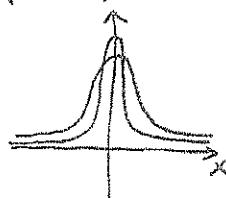# 1 Review of probability distributions

Building block: $X \sim N(\mu, \sigma^2)$  normal distribution

$x \in (-\infty, \infty)$  $\qquad f_X(x) = \dfrac{1}{\sqrt{2\pi}\cdot\sigma}\, e^{-(x-\mu)^2/2\sigma^2}$  density (p.d.f.)



"Bell curve"
for $\mu = 0$
As $\sigma \uparrow$, hump $\downarrow$

$Z$ - standard normal distribution   $Z \sim N(0,1)$   $\boxed{Z = \dfrac{X-\mu}{\sigma} \iff X = \mu + \sigma Z}$

$E[Z] = \displaystyle\int_{-\infty}^{\infty} z \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}\, dz = 0$   (since integrand is an odd function)

explanation

$P(Z \in [a,b]) = \mathbb{P}\left(\dfrac{X-\mu}{\sigma} \in [a,b]\right) = \mathbb{P}\left(X \in [a\sigma+\mu, b\sigma+\mu]\right) =$

$\displaystyle = \int_{a\sigma+\mu}^{b\sigma+\mu} \frac{1}{\sqrt{2\pi}\cdot\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}\, dx = \int_{a}^{b} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}\, dz$

$z = \frac{x-\mu}{\sigma}$
change of variable

$E[Z^2] = 1$.  Why? Start with   $1 = \displaystyle\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}\, dz$ , integrate by parts:

$u = e^{-\frac{z^2}{2}}$   $\qquad du = -z e^{-\frac{z^2}{2}}\, dz$
$dv = dz$   $\qquad v = z$

$\Rightarrow 1 = \dfrac{1}{\sqrt{2\pi}} z e^{-\frac{z^2}{2}}\Big|_{-\infty}^{\infty} - \dfrac{1}{\sqrt{2\pi}}\displaystyle\int_{-\infty}^{\infty} z(-z) e^{-\frac{z^2}{2}}\, dz =$

$\displaystyle = 0 + \int_{-\infty}^{\infty} z^2 \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}\, dz = E[Z^2]$

So, $\mathrm{Var}(Z) = E[Z^2] - (E[Z])^2 = 1 - 0 = 1$

Hence,  $E[X] = \mu + \sigma E[Z] = \mu$
$E[X^2] = E[\mu^2 + 2\mu\sigma Z + \sigma^2 Z^2] = \mu^2 + \sigma^2$   (using linearity of expectation)
$\mathrm{Var}(X) = E[X^2] - (E[X])^2 = \sigma^2$

__FACT1__ If $X_1, X_2, \ldots, X_n$ are independent random variables,

$X_i \sim N(\mu_i, \sigma_i^2)$, $i = 1, 2, \ldots, n$, then $\sum_{i=1}^{n} X_i \sim N\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$

__FACT2__ (__Law of large numbers__)

If $X_1, X_2, \ldots, X_n$ is an i.i.d. sample with $E[X_1] < \infty$, then the

sample average $\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i \longrightarrow E[X_1]$, i.e converges to

common expectation in _probability_, i.e.

$$\forall \varepsilon > 0 \quad \mathbb{P}\left(|\overline{X}_n - E[X_1]| > \varepsilon\right) \longrightarrow 0 \quad \text{as } n \to \infty.$$

__FACT3__ (__Central limit theorem__)

If $X_1, X_2, \ldots, X_n$ is an i.i.d. sample from some continuous distribution

such that $E[X_1] < \infty$ and $\sigma^2 = \text{Var}(X_1) < \infty$, then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} (X_i - E[X_1]) = \sqrt{n}\left(\overline{X}_n - E[X_1]\right) \xrightarrow{d} N(0, \sigma^2), \text{ i.e.}$$

converges in distribution to a $N(0, \sigma^2)$, where convergence in

distribution means that for every interval $[a, b]$:

$$\mathbb{P}\left(\sqrt{n}\left(\overline{X}_n - E[X_1]\right) \in (a, b)\right) \longrightarrow \int_a^b \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{x^2}{2\sigma^2}} dx$$

| Gamma distribution |

has two parameters $\overset{\text{shape}}{\alpha > 0}, \overset{\text{scale}}{\beta > 0}$ Before we define it, let's recall
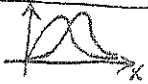
gamma function $\quad \Gamma(\alpha) = \int_0^\infty x^{\alpha - 1} e^{-x} dx$

Divide both sides by $\Gamma(\alpha)$ to get: $\quad 1 = \int_0^\infty \frac{1}{\Gamma(\alpha)} x^{\alpha - 1} e^{-x} dx$, i.e

(after substitution $x = \beta y$) $\quad 1 = \int_0^\infty \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha - 1} e^{-\beta y} dy$

Define $\boxed{f(x\mid\alpha,\beta) = \dfrac{\beta^{\alpha}}{\Gamma(\alpha)}\, x^{\alpha-1} e^{-\beta x}\text{, if } x\geq 0\text{; and } 0 \text{ otherwise}}$

This is the p.d.f. (since it's nonnegative and it integrates to 1). 

Let $X\sim \Gamma(\alpha,\beta)$, i.e. $X$ is a random variable with p.d.f. $f_X(x) = f(x\mid\alpha,\beta)$

Properties of gamma function $\Gamma(\alpha)$

$\Gamma(\alpha) = \displaystyle\int_0^{\infty} x^{\alpha-1} e^{-x}\,dx = x^{\alpha-1}(-e^{-x})\Big|_0^{\infty} - \int_0^{\infty}(-e^{-x})(\alpha-1)x^{\alpha-2}\,dx = (\alpha-1)\int_0^{\infty} x^{(\alpha-1)-1} e^{-x}\,dx$

$\qquad u = x^{\alpha-1} \qquad\quad dv = e^{-x}dx$

$\qquad du = (\alpha-1)x^{\alpha-2}dx \qquad v = -e^{-x}$

i.e. $\Gamma(\alpha) = (\alpha-1)\Gamma(\alpha-1)$. Since $\Gamma(1) = \displaystyle\int_0^{\infty} e^{-x}dx = 1$, iterating the last identity, we get $\Gamma(n) = (n-1)!$

$k^{th}$-moment of Gamma distribution:

$E[X^k] = \displaystyle\int_0^{\infty} x^k\, \frac{\beta^{\alpha}}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}\,dx = \frac{\beta^{\alpha}}{\Gamma(\alpha)}\int_0^{\infty} x^{(\alpha+k)-1} e^{-\beta x}\,dx =$

$= \dfrac{\beta^{\alpha}}{\Gamma(\alpha)}\cdot\dfrac{\Gamma(\alpha+k)}{\beta^{\alpha+k}}\displaystyle\int_0^{\infty}\boxed{\dfrac{\beta^{\alpha+k}}{\Gamma(\alpha+k)} x^{(\alpha+k)-1} e^{-\beta x}}dx = \dfrac{\beta^{\alpha}}{\Gamma(\alpha)}\cdot\dfrac{\Gamma(\alpha+k)}{\beta^{\alpha+k}} =$

p.d.f. of $\Gamma(\alpha+k,\beta)$ so, it integrates to 1

$= \dfrac{\Gamma(\alpha+k)}{\Gamma(\alpha)}\cdot\dfrac{1}{\beta^k} = \dfrac{(\alpha+k-1)(\alpha+k-2)\cdots\alpha}{\beta^k}$

(using the property of Gamma funct.)

So, $\boxed{E[X] = \dfrac{\alpha}{\beta}\ ,\ E[X^2] = \dfrac{(\alpha+1)\alpha}{\beta^2}\ ,\ Var(X) = \dfrac{\alpha}{\beta^2}\ \text{for } X\sim\Gamma(\alpha,\beta)}$

FACT4. If $X_i \sim \Gamma(\alpha_i,\beta)$, $i=1,2,\ldots,n$, are independent random variables, then $\displaystyle\sum_{i=1}^{n} X_i \sim \Gamma\left(\sum_{i=1}^{n}\alpha_i,\beta\right)$.

proof of FACT 4 : First, we find a moment generating function (m.g.f.) of $X \sim \Gamma(\alpha, \beta)$ :

$$E[e^{tX}] = \int_0^\infty e^{tx} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} dx = \int_0^\infty \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-(\beta-t)x} dx =$$

$$= \frac{\beta^\alpha}{(\beta-t)^\alpha} \int_0^\infty \boxed{\frac{(\beta-t)^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-(\beta-t)x}} dx = \left(\frac{\beta}{\beta-t}\right)^\alpha$$

p.d.f. of $\Gamma(\alpha, \beta-t)$
so it integrates to 1

The m.g.f. of $\sum_{i=1}^n X_i$ , $X_i \sim \Gamma(\alpha_i, \beta)$ is :

$$E[e^{t\sum_{i=1}^n X_i}] = E[\prod_{i=1}^n e^{tX_i}] = \prod_{i=1}^n E[e^{tX_i}] = \prod_{i=1}^n \left(\frac{\beta}{\beta-t}\right)^{\alpha_i} = \left(\frac{\beta}{\beta-t}\right)^{\sum_{i=1}^n \alpha_i}$$
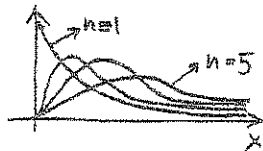
independence of $X_i$

which is again the m.g.f. of a Gamma distribution $\Gamma(\sum_{i=1}^n \alpha_i, \beta)$ $\square$

$\boxed{\chi_n^2\text{-distribution}}$

($n$ degrees of freedom) is the distribution of $\sum_{i=1}^n X_i^2$, where $X_i \sim N(0,1), i=1,...,n$
$X_i$ independent

Relationship with Gamma
$$\chi_n^2 \equiv \Gamma\left(\frac{n}{2}, \frac{1}{2}\right)$$



As $n\uparrow$, flatter w/ hump

proof: Let $X \sim N(0,1)$. Then the cumulative distribution function of $X^2$ is

$$\mathbb{P}(X^2 \leq x) = \mathbb{P}(-\sqrt{x} \leq X \leq \sqrt{x}) = \int_{-\sqrt{x}}^{\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt \text{ ; so, the p.d.f. is}$$

(c.d.f)

$$f_{\chi^2}(x) = \frac{d}{dx}\left(\mathbb{P}(X^2 \leq x)\right) = \frac{d}{dx} \int_{-\sqrt{x}}^{\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = \frac{1}{\sqrt{2\pi}} e^{-\frac{(\sqrt{x})^2}{2}} \frac{d}{dx}(\sqrt{x}) - \frac{1}{\sqrt{2\pi}} e^{-\frac{(-\sqrt{x})^2}{2}} \frac{d}{dx}(-\sqrt{x})$$

i.e. $f_{\chi^2}(x) = \frac{1}{\sqrt{2\pi}} x^{\frac{1}{2}-1} e^{-\frac{x}{2}}$ (after some algebraic manipulation)

Since $\Gamma(1/2) = \sqrt{\pi}$ (why?) , $X^2 \sim \Gamma(1/2, 1/2)$.

