-1)p}{n-p} F_{p,n-p}(\alpha)$. That is they are determined by starting at $\bar{\mathbf{x}}$ and

going $\pm \sqrt{\lambda_i} c / \sqrt{n} = \pm \sqrt{\lambda_i} \sqrt{\frac{(n-1)p}{n(n-p)} F_{p,n-p}(\alpha)}$ units along the eigenvectors \mathbf{e}_i . This is better

than doing p separate t-tests of each variable, since it uses all the information in S, including the fact that some variables are highly correlated.



In summary, do not claim when you study p variables and all of them fit your theory that you are 95% confident in your theory. Apparent confidence is not real confidence. On the other hand, even if one-at-a-time you cannot reject the null, you still may be able to with 95% confidence state, “There are some elements of this theory that must be true, I just cannot tell you which ones.” In the above graph, the true $\boldsymbol{\mu}$ is inside the apparent 95% confidence interval, so you apparently cannot reject any element μ_i , but $\boldsymbol{\mu}$ is outside the 95% confidence ellipsoid, so you reject $\boldsymbol{\mu}$ as a whole.