

Varianza muestral generalizada y Normal multivariada

martes, 8 de febrero de 2022 9:06 a.m.

$$S_n = \frac{1}{n} \sum_j (x_{ji} - \bar{x}_i)^2$$

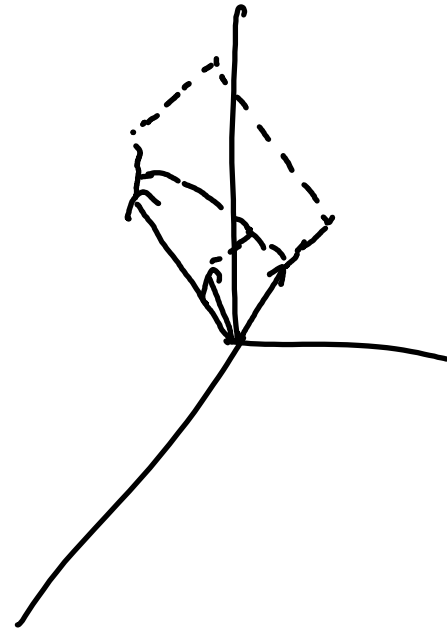
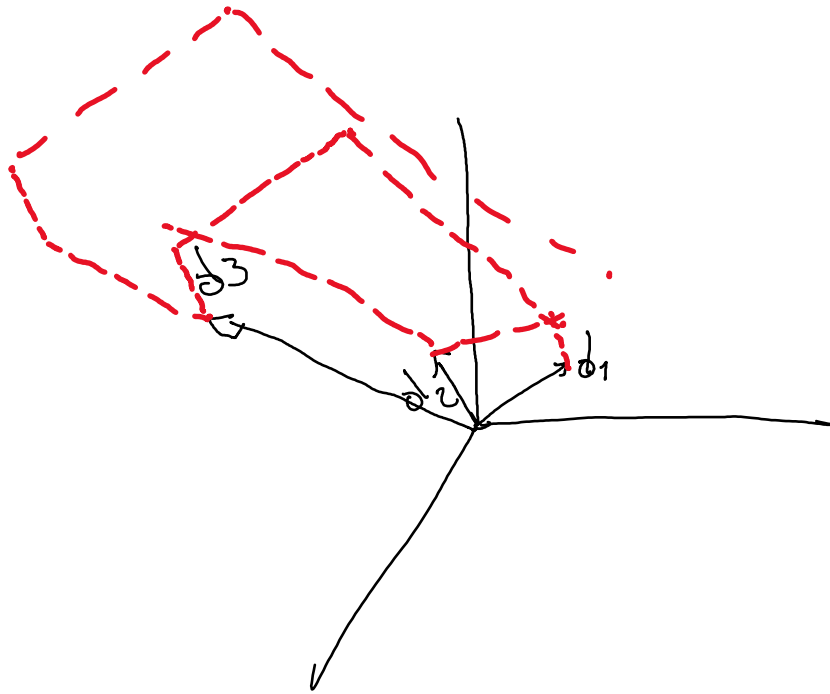
$$S = \frac{1}{n-1} \sum_j (x_{ij} - \bar{x}_i)^2$$

$$S = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

Matriz de varianzas y
covarianzas muestrales

Insegura

Varianza muestral generalizada = $|S|$



La Varianza generalizada es cero cuando y solo cuando al menos un vector de variación yace en el hiperplano formado por todas las combinaciones lineales de los restantes vectores. Es decir cuando al menos una columna de la matriz S es linealmente dependiente de otra.

P

$$n \leq p$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$S = \begin{bmatrix} s_{11} & \dots & s_{1n} \\ \vdots & \ddots & \vdots \\ s_{n1} & \dots & s_{nn} \end{bmatrix}$$

Varianza generalizada determinada por $|R|$

- La varianza muestral generalizada se ve afectada íntegramente por la variabilidad de las mediciones en una sola variable.
- Por ejemplo, si S_{ii} es muy grande (o pequeña). Geométricamente el vector de desviación correspondiente $d_i = (y_i - \bar{x}_i \mathbf{1})$ será muy largo (o muy corto) y

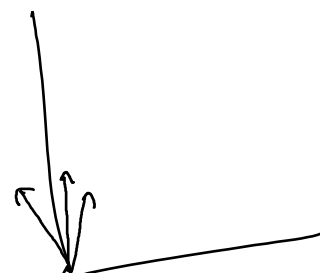
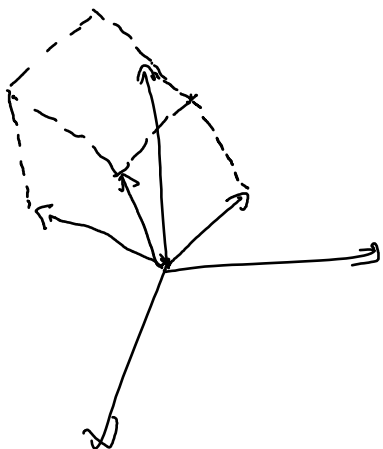
afectar el volumen.