Now, using FACT 4 , $\sum_{i=1}^n X_i^2 \sim \Gamma(n/2, 1/2)$ where $X_i \sim N(0,1)$ i.i.d. $i=1,...,n$

So, $\chi_n^2 \equiv \Gamma(n/2, 1/2)$

-4-

## Fisher F-distribution

Let $X \sim \chi_k^2 \equiv \Gamma(k/2, 1/2)$, $Y \sim \chi_m^2 \equiv \Gamma(m/2, 1/2)$; $X, Y$ independent rand. var's.
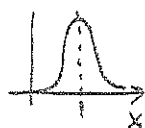
Let $Z \equiv \dfrac{X/k}{Y/m}$. Then, $Z$ is said to have a Fisher distribution with degrees of freedom $k$ and $m$, and is denoted $Z \sim F_{k,m}$.

**Observation 1.** Since $X \sim \chi_k^2$; then $X \equiv \sum_{i=1}^{k} X_i^2$, $X_i \sim N(0,1)$ i.i.d. By Law of Large numbers (FACT 2): $\dfrac{1}{k} \sum_{i=1}^{k} X_i^2 \to E[X_i^2] = 1$, as $k \to \infty$.

So, as $k, m \to \infty$, $X/k \to 1$, $Y/m \to 1$ $\Rightarrow$ $Z$ will concentrate around $1$.

p.d.f. of $Z$ is



As $k, m \to \infty$, hump $\uparrow$ and $\to \leftarrow$ (narrower).

**Observation 2** $F_{k,m}(c, \infty) = F_{m,k}(0, \frac{1}{c})$. Why? $F_{k,m}(c, \infty) = \mathbb{P}\left(\dfrac{X/k}{Y/m} \geq c\right) = \mathbb{P}\left(\dfrac{Y/m}{X/k} \leq \dfrac{1}{c}\right) = F_{m,k}(0, \frac{1}{c})$

What is the p.d.f. of $Z \sim F_{k,m}$? First, we compute the p.d.f. of $\dfrac{X}{Y} = \dfrac{k}{m} Z$.

$$f_X(x) = \frac{(1/2)^{k/2}}{\Gamma(k/2)} x^{k/2 - 1} e^{-x/2}, \quad x \geq 0$$

$$f_Y(y) = \frac{(1/2)^{m/2}}{\Gamma(m/2)} y^{m/2 - 1} e^{-y/2}, \quad y \geq 0$$

$\Bigg\}$ the p.d.f.'s of $X$, $Y$, respectively.

To find the p.d.f. of $X/Y$, first write the c.d.f. Since $X, Y > 0$, $X/Y > 0$, so for $t \geq 0$: $\mathbb{P}(X/Y \leq t) = \mathbb{P}(X \leq tY) = \int_0^\infty \int_0^{ty} f_{X,Y}(x,y) \, dx \, dy$, where

$f_{X,Y}(x,y)$ is the joint density of $X$ and $Y$.

**But**, $X, Y$ are independent, so $f_{X,Y}(x,y) = f_X(x) f_Y(y)$. Hence,

$$\mathbb{P}(X/Y \leq t) = \int_0^\infty \int_0^{ty} f_X(x) f_Y(y) \, dx \, dy, \quad \text{and}$$

$$f_{X/Y}(t) = \frac{d}{dt} \mathbb{P}(X/Y \leq t) = \frac{d}{dt} \int_0^\infty \int_0^{ty} f_X(x) f_Y(y) \, dx \, dy = \int_0^\infty f_X(ty) \cdot f_Y(y) y \, dy =$$

$$= \int_0^\infty \frac{(1/2)^{k/2}}{\Gamma(k/2)} (ty)^{k/2 - 1} e^{-(ty)/2} \frac{(1/2)^{m/2}}{\Gamma(m/2)} y^{m/2 - 1} e^{-y/2} y \, dy =$$

$$= \frac{(1/2)^{(k+m)/2}}{\Gamma(k/2) \Gamma(m/2)} t^{k/2 - 1} \int_0^\infty y^{(k+m)/2 - 1} e^{-((t+1)y)/2} \, dy =$$

-5-

$$= \frac{(1/2)^{(k+m)/2}}{\Gamma(k/2)\Gamma(m/2)} \cdot t^{k/2-1} \cdot \frac{\Gamma\left(\frac{k+m}{2}\right)}{\left(\frac{t+1}{2}\right)^{(k+m)/2}} \int_0^\infty \frac{\left(\frac{t+1}{2}\right)^{(k+m)/2}}{\Gamma\left(\frac{k+m}{2}\right)} \cdot y^{(k+m)/2-1} e^{-\left(\frac{t+1}{2}\right)y} \, dy$$

p.d.f. of $\Gamma\left(\frac{k+m}{2}, \frac{t+1}{2}\right)$ So integral $=1$

Hence, $f_{X/Y}(t) = \dfrac{\Gamma\left(\frac{k+m}{2}\right)}{\Gamma\left(\frac{k}{2}\right)\Gamma\left(\frac{m}{2}\right)} \cdot t^{\frac{k}{2}-1} \cdot (1+t)^{-\frac{k+m}{2}}$

Since $\mathbb{P}(Z \leq t) = \mathbb{P}\left(\dfrac{X}{Y} \leq \dfrac{kt}{m}\right) \Rightarrow f_Z(t) = \dfrac{d}{dt}\mathbb{P}(Z \leq t) = f_{X/Y}\left(\dfrac{kt}{m}\right) \cdot \dfrac{k}{m}$

Finally, the p.d.f. of $Z \sim F_{k,m}$ is:

$$f_Z(t) = \frac{\Gamma\left(\frac{k+m}{2}\right)}{\Gamma(k/2)\Gamma(m/2)} \cdot k^{\frac{k}{2}} m^{\frac{m}{2}} t^{\frac{k}{2}-1} (m+kt)^{-\frac{k+m}{2}} = f_{k,m}(t)$$

---

Student $t_n$-distribution

distribution of a random variable $T \equiv \dfrac{X_1}{\sqrt{\left(\sum_{i=1}^n Y_i^2\right)/n}}$ , where $X_1, Y_1, \ldots, Y_n$ are all i.i.d. $N(0,1)$

As $n \uparrow$, hump goes up, and it approaches the standard normal distribution

What is the p.d.f. of $T$?   $\mathbb{P}(-t \leq T \leq t) = \mathbb{P}(T^2 \leq t^2) = \mathbb{P}\left(\dfrac{X_1^2}{\frac{1}{n}\sum_{i=1}^n Y_i^2} \leq t^2\right)$

$$\int_{-t}^t f_T(x)\,dx \quad\quad = \quad\quad \int_0^{t^2} f_{1,n}(x)\,dx$$

Since $\dfrac{X_1^2}{\frac{1}{n}\sum_{i=1}^n Y_i^2} \sim F$

Taking $\frac{d}{dt}$ of both sides:   $f_T(t) + f_T(-t) = f_{1,n}(t^2) \cdot 2t$

$t_n$-distribution is symmetric (because the numerator has symmetric distribution $N(0,1)$); hence,

$$f_T(t) = f_T(-t) \text{ and thus}$$

$$f_T(t) = f_{1,n}(t^2) \cdot t, \text{ i.e. } f_T(t) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma(1/2)\Gamma(n/2)} \frac{1}{\sqrt{n}}\left(1+\frac{t^2}{n}\right)^{-\frac{n+1}{2}}$$

# Sample mean and sample variance

Given a sample $y_1, y_2, \ldots, y_n$ (independent $y_i$'s) from an (unknown) distribution,

the <u>sample mean</u> is : $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

the <u>sample variance</u> is $s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (y_i - \bar{y})^2$

<u>Theorem</u> If $y_1, y_2, \ldots, y_n$ are independent sample from $N(\mu, \sigma^2)$,

then $\quad \bar{y} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

$$(n-1)\frac{s^2}{\sigma^2} \sim \chi^2_{n-1} \quad \longleftarrow \text{Comment: "Loss" of one degree of freedom is}$$
due to the linear constraint $\sum_{i=1}^{n}(y_i - \bar{y}) = 0$
in estimating $\mu$ by $\bar{y}$.

$\bar{y}$ and $s^2$ are independent

proof: This will be a special case of a much more general result on multiple linear regression (next time). It is also proved in Section 3 (Theorem 1)

## Multivariate distributions

If $X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{pmatrix}$ is a $k \times 1$ random vector, then its expectation is $E[X] = \begin{pmatrix} E[X_1] \\ E[X_2] \\ \vdots \\ E[X_k] \end{pmatrix}$

and its covariance matrix is $Cov(X) = E\left[(X - E[X])(X - E[X])'\right]$   ($k \times k$ matrix)

Covariance matrix is always symmetric, i.e. $(Cov(X))' = Cov(X)$ and nonnegative definite, i.e.

for any $k \times 1$ non-random vector $a$, we have:

$$a' Cov(X) a = E\left[a'(X - E[X])(X - E[X])'a\right] = E\left[\|a'(X - E[X])\|^2\right] \geq 0$$

$\underbrace{\text{using } \|v\|^2 = v'v}$

How does the covariance of $X$ change once we multiply $X$ by a non-random $n \times k$ matrix $A$?
Let $Y = AX$ ($Y$ is $n \times 1$ vector). The covariance of $Y$ will be an $n \times n$ matrix:

$$Cov(Y) = E\left[(Y - E[Y])(Y - E[Y])'\right] = E\left[(AX - E[AX])(AX - E[AX])'\right] =$$

$$= E\left[A(X - E[X])(A(X - E[X]))'\right] = E\left[A(X - E[X])(X - E[X])'A'\right] =$$   ($A$ non-random)

$$(CD)' = D'C'$$

$$= A E\left[(X - E[X])(X - E[X])'\right] A' = A\, Cov(X)\, A'$$

Hence,   $\boxed{Cov(AX) = A\, Cov(X)\, A'}$

$\boxed{\text{k-variate normal distribution}}$

A $k \times 1$ random vector $Z = (Z_1, \ldots, Z_k)'$ is said to have the <u>k-variate standard normal distri</u>
if $Z_1, Z_2, \ldots, Z_k$ are independent $N(0, 1)$. The density of $Z$ is given by:

$$f_Z(z) = \prod_{i=1}^{k} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_i^2} = (2\pi)^{-k/2} e^{-\|z\|^2/2}, \quad z = (z_1, \ldots, z_k) \in \mathbb{R}^k$$

The distribution of a $n \times 1$ vector $AZ$, where $A$ is a non-random $n \times k$ matrix is
called a <u>k-variate normal distribution</u> with mean $0_{n \times 1}$ and covariance

$$\Sigma = Cov(AZ) = A\, Cov(Z) A' = AA'   \quad (\text{since } Cov(Z) = I_{k \times k})$$

and is denoted simply by $N(0, \Sigma)$.    <u>Note:</u> $0$ here denotes an $n \times 1$ vector of zeroes
$\Sigma$ is an $n \times n$ matrix

One can also shift this distribution by an $n \times 1$ vector $\mu$. Letting $Y = \mu + AZ$, $Y$ is
said to have a k-variate normal distribution $N(\mu, \Sigma)$. (here, again $\Sigma = AA'$)

Notice that in the definition and final notation for, say, $N(0, \Sigma)$ we assumed that the distribution depends only on a covariance matrix $\Sigma$ and does not depend on the construction, i.e. does not depend on the choice of $Z$ and $A$. We could have started with an $m$-variate standard normal vector $\tilde{Z}$ and a non-random $n \times m$ matrix $B$ so that the covariance matrix of $B\tilde{Z}$ again happens to be equal to $\Sigma$, i.e. so that $Cov(B\tilde{Z}) = BB' = \Sigma \ (= AA')$.

Both constructions give the same multivariate normal distribution $N(0, \Sigma)$ according to our definition. Why are the distributions of $AZ$ and $B\tilde{Z}$ even the same? We show the proof here in the case when $A$ and $B$ are both $n \times n$ invertible matrices (and $Z, \tilde{Z}$ are $n$-variate standard normal vectors); the proof in general is a bit more complicated.

Let's calculate the p.d.f. of $AZ$. For every set $\Omega \subseteq \mathbb{R}^n$, we can write:

$$\mathbb{P}(AZ \in \Omega) = \mathbb{P}(Z \in A^{-1}\Omega) = \int_{A^{-1}\Omega} (2\pi)^{-n/2} e^{-\frac{1}{2}\|z\|^2} dz = \int_{\Omega} (2\pi)^{-n/2} \frac{1}{|det(A)|} e^{-\frac{1}{2}\|A^{-1}y\|^2} dy$$

$$\boxed{\begin{array}{l} y = Az \\ z = A^{-1}y \\ \text{change of variables} \end{array}}$$

Now, $det(\Sigma) = det(AA') = det(A) det(A') = (det(A))^2$

$$\|A^{-1}y\|^2 = (A^{-1}y)'(A^{-1}y) = y'(A^{-1})'(A^{-1}y) = y'(A')^{-1}A^{-1}y = y'(AA')^{-1}y = y'\Sigma^{-1}y$$

$$\boxed{\begin{array}{l} (A^{-1})' = (A')^{-1} \\ \text{when } A \text{ invertible} \end{array}} \qquad \boxed{(CD)^{-1} = D^{-1}C^{-1}}$$

So, $\mathbb{P}(AZ \in \Omega) = \int_{\Omega} (2\pi)^{-n/2} \frac{1}{\sqrt{det(\Sigma)}} e^{-\frac{1}{2}y'\Sigma^{-1}y} dy$

So, random vector $AZ$ has the density $(2\pi)^{-n/2} \cdot \frac{1}{\sqrt{det(\Sigma)}} e^{-\frac{1}{2}y'\Sigma^{-1}y}$, which depends only on $\Sigma$, and not on $A$! Hence, $AZ$ and $B\tilde{Z}$ must have the same density/distribution, and our definition of multivariate normal distribution is valid, since it depends only on $\Sigma$ not on particular choice of $A$ (and $Z$).

One nice consequence of this discussion is the density function of a $k$-variate normal distribut Let $Y \sim N(\mu, \Sigma)$. Then,

$$\boxed{f_Y(y) = (2\pi)^{-k/2} \cdot \frac{1}{\sqrt{det(\Sigma)}} e^{-\frac{1}{2}(y-\mu)'\Sigma^{-1}(y-\mu)}, \qquad y \in \mathbb{R}^k}$$

For $k = 2$, this can be written as:

$$f_Y(y_1, y_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot e^{-\frac{\frac{(y_1-\mu_1)^2}{\sigma_1^2} - 2\rho \cdot \frac{(y_1-\mu_1)(y_2-\mu_2)}{\sigma_1\sigma_2} + \frac{(y_2-\mu_2)^2}{\sigma_2^2}}{2(1-\rho^2)}}$$

known as the bivariate normal density function
$\sigma_i^2 = Var(Y_i)$
($\rho$ correlation coeff. between $Y_1$ and $Y_2$)

**Question:** Given a symmetric non-negative definite $n \times n$ matrix $\Sigma$, how does one find a matrix $A$ such that $\Sigma = AA'$? One can use, for example, the eigenvalue decomposition $\Sigma = QDQ'$, where $Q$ is orthogonal, $D$ is diagonal (with eigenvalues $\lambda_1, \dots \lambda_n$ of $\Sigma$ on its diagonal). If $D^{1/2}$ denotes the diagonal matrix with $\sqrt{\lambda_i}$ on the diagonal, one can take $A = QD^{1/2}$ or $A = QD^{1/2}Q'$ (so that $AA' = \Sigma$).

**FACT 5.** Let $Y \sim N(0_{k \times 1}, \Sigma_{k \times k})$. Let $M$ be an $m \times k$ non-random matrix. Then $MY \sim N(0_{m \times 1}, M\Sigma M')$.

*"linear transfor of normal is again normal"*

proof: $Y = AZ$ for some $k \times k$ matrix $A$ such that $\Sigma = AA'$ and a $k$-variate standard normal $Z$. Then $MY = M(AZ) = (MA)Z$ is, by definition, $m$-variate normal with mean $0_{m \times 1}$ and $cov(MY) = (MA)(MA)' = MAA'M' = M\Sigma M'$.

**FACT 6.** Let $Z \sim N(0_{k \times 1}, I_{k \times k})$ and let $Q$ be an orthogonal $k \times k$ matrix. Then $QZ \sim N(0_{k \times 1}, I_{k \times k})$.

*"orthogonal transfor of a standard norm is again standard norm"*

proof: Recall that a $k \times k$ matrix $Q$ is orthogonal when one of the following properties hold

A) $Q^{-1} = Q'$ (and hence $|\det(Q)| = 1$)

B) rows/columns of $Q$ form an orthonormal basis in $\mathbb{R}^k$

C) for any $x \in \mathbb{R}^k$ we have $\|Qx\| = \|x\|$, i.e. $Q$ preserves the length of vectors

Geometrically, orthogonal transformation represent linear transformations that preserve distance between points, such as rotations and reflections.

$$\forall \Omega \subseteq \mathbb{R}^k: \mathbb{P}(QZ \in \Omega) = \mathbb{P}(Z \in Q^{-1}\Omega) = \int_{Q^{-1}\Omega} f_Z(z)\,dz = \int_\Omega \frac{f_Z(Q^{-1}x)}{|\det(Q)|}\,dx$$

Change of var
$x = Qz$
$z = Q^{-1}x$

Since $|\det(Q)| = 1$ and $\|Q^{-1}x\| = \|x\|$, we get:

$$f_Z(Q^{-1}x) = (2\pi)^{-k/2} e^{-\|Q^{-1}x\|^2/2} = (2\pi)^{-k/2} e^{-\|x\|^2/2} = f_Z(x)$$

Hence, $\mathbb{P}(QZ \in \Omega) = \int_\Omega f_Z(x)\,dx = \mathbb{P}(Z \in \Omega) \qquad \forall \Omega \subseteq \mathbb{R}^k$

So, $QZ \sim N(0_{k \times 1}, I_{k \times k})$.

FACT 7  Uncorrelated components of a multivariate normal vector are independent.

FACT 8  Multivariate CLT (central limit theorem)

Suppose $X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{pmatrix}$ is a random $k \times 1$ vector with covariance $\Sigma$ (and $E[X_i^2] < \infty$)

Let $Y_1, Y_2, \ldots, Y_n$ be a sequence of i.i.d. copies of $X$. Then

$$S_n := \frac{1}{\sqrt{n}} \sum_{i=1}^{n} (Y_i - E[Y_i]) \xrightarrow{d} N(\underset{k \times 1}{0}, \Sigma) \quad \text{as } n \to \infty$$

where the convergence $\xrightarrow{d}$ in distribution means that for any set $\Omega \subseteq \mathbb{R}^t$

$$\lim_{n \to \infty} \mathbb{P}(S_n \in \Omega) = \mathbb{P}(Z \in \Omega) \text{ for a random vector } Z \sim N(0, \Sigma,$$

## Sample mean and covariance matrix from a multivariate normal distribution

Let $Y_1, \ldots, Y_n$ be independent $m \times 1$ random $N(\mu, \Sigma)$ vectors with $n > m$
and positive definite $\Sigma$. Define

the sample mean vector $\quad \overline{Y} = \frac{1}{n} \sum_{i=1}^{n} Y_i$

the sample covariance matrix $\quad \dfrac{W}{n-1}$, where $W = \sum_{i=1}^{n} (Y_i - \overline{Y})(Y_i - \overline{Y})'$

Generalizing the corresponding results $\overset{\text{on page 7}}{\text{in the case } m=1}$, the following facts are known

FACT 9  $\overline{Y} \sim N(\mu, \frac{\Sigma}{n})$

$\overline{Y}$ and $\dfrac{W}{n-1}$ are independent

Question: How do we generalize $\sum_{i=1}^{n} (Y_i - \overline{Y})^2 / \sigma^2 = (n-1) \dfrac{s^2}{\sigma^2} \sim \chi_{n-1}^2$ to the multivariate case

We need to generalize the $\chi^2$-distribution to the multivariate case.