- Entonces, es útil escalar todos los vectores de desviación para que tengan la misma longitud.
- Equivale a reemplazar la observación original x_{ji} por $(x_{ji} - \bar{x}_i) / \sqrt{s_{ii}}$

Entonces, la matriz de varianzas y covarianzas $R \rightarrow$ La matriz de correlación

$|R|$ = Varianza muestral generalizada de los n variables estandarizadas

variance



Varianza total de la muestra $= S_{11} + S_{22} + \dots + S_{pp}$

Normal multivariada

Recordemos:

La PDF de una v.a. normal con media μ y var σ^2 es:

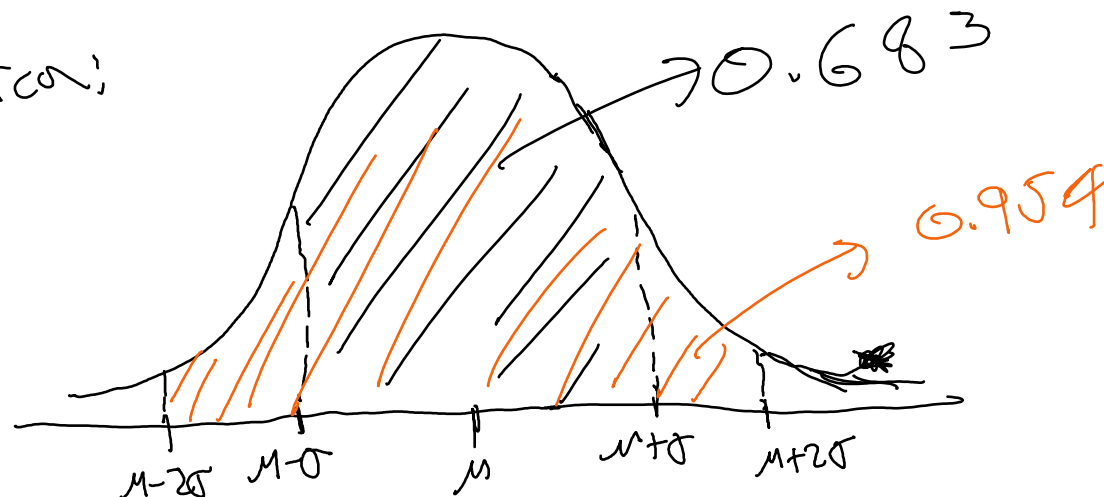
$$\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$x \in \mathbb{R}$$

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$X \sim N(\mu, \sigma^2)$

Ley empírica:



Definimos la distr. normal con media μ y var $\sigma^2 \rightarrow N(\mu, \sigma^2)$

X es normal con media μ y var σ^2

$$X \sim N(\mu, \sigma^2)$$

$$\frac{(x-\mu)^2}{\sigma^2}$$

→ Distancia estadística entre x y μ

$$= (x-\mu)(\sigma^2)^{-1}(x-\mu)$$

Generalizando \vec{x} vector de observaciones de X
 ↪ una fila de la matriz X

distancia estadística x a μ

$$(x-\mu)' \Sigma^{-1} (x-\mu)$$

Entonces podemos generalizar y obtener la PDF de la normal multivariada.

$$= (x-\mu)' \Sigma^{-1} (x-\mu) / 2$$

$$f(x) = \frac{1}{\dots} \cdot e$$

$$J \sim \underbrace{(2\pi)^{p/2} |\Sigma|^{1/2}}$$

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$$

Es un vector aleatorio $X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$ v.a.

Denotamos la normal multivariada como

$$X \sim N_p(\mu, \Sigma)$$

vector matriz.

Obs:

$$\text{Si } p=2 \quad \gamma \quad \rho_{12}=0$$

$$(\text{cov}(X_1, X_2) = 0)$$


 X_1, X_2 indep.

Teorema:

Si Σ es positiva (Σ^{-1} existe)

entonces si e es un vector propio de Σ con valor propio asociado λ , e es un vect propio de Σ^{-1} con val propio asociado $\frac{1}{\lambda}$. Además

Σ^{-1} es positiva:

chi-cuadrado.

Obs:

$$1) (X - \mu)' \Sigma^{-1} (X - \mu) \leq \chi^2_p(\alpha)$$

con probabilidad $1-\alpha$

gu.