## Wishart distribution

Let $Y_1, \ldots, Y_n$ be independent $N(0_{m \times 1}, \Sigma)$. The random matrix $W = \sum_{i=1}^{n} Y_i Y_i^T$ is said to have a Wishart distribution, denoted by $W_m(\Sigma, n)$.

Recall that $\chi_n^2 \equiv \Gamma(n/2, 1/2)$, so the density of $W_1(\sigma^2, n) \equiv \chi_n^2$ is

$$f_{W_1}(\omega) = \omega^{(n-2)/2} e^{-\frac{\omega}{2\sigma^2}} \cdot \frac{1}{(2\sigma^2)^{n/2} \Gamma(n/2)}$$

The density of the Wishart distribution $W_m(\Sigma, n)$ generalizes this to:

$$f(W) = \frac{(\det(W))^{(n-m-1)/2} e^{-\frac{1}{2} tr(\Sigma^{-1} W)}}{(2^m \det(\Sigma))^{n/2} \Gamma_m(n/2)} \quad \text{for all positive definite matrices } W$$

where $\Gamma_m(\cdot)$ is the multivariate gamma function.

$$\Gamma_m(t) = \pi^{m(m-1)/4} \cdot \prod_{i=1}^{m} \Gamma\left(t - \frac{i-1}{2}\right)$$

Wishart distribution has many applications in CAPM testing, and in Bayesian statistics.

<u>FACT 10</u> - If $W \sim W_m(\Sigma, n)$, then $E[W] = n\Sigma$.

- If $W_1, W_2, \ldots, W_k$ are independent with $W_i \sim W_m(\Sigma, n_i)$, then
$$\sum_{i=1}^{k} W_i \sim W_m\left(\Sigma, \sum_{i=1}^{k} n_i\right)$$

- If $W \sim W_m(\Sigma, n)$ and $A$ is a nonrandom $m \times m$ nonsingular matrix, then $AWA' \sim W_m(A\Sigma A', n)$.

- $W = \sum_{i=1}^{n} (Y_i - \bar{Y})(Y_i - \bar{Y})' \sim W_m(\Sigma, n-1)$.

$\boxed{\text{Multivariate } t\text{-distribution}}$

Let $Z \sim N(O_{m \times 1}, \Sigma)$ and $W \sim W_m(\Sigma, k)$ be independent. Then $(W/k)^{-1/2} \cdot Z$ is said to have the $m$-variate $t$-distribution with $k$ degrees of freedom.

The density function is
$$f(t) = \frac{\Gamma\left(\frac{k+m}{2}\right)}{(\pi k)^{\frac{m}{2}} \Gamma(k/2)} \left(1 + \frac{\|t\|^2}{k}\right)^{-\frac{k+m}{2}} \quad , \quad t \in \mathbb{R}^m$$

(which reduces to the expression on bottom of page 6. in the case $m=1$).

$m$-variate $t$-distribution is used in risk management (statistical models for VaR) and $t$-copulas.

The square of a $t_k$-distributed random variable (i.e. univariate with $k$ degrees of freedom) is actually $\underline{F_{1,k}}$-distributed.

More generally, if $t$ has the $m$-variate $t$-distribution with $k$ degrees of freedom ($k \geq m$), then $\frac{k-m+1}{km} \|t\|^2$ has the $F_{m, k-m+1}$ distribution $\qquad \circledast$

Now, let's go back to the sample setup:

Let $Y_1, \ldots, Y_n$ independent $m \times 1$ random $N(\mu, \Sigma)$ vectors with $n > m$ and positive definite $\Sigma$. We know that $\bar{Y} = \frac{1}{n} \sum_{i=1}^{n} Y_i$ and $\frac{W}{n-1} = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i - \bar{Y})(Y_i - \bar{Y})'$ are independent; $\bar{Y} \sim N\left(\mu, \frac{\Sigma}{n}\right)$; $W \sim W_m(\Sigma, n-1)$.

Define the $\underline{\text{Hotelling's } T^2\text{-statistic}}$ (famous in multivariate hypothesis testing)
$$T^2 = n(\bar{Y} - \mu)' \left(\frac{W}{n-1}\right)^{-1} (\bar{Y} - \mu)$$

Note that $\left(\frac{W}{n-1}\right)^{-1/2}(\sqrt{n}(\bar{Y} - \mu)) \sim \left(\frac{W_m(\Sigma, n-1)}{n-1}\right)^{-1/2} N(0, \Sigma)$ has the $m$-variate $t$-distribution with $n-1$ degrees of freedom

Then, according to $\circledast$, $\frac{n-1-m+1}{(n-1)m} T^2 = \frac{n-m}{(n-1)m} T^2 \sim F_{m, n-1-m+1}$, i.e $\boxed{\frac{n-m}{(n-1)m} T^2 \sim F_{m, n-m}}$

-13-

# 2 Method of maximum likelihood

Given data of any kind, we're often faced with the following questions:

A) How to estimate the unknown parameters of a distribution given the data from it?

B) How good are these estimates; are they close to the actual "true" parameters?

C) Does the data come from a particular type of distribution; for example, normal or gamma

First, we'll keep it simple and study Questions A) and B), while assuming that we know what type of distribution the sample comes from (so we only do not know the parameters of this distribution

Consider a family of distributions $\mathbb{P}_\Theta$ indexed by a parameter (which, in general, could be a vector of parameters) $\Theta$ that belongs to a set $\Theta$. For example, we could be considering a family of normal distributions $N(\mu, \sigma^2)$ in which case $\Theta = (\mu, \sigma^2)'$.

Let $f_\Theta(x_1, x_2, ..., x_n)$ be the joint density function of $X_1, ..., X_n$. The <u>likelihood function</u> based on the observations $X_1, ..., X_n$ is $L(\Theta) = f_\Theta(X_1, ..., X_n)$ and the MLE (<u>maximum likelihood estimate</u>) $\hat{\Theta}$ of $\Theta$ is the value of $\Theta$ that maximizes $L(\Theta)$, over all $\Theta \in \Theta$.

More often than not, the sample $X_1, ..., X_n$ is assumed to be independent[*], so

$$L(\Theta) = f_\Theta(X_1) \cdot f_\Theta(X_2) \cdots f_\Theta(X_n),$$ where $f_\Theta(x)$ is the p.d.f. of the distribution $\mathbb{P}_\Theta$

(Make sure you understand that $X_1, ..., X_n$ are given; so $L$ is a function of $\Theta$ only!)

Intuitively, the likelihood function is the probability to observe the sample $X_1, ..., X_n$ when the unknown parameters of the distribution are equal to $\Theta$.

When finding the MLE, it is sometimes easier to maximize the <u>log-likelihood function</u>

$$\ell(\Theta) = \log f_\Theta(X_1, ..., X_n) \text{ instead } \left(\underline{\text{Note:}} \ \log x \text{ is an increasing function}\right)$$

When $X_1, ..., X_n$ are independent, then $\ell(\Theta) = \sum_{i=1}^{n} \ell_\Theta(X_i)$, where $\ell_\Theta(x) = \log f_\Theta(x)$

Let's do several examples of calculating the MLE!

---

[*] Without assuming that $X_i$ are independent, laws of large numbers and CLT could not be applied. However, one could still use martingale strong laws and central limit theorems; and most of the results here would still hold, under some regularity conditions.

−14−

<u>Example 1.</u>  <u>Bernoulli distribution $B(p)$</u>

$X \sim B(p)$ $\qquad 0 \le p \le 1$ $\qquad\qquad \mathbb{P}(X=1) = p$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \mathbb{P}(X=0) = 1-p$

p.d.f $\quad f_p(x) = \begin{cases} p, & \text{if } x=1 \\ 1-p, & \text{if } x=0 \end{cases}$ $\quad$ OR $\qquad f_p(x) = p^x(1-p)^{1-x}$

likelihood function $\quad L(p) = \prod\limits_{i=1}^{n} f_p(X_i) = p^{\overset{\#\text{of 1's in } X_1,\dots,X_n}{}} \cdot (1-p)^{\overset{\#\text{of 0's in } X_1,\dots,X_n}{}}$

$$L(p) = p^{X_1+\dots+X_n}(1-p)^{n-(X_1+\dots+X_n)}$$

log-likelihood function $\quad \ell(p) = \left(\sum\limits_{i=1}^{n} X_i\right)\log p + \left(n - \sum\limits_{i=1}^{n} X_i\right)\log(1-p)$

$\dfrac{d}{dp}(\ell(p)) = 0 \implies \dfrac{1}{p}\sum\limits_{i=1}^{n} X_i - \left(n - \sum\limits_{i=1}^{n} X_i\right)\cdot\dfrac{1}{1-p} = 0$ $\qquad$ solve for $p$ to get:

$$p = \frac{X_1+\dots+X_n}{n} = \bar{X}$$

Therefore, the proportion of successes $\hat{p} = \bar{X}$ in the sample is the MLE for $p$, which is perfectly intuitive. Note that by the law of large numbers (FACT 2), we have $\hat{p} = \bar{X} \to \mathbb{E}[X_1] = p$ (in probability), which means that our MLE will approximate the unknown parameter $p$ well when we get more and more data. More about this in a second, once we start talking about consistency of the MLE

<u>Example 2.</u>  <u>Normal distribution $N(\mu, \sigma^2)$</u> $\qquad$ Let $\theta = (\mu, \sigma^2)'$

p.d.f. $\quad f_\theta(x) = \dfrac{1}{\sqrt{2\pi}\cdot\sigma}\, e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$L(\theta) = \prod\limits_{i=1}^{n} \dfrac{1}{\sqrt{2\pi}\cdot\sigma}\, e^{-\frac{(X_i-\mu)^2}{2\sigma^2}}$

$\ell(\theta) = \sum\limits_{i=1}^{n}\left(\log\left(\dfrac{1}{\sqrt{2\pi}}\right) - \log\sigma - \dfrac{(X_i-\mu)^2}{2\sigma^2}\right)$

$\ell(\theta) = n\log\left(\dfrac{1}{\sqrt{2\pi}}\right) - n\log\sigma - \dfrac{1}{2\sigma^2}\sum\limits_{i=1}^{n}(X_i-\mu)^2$

We need to maximize $\ell(\theta)$ over $\theta = \left\{(\mu,\sigma^2)' \mid \mu\in(-\infty,\infty),\ \sigma^2 > 0\right\} \subseteq \mathbb{R}^2$

The usual approach would be to find the MLE by solving the equation $\nabla\ell(\theta) = 0_{2\times 1}$, where $\nabla\ell$ is the gradient vector of partial derivatives

$\nabla \ell = \begin{pmatrix} \frac{d\ell}{d\theta_1} \\ \frac{d\ell}{d\theta_2} \end{pmatrix} \equiv \begin{pmatrix} \frac{d\ell}{\partial \mu} \\ \frac{d\ell}{\partial \sigma^2} \end{pmatrix}$ ; but in this case, we the "conditional log likelihood method".

First, for any $\sigma^2$, we minimize $\sum_{i=1}^{n}(X_i - \mu)^2$ over $\mu$:

$$\frac{d}{d\mu}\sum_{i=1}^{n}(X_i - \mu)^2 = 0 \text{ gives } -2\sum_{i=1}^{n}(X_i - \mu) = 0, \text{ i.e. } \boxed{\hat{\mu} = \bar{X}}.$$

We plug in this estimate in the log-likelihood function to obtain the conditional log-likelihood function :

$$n\log\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{n}{2}\log(\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^{n}(X_i - \bar{X})^2 \text{, which needs to be maximized}$$

over $\sigma^2$. Letting $\gamma \equiv \sigma^2$ in the last expression and taking $\frac{d}{d\gamma}$, we get:

$$-\frac{n}{2\gamma} + \frac{1}{2\gamma^2}\sum_{i=1}^{n}(X_i - \bar{X})^2 = 0 \text{, i.e. } \hat{\gamma} \equiv \boxed{\hat{\sigma^2} = \frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X})^2},$$

which is only slightly different from the sample variance $s^2 = \frac{1}{n-1}\sum_{i=1}^{n}(X_i - \bar{X})^2$

Example 3. Uniform distribution $U[0,\theta]$

p.d.f. $f_\theta(x) = \begin{cases} 1/\theta & \text{, if } 0 \leq x \leq \theta \\ 0 & \text{, otherwise} \end{cases}$

$$L(\theta) = \prod_{i=1}^{n} f_\theta(X_i) = \frac{1}{\theta^n} \cdot \mathcal{I}\left(X_1 \in [0,\theta] \text{ and } X_2 \in [0,\theta] \text{ and } \dots \text{ and } X_n \in [0,\theta]\right)$$

indicator random variable for the event in parentheses

Simpler way to write it:

$$L(\theta) = \frac{1}{\theta^n}\mathcal{I}\left(\max(X_1, \dots, X_n) \leq \theta\right) = \begin{cases} 0 & \text{, if } \theta < \max(X_1, \dots, X_n) \\ \frac{1}{\theta^n} & \text{, if } \theta \geq \max(X_1, \dots, X_n) \end{cases}$$



In this example, no need to go to $\ell(\theta)$. Also, we cannot differentiate $L(\theta)$ w.r.t $\theta$.

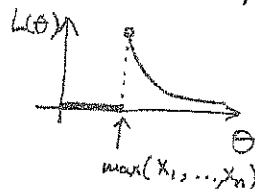Nonetheless, it's easy to see how to maximize the likelihood function!

Set $\hat{\theta} = \max(X_1, \dots, X_n)$. This is the MLE for $\theta$!

Note: It is often not easy to find the MLE as in the above examples, so numerical procedur (such as Newton-Raphson) need to be used. Also, MLE does not always exist!

Here is an artificial example based on Example 3.

Consider $\mathbb{P}_\theta$ to be $U[0,\theta)$ uniform on $[0,\theta)$ (where $\theta$ is unknown).

Then, similarly as before, $L(\theta) = \frac{1}{\theta^n} I(\max(X_1,\ldots,X_n) < \theta)$ and the maximum at the point $\hat{\theta} = \max(X_1,\ldots,X_n)$ is not achieved.

Question: Why are MLE's good?

Because of Consistency and asymptotic normality. Next, we explain these concepts in the univariate case, i.e. when $\theta$ is just 1-dimensional.

The multivariate case (when $\theta$ is a vector of unknown parameters) is very similar and will be mentioned at the end of this section.

Consistency ("no bias") We say that the MLE $\hat{\theta}$ is consistent if $\hat{\theta} \to \theta_0$ in probability, as $n \to \infty$, where $\theta_0$ is the true value of the unknown parameter of the distribution of the sample.

("$\hat{\theta} \to \theta_0$ in prob. as $n \to \infty$" means "$\forall \varepsilon > 0$ $\mathbb{P}(|\hat{\theta} - \theta_0| > \varepsilon) \to 0$ as $n \to \infty$")

Asymptotic normality We say that $\hat{\theta}$ is asymptotically normal if as $n \to \infty$
$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \sigma_{\theta_0}^2) \quad \text{for some } \sigma_{\theta_0}^2 \text{ which is called the}$$
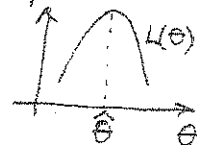asymptotic variance of the estimator $\hat{\theta}$.

$\sigma_{\theta_0}$ is also known as the standard error of $\hat{\theta}$.

Actually, we will show that $\sigma_{\theta_0}^2 = \frac{1}{I(\theta_0)}$, where $I(\theta_0)$ will be defined later as the Fisher information.

Let's attempt to prove Consistency, at least intuitively. Assume that the likelihood function $L(\theta) = \prod_{i=1}^{n} f_\theta(X_i)$ is smooth and its maximum is achieved at a unique point $\hat{\theta}$.

$\hat{\theta}$ also maximizes $\ell(\theta)$ (the log-likelihood), as well as

$\ell_n(\theta) := \frac{1}{n}\ell(\theta) = \frac{1}{n}\sum_{i=1}^{n} \ell_\theta(X_i)$ (the log-likelihood "normalized" by $1/n$)

Aside from $l_n(\theta)$, we also introduce $\mathcal{L}(\theta) := E_{\theta_0}[l_\theta(x)]$, i.e.

$$\mathcal{L}(\theta) = \int l_\theta(x) f_{\theta_0}(x)\, dx$$

Note that while $l_n$ depends on the sample $X_1, X_2, \ldots, X_n$; $\mathcal{L}(\theta)$ only depends on $\theta$.

By FACT2 (Law of large numbers), for any $\theta$, we have

$$l_n(\theta) \longrightarrow E_{\theta_0}[l_\theta(x)] = \mathcal{L}(\theta). \quad (*)$$

<u>Claim</u>  $\mathcal{L}(\theta) \leq \mathcal{L}(\theta_0)$ for every $\theta$.

proof: Recall Jensen's inequality: Suppose $Z$ has a finite mean and $\psi: \mathbb{R} \to \mathbb{R}$ is convex, i.e.
$$\lambda \psi(x) + (1-\lambda)\psi(y) \geq \psi(\lambda x + (1-\lambda)y) \quad \forall\, 0 < \lambda < 1, \forall\, x, y \in \mathbb{R}.$$
Then $E[\psi(Z)] \geq \psi(E[Z])$.

Since $\log$ is a concave function (i.e. $-\log$ is convex), then by Jensen's inequality for
$$Z \equiv \frac{f_\theta(x)}{f_{\theta_0}(x)}, \text{ we get: } E_{\theta_0}[-\log Z] \geq -\log(E_{\theta_0}[Z]), \text{ i.e.}$$

$$E_{\theta_0}[\log Z] \leq \log(E_{\theta_0}[Z]), \text{ i.e.}$$

$$E_{\theta_0}\left[\log\left(\frac{f_\theta(x)}{f_{\theta_0}(x)}\right)\right] \leq \log\left(\int \frac{f_\theta(x)}{f_{\theta_0}(x)} f_{\theta_0}(x)\, dx\right) = \log\left(\int f_\theta(x)\, dx\right) = \log 1 = \overset{\overset{\displaystyle 1}{\|}}{}$$

(integrating the p.d.f.)

$$\underset{\|\ \text{linearity of expectation}}{E_{\theta_0}[\log(f_\theta(x))] - E_{\theta_0}[\log(f_{\theta_0}(x))]}$$

$$\underset{\|\ \text{definition of } l}{}$$

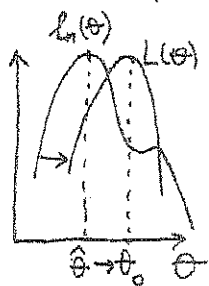$$\underset{\|\ \text{definition of } \mathcal{L}}{E_{\theta_0}[l_\theta(x)] - E_{\theta_0}[l_{\theta_0}(x)]}$$

$$\mathcal{L}(\theta) - \mathcal{L}(\theta_0)$$

So, $\mathcal{L}(\theta) - \mathcal{L}(\theta_0) \leq 0$, i.e. $\mathcal{L}(\theta) \leq \mathcal{L}(\theta_0)$ □.

<u>Theorem</u> Under some regularity conditions on the family of distributions (which we don't state here), the MLE $\hat{\theta}$ is consistent, i.e. $\hat{\theta} \to \theta_0$ as $n \to \infty$.

(intuitive) proof: Let's summarize what we discussed so far:



1. $\hat{\theta}$ is the maximizer of $l_n(\theta)$ (by definition of MLE)
2. $\theta_0$ is the maximizer of $\mathcal{L}(\theta)$ by previous claim.
3. $\forall \theta$, we have $l_n(\theta) \to \mathcal{L}(\theta)$ by (*).