Proposición 1)

Si X vector aleatorio normal multivariado.

$$X \sim N_p(\mu, \Sigma) \text{ Entonces}$$

$$\text{Si } a \in \mathbb{R}^p, \quad a'X = a_1X_1 + a_2X_2 + \dots + a_pX_p \\ \sim N(a'\mu, a'\Sigma a)$$

Prop 2) Si $a'X \sim N(a'\mu, a'\Sigma a), \quad \forall a \in \mathbb{R}^p$
entonces $X \sim N_p(\mu, \Sigma)$

Prop 2) $X \sim N_p(\mu, \Sigma)$

Sea $A \in \mathbb{R}^{q \times p}$

$$AX \sim N_q(A\mu, A\Sigma A')$$

Prop 4) $X \sim N_p(\mu, \Sigma) \quad d \in \mathbb{R}^p$

$$X + d \sim N_p(\mu + d, \Sigma)$$

Prop 5) Todas las particiones de vect. aleatorios normales resultan en vect. aleat normales.

$$X \sim N_p(\mu, \Sigma)$$

$$X = \begin{pmatrix} X_{(1)} \\ X_{(2)} \end{pmatrix} \begin{matrix} \} q \\ \} p-q \end{matrix}$$

$$\mu = \begin{pmatrix} \mu_{(1)} \\ \vdots \\ \mu_{(2)} \end{pmatrix}$$

$$\mu_{(1)} \leq \mu_{(2)}$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

$$X_{(1)} \sim N_q(\mu_{(1)}, \Sigma_{11})$$

$$X_{(2)} \sim N_{p-q}(\mu_{(2)}, \Sigma_{22})$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Prop 6) Si $X_{(1)}^{q_1 \times 1}, X_{(2)}^{q_2 \times 1}$ normales multivariadas e indep.
 entonces $\underbrace{\text{cov}(X_{(1)}, X_{(2)})}_{\text{matriz } q_1 \times q_2} = 0 \rightarrow \text{Matriz } 0$

$$2) \text{ Si } X = \begin{pmatrix} X_{(1)} \\ X_{(2)} \end{pmatrix} \sim N_{q_1+q_2} \left(\begin{pmatrix} \mu_{(1)} \\ \mu_{(2)} \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

$X_{(1)}$ y $X_{(2)}$ indep si y solo si $\Sigma_{12} = 0$

Covarianza 0 \Rightarrow Sí implica independencia

Si HA y normalizado,

3) $X_{(1)}, X_{(2)}$ normales indep si y solo si

$$X = \begin{pmatrix} X_{(1)} \\ X_{(2)} \end{pmatrix} \sim N_{q_1+q_2} \left(\begin{pmatrix} \mu_{(1)} \\ \mu_{(2)} \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix} \right)$$

Prop 7) Sea $X \sim N_p(\mu, \Sigma)$, $|\Sigma| > 0$

Entonces

$$1) (X - \mu)' (\Sigma^{-1}) (X - \mu) \sim \chi_p^2$$

la distancia de los valores a la media tiene
trib χ^2
La probabilidad del evento.

$$\{ (X-\mu)'(\Sigma^{-1})(X-\mu) \leq \chi_p^2(\alpha) \} \text{ es } 1-\alpha.$$

Muestras de la normal multivariada.

Suponga que X_1, \dots, X_n son una muestra aleatoria de una población normal multivariada con media

μ y cov Σ

Densidad conjunta de X_1, \dots, X_n

$$\prod_{j=1}^n \left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}[(X_j - \mu)' \Sigma^{-1} (X_j - \mu)]} \right)$$

función de verosimilitud (función de μ y Σ)

1
 dadas las observaciones $X_1, \dots, X_n : L(\mu, \Sigma)$

Método de máxima-verosimilitud:
 Utilizar como estimaciones de parámetros poblacionales desconocidos los valores que maximizan $L(\mu, \Sigma)$

Los valores que "mejor" explican los datos

Teo: Sean X_1, \dots, X_n una muestra aleatoria de una población normal con media μ y cov. Σ .

Entonces los MLE son:

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

$$\hat{\Sigma} = \frac{1}{n} \left(\sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})' \right)$$

$$= S = \frac{n-1}{n} S$$

Obs: $\hat{\mu}$, $\hat{\Sigma}$ son VECTORES
 ALEATORIOS antes de
 tomar la data

$$2) L(\hat{\mu}, \hat{\Sigma}) = \frac{1}{(2\pi)^{np/2}} e^{-n/2} \times \frac{1}{|\hat{\Sigma}|^{n/2}}$$

↓
n-1 S

$$\text{de } (|S|)^{-n/2}$$

la varianza generalizada determina el

máximo de la función de verosimilitud.