Since two functions $l_n$ and $\mathcal{L}$ are getting closer, the points of maximum should also get closer, which means $\hat{\theta} \to \theta_0$ $\square$

Before we discuss asymptotic normality, let's define $\underline{\text{Fisher information}}$:

$\underline{\text{Definition}}$ Fisher information of a random variable $X$ with distribution $\mathbb{P}_{\theta_0}$
$\underline{\text{(univariate case)}}$ from the family $\{\mathbb{P}_\theta : \theta \in \Theta\}$ is defined by:

$$I(\theta_0) = E_{\theta_0}\left[\left(\frac{d}{d\theta} l_\theta(X)\Big|_{\theta=\theta_0}\right)^2\right]$$

$\underline{\text{Question}}$: How to think of Fisher information intuitively?

It can be interpreted as a measure of how quickly the p.d.f. will change when one changes the parameter $\theta$ near $\theta_0$ slightly.

$$\frac{d}{d\theta} l_\theta(X)\Big|_{\theta=\theta_0} = \frac{d}{d\theta}\left(\log f_\theta(X)\right)\Big|_{\theta=\theta_0} = \frac{\left(\frac{d}{d\theta} f_\theta(x)\right)\Big|_{\theta=\theta_0}}{f_{\theta_0}(x)}$$

When we square this and take expectation, i.e. average over all $x$, we get an "average" version of this measure. So, if Fisher information is large, this means that the distribution will change quickly when we move the parameter away from its true value $\theta_0$, so the distribution with parameter $\theta_0$ is "quite different" and can be "well distinguished" from the distributions with parameters away from $\theta_0$. If Fisher information is small, this means that the distribution is "very similar" to the distributions with parameters away from $\theta_0$ and thus more difficult to distinguish; so our estimation will be worse.

There is an alternative formula for Fisher information, which is actually used more often than the definition above.

-19-

Formula for $I(\theta_0)$: $\quad I(\theta_0) = -E_{\theta_0}\left[\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(x)\right)\Big|_{\theta=\theta_0}\right]$

proof: First, let's find an "appropriate" expression for $\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(x)\right)\Big|_{\theta=\theta_0}$ to be used in this proof

$$\left(\frac{\partial}{\partial\theta}\ell_\theta(x)\right)\Big|_{\theta=\theta_0} = \frac{\left(\frac{\partial}{\partial\theta}f_\theta(x)\right)\big|_{\theta=\theta_0}}{f_{\theta_0}(x)} \qquad \text{using the definition of } \ell_\theta(x) \text{ as } \log f_\theta(x).$$

Quotient rule then gives:

$$\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(x)\right)\Big|_{\theta=\theta_0} = \frac{\left(\frac{\partial^2}{\partial\theta^2}f_\theta(x)\right)\big|_{\theta=\theta_0}}{f_{\theta_0}(x)} - \left(\frac{\left(\frac{\partial}{\partial\theta}f_\theta(x)\right)\big|_{\theta=\theta_0}}{f_{\theta_0}(x)}\right)^2 \text{, i.e.}$$

$$\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(x)\right)\Big|_{\theta=\theta_0} = \frac{\left(\frac{\partial^2}{\partial\theta^2}f_\theta(x)\right)\big|_{\theta=\theta_0}}{f_{\theta_0}(x)} - \left(\left(\frac{\partial}{\partial\theta}\ell_\theta(x)\right)\Big|_{\theta=\theta_0}\right)^2 \quad \oplus$$

Now,

$$E_{\theta_0}\left[\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(x)\right)\Big|_{\theta=\theta_0}\right] = \int\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(x)\right)\Big|_{\theta=\theta_0}\cdot f_{\theta_0}(x)\,dx =$$

$$= \int\frac{\left(\frac{\partial^2}{\partial\theta^2}f_\theta(x)\right)\big|_{\theta=\theta_0}}{f_{\theta_0}(x)}\cdot f_{\theta_0}(x)\,dx - \int\left(\left(\frac{\partial}{\partial\theta}\ell_\theta(x)\right)\Big|_{\theta=\theta_0}\right)^2 f_{\theta_0}(x)\,dx =$$

$$= \int\left(\frac{\partial^2}{\partial\theta^2}f_\theta(x)\right)\Big|_{\theta=\theta_0}\,dx - E_{\theta_0}\left[\left(\frac{\partial}{\partial\theta}\ell_\theta(x)\Big|_{\theta=\theta_0}\right)^2\right] = 0 - I(\theta_0) = -I(\theta_0)$$

Why is $\int\left(\frac{\partial^2}{\partial\theta^2}f_\theta(x)\right)\Big|_{\theta=\theta_0}\,dx = 0$?

Since $f_\theta(x)$ is a p.d.f., one has: $\int f_\theta(x)\,dx = 1$, so $\frac{\partial}{\partial\theta}\int f_\theta(x)\,dx = 0$, i.e.

$\int\frac{\partial}{\partial\theta}f_\theta(x)\,dx = 0$ and $\int\left(\frac{\partial^2}{\partial\theta^2}f_\theta(x)\right)\Big|_{\theta=\theta_0}\,dx = 0$. $\qquad\square$

Theorem 2  $\sqrt{n}\,(\hat{\Theta}-\Theta_0) \xrightarrow{d} N\!\left(0, \frac{1}{I(\Theta_0)}\right)$  as $n \to \infty$

Note:  Larger $I(\Theta_0)$ means better MLE, since its asymptotic variance/dispersion around the true value of the parameter is smaller.

proof: Recall from the proof of Theorem 1 that $\hat{\Theta}$ maximizes $\ell_n(\Theta) = \frac{1}{n}\sum_{i=1}^{n}\log f_\Theta(x_i)$

So:  $\left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\Big|_{\Theta=\hat{\Theta}} = 0$.

By Taylor's theorem (or by mean value theorem: $f(a) = f(b) + f'(c)(a-b)$ for some $c \in [a,b]$)

we get that for some $\Theta^* \in [\hat{\Theta},\Theta_0]$:  $\longleftarrow$ $f(\Theta) = \frac{\partial}{\partial\Theta}\ell_n(\Theta)$, $a=\hat{\Theta}$, $b=\Theta_0$

$$0 = \left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\Big|_{\Theta=\hat{\Theta}} = \left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\Big|_{\Theta=\Theta_0} + \left(\frac{\partial^2}{\partial\Theta^2}\ell_n(\Theta)\right)\Big|_{\Theta=\Theta^*}\cdot(\hat{\Theta}-\Theta_0)$$

So  $\hat{\Theta}-\Theta_0 = -\dfrac{\left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\big|_{\Theta=\Theta_0}}{\left(\frac{\partial^2}{\partial\Theta^2}\ell_n(\Theta)\right)\big|_{\Theta=\Theta^*}}$ , i.e. $\boxed{\sqrt{n}\,(\hat{\Theta}-\Theta_0) = -\dfrac{\sqrt{n}\left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\big|_{\Theta=\Theta_0}}{\left(\frac{\partial^2}{\partial\Theta^2}\ell_n(\Theta)\right)\big|_{\Theta=\Theta^*}}}$ $\circledast$

Recall from the proof of Theorem 1 that $\Theta_0$ maximizes $\mathcal{Z}(\Theta) = E_{\Theta_0}\!\left[\ell_\Theta(x)\right]$, so

$$\left(\frac{\partial}{\partial\Theta}\mathcal{Z}(\Theta)\right)\Big|_{\Theta=\Theta_0} = E_{\Theta_0}\!\left[\left(\frac{\partial}{\partial\Theta}\ell_\Theta(x)\right)\Big|_{\Theta=\Theta_0}\right] = 0 \qquad \oplus$$

Consider the numerator in $\circledast$. We have:

$$\sqrt{n}\left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\Big|_{\Theta=\Theta_0} = \sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial}{\partial\Theta}\ell_\Theta(x_i)\right)\Big|_{\Theta=\Theta_0} - 0\right) =$$

$\left(\text{using } \oplus \text{ here}\right) \longrightarrow = \sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial}{\partial\Theta}\ell_\Theta(x_i)\right)\Big|_{\Theta=\Theta_0} - E_{\Theta_0}\!\left[\left(\frac{\partial}{\partial\Theta}\ell_\Theta(x_1)\right)\Big|_{\Theta=\Theta_0}\right]\right)$

By FACT3 (central limit theorem)

$$\sqrt{n}\left(\frac{\partial}{\partial\Theta}\ell_n(\Theta)\right)\Big|_{\Theta=\Theta_0} \xrightarrow{d} N\!\left(0, \operatorname{Var}_{\Theta_0}\!\left(\left(\frac{\partial}{\partial\Theta}\ell_\Theta(x_1)\right)\Big|_{\Theta=\Theta_0}\right)\right)$$

-21-

Consider the denominator in ⊛. By FACT 2 (Law of large numbers), for every $\theta$:

$$\frac{\partial^2}{\partial\theta^2} \ell_n(\theta) = \frac{1}{n}\sum_{i=1}^n \frac{\partial^2}{\partial\theta^2}\ell_\theta(x_i) \longrightarrow E_{\theta_0}\left[\frac{\partial^2}{\partial\theta^2}\ell_\theta(X_1)\right]$$

Since $\theta^* \in [\hat\theta, \theta_0]$ and since (by Theorem 1) $\hat\theta \longmapsto \theta_0$, then also $\theta^* \to \theta_0$, and:

$$\left(\frac{\partial^2}{\partial\theta^2}\ell_n(\theta)\right)\Big|_{\theta=\theta^*} \longrightarrow E_{\theta_0}\left[\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(X_1)\right)\Big|_{\theta=\theta^*}\right] \to E_{\theta_0}\left[\left(\frac{\partial^2}{\partial\theta^2}\ell_\theta(X_1)\right)\Big|_{\theta=\theta_0}\right], \ i.e.$$

by Formula for $I(\theta_0)$:

$$\left(\frac{\partial^2}{\partial\theta^2}\ell_n(\theta)\right)\Big|_{\theta=\theta^*} \longrightarrow -I(\theta_0).$$

<u>Slutsky's theorem</u>: If $Z_n \xrightarrow{d} Z$; $W_n \to \omega$ in probability (where $\omega$ is nonrandom),

$$\text{then} \quad W_n Z_n \xrightarrow{d} \omega Z$$

Using this theorem and our convergence results for the numerator & denominator of ⊛, we get:

$$\sqrt{n}(\hat\theta - \theta_0) \xrightarrow{d} N\left(0, \frac{Var_{\theta_0}\left(\left(\frac{\partial}{\partial\theta}\ell_\theta(X_1)\right)\Big|_{\theta=\theta_0}\right)}{(I(\theta_0))^2}\right)$$

Finally, $Var_{\theta_0}\left(\left(\frac{\partial}{\partial\theta}\ell_\theta(X_1)\right)\Big|_{\theta=\theta_0}\right) = E_{\theta_0}\left[\left(\left(\frac{\partial}{\partial\theta}\ell_\theta(X_1)\right)\Big|_{\theta=\theta_0}\right)^2\right] - \left(E_{\theta_0}\left[\left(\frac{\partial}{\partial\theta}\ell_\theta(X_1)\right)\Big|_{\theta=\theta_0}\right]\right)^2 =$

$$= I(\theta_0) - 0 = I(\theta_0).$$

by definition of Fisher information and by ⊕

So, $\sqrt{n}(\hat\theta - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right)$ □

Let's do some examples.

Example 4. The family of Bernoulli distributions $B(p)$ has p.d.f. $f_p(x) = p^x(1-p)^{1-x}$.

So, $\ell_p(x) = \log f_p(x) = x \log p + (1-x) \log(1-p)$

$$\frac{\partial}{\partial p} \ell_p(x) = \frac{x}{p} - \frac{1-x}{1-p} \quad ; \quad \frac{\partial^2}{\partial p^2} \ell_p(x) = -\frac{x}{p^2} - \frac{1-x}{(1-p)^2}$$

Fisher information $I(p)$ by formula on page 20. is:

$$I(p_0) = -\mathop{E}_{p_0}\left[\left(\frac{\partial^2}{\partial p^2} \ell_p(x)\right)\Big|_{p=p_0}\right] = \frac{E_{p_0}[X]}{p_0^2} + \frac{1 - E_{p_0}[X]}{(1-p_0)^2} = \frac{p_0}{p_0^2} + \frac{1-p_0}{(1-p_0)^2}$$

$$\Rightarrow I(p_0) = \frac{1}{p_0(1-p_0)}$$

We know from example 1 that the MLE $\hat{p} = \bar{X}$. Asymptotic normality of $\hat{p}$ now states that:

$$\sqrt{n}(\hat{p} - p_0) \xrightarrow{d} N(0, p_0(1-p_0)) \qquad \text{as } n \to \infty$$

(which, of course, also follows directly from the CLT! Check!).

Example 5. The family of exponential distributions $\Xi(\alpha)$ has p.d.f. :

$$f_\alpha(x) = \begin{cases} \alpha e^{-\alpha x} & \text{, if } x \geq 0 \\ 0 & \text{, if } x < 0 \end{cases}$$

Check that the MLE $\hat{\alpha} = \frac{1}{\bar{X}}$ !

$$\ell_\alpha(x) = \log f_\alpha(x) = \log \alpha - \alpha x$$

$$\frac{\partial}{\partial \alpha} \ell_\alpha(x) = \frac{1}{\alpha} - x \quad ; \quad \frac{\partial^2}{\partial \alpha^2} \ell_\alpha(x) = -\frac{1}{\alpha^2} \quad (\text{does not depend on } X)$$

$$I(\alpha_0) = -\mathop{E}_{\alpha_0}\left[\left(\frac{\partial^2}{\partial \alpha^2} \ell_\alpha(x)\right)\Big|_{\alpha=\alpha_0}\right] = -\mathop{E}_{\alpha_0}\left[-\frac{1}{\alpha_0^2}\right] = \frac{1}{\alpha_0^2}$$

So, the asymptotic normality gives:

$$\sqrt{n}\left(\frac{1}{\bar{X}} - \alpha_0\right) \xrightarrow{d} N(0, \alpha_0^2) \qquad \text{as } n \to \infty$$

**Example 6** The family of normal distributions $N(\mu, \sigma^2)$ with known $\sigma^2$ (but unknown $\mu$).

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

$$l_\mu(x) = \log f_\mu(x) = -\frac{1}{2}\left(\log(2\pi) + \log(\sigma^2)\right) - \frac{1}{2\sigma^2}(x-\mu)^2$$

$$\frac{\partial}{\partial\mu} l_\mu(x) = \frac{1}{\sigma^2}(x-\mu) \quad ; \quad \frac{\partial^2}{\partial\mu^2} l_\mu(x) = -\frac{1}{\sigma^2} \quad \text{(does not depend on } x)$$

$$I(\mu_0) = -E_{\mu_0}\left[\left(\frac{\partial^2}{\partial\mu^2} l_\mu(x)\right)\Big|_{\mu=\mu_0}\right] = -E_{\mu_0}\left[-\frac{1}{\sigma^2}\right] = \frac{1}{\sigma^2}$$

So, the asymptotic normality gives:

$$\sqrt{n}(\hat\mu - \mu_0) \xrightarrow{d} N(0, \sigma^2). \qquad \text{Now, the MLE } \hat\mu = \bar{X} \text{ (check! why?)}$$

$$\sqrt{n}(\bar{X} - \mu_0) \xrightarrow{d} N(0, \sigma^2), \text{ i.e. } \bar{X} \xrightarrow{d} N\left(\mu_0, \frac{\sigma^2}{n}\right), \text{ as } n\to\infty$$

This matches the usual statistician's statement that the standard error of the MLE is $\frac{\sigma}{\sqrt{n}}$, when $\sigma$ is known.

**Question:** What if $\theta$, the vector of unknown parameters, is, say, $p$-dimensional?
(multivariate case) Not much changes! The proofs and results are very similar.

<u>Fisher information matrix is a $p \times p$ matrix</u>

$$I(\theta_0) = -\left(E_{\theta_0}\left[\left(\frac{\partial^2}{\partial\theta_i \partial\theta_j} l_\theta(x)\right)\Big|_{\theta=\theta_0}\right]\right)_{1\le i,j\le p}$$

Asymptotic variance now states that

$$\sqrt{n}(\hat\theta - \theta_0) \xrightarrow{d} N\left(0, (I(\theta_0))^{-1}\right) \qquad \text{as } n\to\infty$$

Here, $\theta_0$ denotes the $p\times 1$ vector of true values of unknown parameters

<u>FACT 11</u> (no proof here)

$$(\hat\theta - \theta_0)'\left(-\left(\nabla^2 l(\theta)\right)\Big|_{\theta=\hat\theta}\right)(\hat\theta - \theta_0) \xrightarrow{d} \chi^2_p \qquad \text{as } n\to\infty$$

where $\nabla^2 l(\theta) = \left(\frac{\partial^2 l}{\partial\theta_i \partial\theta_j}\right)_{1\le i,j\le p}$ is the Hessian $p \times p$ matrix of second partial derivatives

$-\nabla^2(l(\theta))\big|_{\theta=\hat\theta}$ is called the <u>observed Fisher information matrix</u>

Set-up: MLE's for normal distribution and their distribution

Consider a random sample $X_1, X_2, \ldots, X_n \sim N(\mu, \sigma^2)$, $\mu, \sigma^2$ unknown.

From Example 2, Section 2:

MLE for $\mu$: $\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ (some algebra)

MLE for $\sigma^2$: $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^2 = \overline{X^2} - (\bar{X})^2$, where $\overline{X^2} = \frac{1}{n} \sum X_i^2$

Question: How close are these estimates to the true values?

Note: By Law of Large Numbers $\bar{X} \to \mu$, $\overline{X^2} - (\bar{X})^2 \to \sigma^2$ as $n \to \infty$ but how quick is the convergence? Can we say more?

Answer: Want to construct intervals of values that will with certain pre-specified probability contain the true value of the parameter.
                                                                    (unknown)

Note: Unlike in Bayesian theory where the parameter is considered random and interval fixed, and where we talk about a probability that the unknown parameter belongs to a fixed ("credible") interval; HERE, the unknown parameter is fixed, while the endpoints of the interval are random and have a probability distribution

We start with

Theorem 1: If $X_1, \ldots, X_n$ are i.i.d. $\sim N(0,1)$, then the sample mean $\bar{X}$ and the MLE variance $\overline{X^2} - (\bar{X})^2$ are independent and

$$\sqrt{n}\, \bar{X} \sim N(0,1) \quad \text{and} \quad n\left(\overline{X^2} - (\bar{X})^2\right) \sim \chi^2_{n-1}$$

In other words, $\hat{\mu}$ and $\hat{\sigma}^2$ are independent,

$$\sqrt{n}\, \hat{\mu} \sim N(0,1) \quad \text{and} \quad n\hat{\sigma}^2 \sim \chi^2_{n-1}$$

(Chi-squared distribution with $n-1$ degrees of freedom)

Note: This is actually Theorem on page 7 (Section 1)! Notice that there is a slight difference between the sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$ and the MLE variance $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^2$

proof: Consider $Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = VX = \begin{pmatrix} 1/\sqrt{n} & \cdots & 1/\sqrt{n} \\ & ? & \end{pmatrix} \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$

some orthogonal transformation of $X$, where the first row of $V$ is taken to be $(1/\sqrt{n} \ldots 1/\sqrt{n})'$ and the remaining rows are to be any orthogonal basis in the hyperplane orthogonal to this unit vector.

Now, $Y_1, \ldots, Y_n$ are also i.i.d. standard normal (FACT 6.) and moreover, $Y_1 = \frac{1}{\sqrt{n}} X_1 + \ldots + \frac{1}{\sqrt{n}} X_n = \sqrt{n}\, \overline{X}$, i.e. $\boxed{\overline{X} = \frac{1}{\sqrt{n}} Y_1}$ (1)

Also, $n(\overline{X^2} - (\overline{X})^2) = X_1^2 + \ldots + X_n^2 - \left(\frac{1}{\sqrt{n}}(X_1 + \ldots + X_n)\right)^2 = X_1^2 + \ldots + X_n^2 - Y_1^2$

$V$ is orthogonal, so $\|Y\| = \|VX\| = \|X\|$, i.e. $Y_1^2 + \ldots + Y_n^2 = X_1^2 + \ldots + X_n^2$,

So $\boxed{n(\overline{X^2} - (\overline{X})^2) = Y_1^2 + \ldots + Y_n^2 - Y_1^2 = Y_2^2 + \ldots + Y_n^2}$ (2)

So, (1) & (2) $\Rightarrow$ $\overline{X}$ and $\overline{X^2} - (\overline{X})^2$ are independent,

$\sqrt{n}\, \overline{X} = Y_1 \sim N(0,1)$ and $n(\overline{X^2} - (\overline{X})^2) \sim \chi^2_{n-1}$ $\qquad \square$

Corollary 1 If $X_1, \ldots, X_n$ are i.i.d. $\sim N(\mu, \sigma^2)$, then the MLE's $\hat\mu = \overline{X}$ and $\hat\sigma^2 = \overline{X^2} - (\overline{X})^2$ are independent and

$\dfrac{\sqrt{n}(\hat\mu - \mu)}{\sigma} \sim N(0,1)$, $\qquad \dfrac{n\hat\sigma^2}{\sigma^2} \sim \chi^2_{n-1}$

Note: We know the complete joint distribution of $\hat\mu$ and $\hat\sigma^2$.

OK! So, let's observe a sample $X_1, \ldots, X_n$ with distribution $\mathbb{P}_{\theta_0}$ from a parametric family $\{\mathbb{P}_\theta : \theta \in \Theta\}$ and $\theta_0$ is unknown

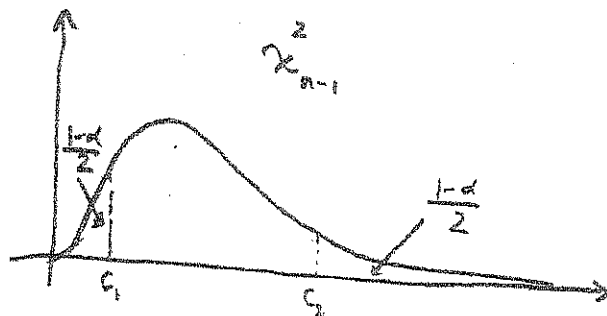<u>Given</u> a <u>confidence parameter</u> $\alpha \in [0,1]$ (usually $\alpha = 0.95$), if there exist two <u>statistics</u> $S_1 = S_1(X_1, \ldots, X_n)$ & $S_2 = S_2(X_1, \ldots, X_n)$ such that probability $P_{\theta_0}(S_1 \leq \theta_0 \leq S_2) \geq \alpha$, then we call $[S_1, S_2]$ <u>a confidence interval</u> for the unknown parameter $\theta_0$.

Let $X_1, \ldots, X_n$ i.i.d $N(\mu, \sigma^2)$; $\mu, \sigma^2$ unknown.

<u>Corollary 1</u> $\Rightarrow$ $A = \dfrac{\sqrt{n}(\hat{\mu} - \mu)}{\sigma} \sim N(0,1)$ and $B = \dfrac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$,

$A$ & $B$ independent

So, we can represent $A$ & $B$ as $A = Y_1$, $B = Y_2^2 + \ldots + Y_n^2$ for some $Y_1, \ldots, Y_n$

i.i.d. $N(0,1)$

Choose points $c_1, c_2$ so that $P(c_1 \leq B \leq c_2) = \alpha$, i.e. the area btw. $c_1$ and $c_2$ is $\alpha$, i.e. the area in each tail is $(1-\alpha)/2$. (see the <u>Note</u> on the bottom)



For these values of $c_1$ and $c_2$, we can GUARANTEE w/ confidence $\alpha$ that

$$c_1 \leq B = \frac{n\hat{\sigma}^2}{\sigma^2} \leq c_2$$

Solve this for the unknown parameter $\sigma^2$ $\Rightarrow$ the $\alpha$-confidence interval for $\sigma^2$

is $\left[ \dfrac{n\hat{\sigma}^2}{c_2}, \dfrac{n\hat{\sigma}^2}{c_1} \right]$ where $c_1, c_2$ are such that

$$P(c_1 \leq Z \leq c_2) = \alpha \quad \text{where } Z \sim \chi^2_{n-1}$$

<u>Note</u> <u>Def</u> the $q^{th}$ <u>quantile</u> $u$ of a probability distribution of a continuous random variable $U$ is defined by $P(U \leq u) = q$.

So, $c_1$ is nothing else but the $(1-\alpha)/2$-quantile of $\chi^2_{n-1}$; while $c_2$ is the $\frac{1+\alpha}{2}$-quantile. Sometimes one writes $c_1 \equiv \chi^2_{n-1; (1-\alpha)/2}$ and $c_2 = \chi^2_{n-1; (1+\alpha)/2}$

Next, let's find an $\alpha$-confidence interval for $\mu$.

Consider $\dfrac{A}{\sqrt{\frac{B}{n-1}}} = \dfrac{Y_1}{\sqrt{\frac{1}{n-1}(Y_2^2 + \ldots + Y_n^2)}} \sim t_{n-1}$ ← student $t$-distr. w/ $n-1$ degrees of freedom
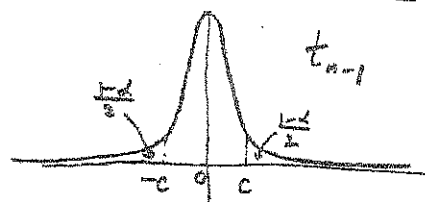
But, $\dfrac{A}{\sqrt{\frac{B}{n-1}}} = \dfrac{\sqrt{n} \cdot \frac{\hat{\mu} - \mu}{\sigma}}{\sqrt{\frac{1}{n-1} \cdot \frac{n \hat{\sigma}^2}{\sigma^2}}} = \underbrace{\dfrac{\sqrt{n-1}}{\hat{\sigma}}}_{\text{a constant}} (\hat{\mu} - \mu)$

So, choose $c$ such that the area in each tail of the $t_{n-1}$-distr. is $\frac{1-\alpha}{2}$.

Then w/ probability $\alpha$, we have

$$-c \leq \frac{\sqrt{n-1}}{\hat{\sigma}} (\hat{\mu} - \mu) \leq c.$$



So, the $\alpha$-confidence interval for $\mu$ is $\left[ \hat{\mu} - c \dfrac{\hat{\sigma}}{\sqrt{n-1}} \ , \ \hat{\mu} + c \dfrac{\hat{\sigma}}{\sqrt{n-1}} \right]$

**Example** sample of size $n=10$ from $N(\mu, \sigma^2)$; $\mu, \sigma^2$ unknown

$$X \begin{bmatrix} 0.86 \\ 1.53 \\ 1.57 \\ 1.81 \\ 0.99 \\ 1.09 \\ 1.29 \\ 1.78 \\ 1.29 \\ 1.58 \end{bmatrix}$$

$\hat{\mu} = \bar{X} = 1.379$

$\hat{\sigma}^2 = \overline{X^2} - (\bar{X})^2 = 0.0866$

choose $\alpha = 0.95$ (95%)

Find $c$ such that $t_9(-\infty, c) = 0.975 \rightarrow c = 2.262$

find $c_1, c_2$ s.t. $\chi_9^2([0, c_1]) = 0.025 \rightarrow c_1 = 2.7$

$\chi_9^2([0, c_2]) = 0.975 \rightarrow c_2 = 19.02$

So, w/ prob. 95% we can guarentee that

$$\bar{X} - c \sqrt{\tfrac{1}{9}(\overline{X^2} - (\bar{X})^2)} \leq \mu \leq \bar{X} + c \sqrt{\tfrac{1}{9}(\overline{X^2} - (\bar{X})^2)}$$

i.e. $1.1446 \leq \mu \leq 1.6134$

and similarly for $\sigma^2$: $0.0508 \leq \sigma^2 \leq 0.3573$

Testing hypotheses (about normal distr. for now)

Setup: Given an i.i.d. sample $(X_1, \ldots, X_n) \in S$ from $N(\mu, \sigma^2)$; $\mu, \sigma^2$ unknown
We need to decide between two hypotheses about the unknown parameters,
say $\mu$. Hypotheses will be one of the following:

① $\begin{cases} H_0 : \mu = \mu_0 \quad (\text{null}) \\ H_1 : \mu \neq \mu_0 \quad (\text{alternative}) \end{cases}$  OR ② $\begin{cases} H_0 : \mu \geq \mu_0 \\ H_1 : \mu < \mu_0 \end{cases}$  -OR ③ $\begin{cases} H_0 : \mu \leq \mu_0 \\ H_1 : \mu > \mu_0 \end{cases}$

Want to construct a TEST $\delta : S \to \{H_0, H_1\}$ that given an i.i.d
$(X_1, \ldots, X_n) \in S$ either accepts $H_0$ or rejects $H_0$ (i.e. accepts $H_1$)

Choose $\alpha \in (0,1)$ : level of significance for $\delta$ and we want to devise $\delta$
so that $\delta$ rejects $H_0$ when it's true w/ prob $\leq \alpha$, i.e.
$$\mathbb{P}(\delta = H_1 | H_0) \leq \alpha.$$
$$\|$$
$$\sup_{(\mu, \sigma^2) \in \Omega_0} \mathbb{P}(\delta = H_1 | \mu, \sigma^2)$$

usually 0.05

$\Omega_0$ : subset of parameter space $\Omega$ s.t. $H_0 : \theta \in \Omega_0$
unknown parameter

## 4.1. Hypothesis about mean of one normal sample       (t-test in Matlab)

We know that $\sqrt{n-1} \cdot \dfrac{\hat{\mu} - \mu}{\hat{\sigma}} \sim t_{n-1}$

Consider a t-statistic $\boxed{T = \sqrt{n-1} \cdot \dfrac{\hat{\mu} - \mu_0}{\hat{\sigma}}}$. It behaves differently depending

(another random var)

on whether the true unknown mean $\mu = \mu_0$, $\mu < \mu_0$ or $\mu > \mu_0$. Why?

If $\mu = \mu_0$, then $T \sim t_{n-1}$. If $\mu < \mu_0$ then

$$T = \sqrt{n-1} \cdot \frac{\hat{\mu} - \mu}{\hat{\sigma}} + \sqrt{n-1} \cdot \frac{\mu - \mu_0}{\hat{\sigma}} \to -\infty \quad \text{since the first term has}$$

$t_{n-1}$-distribution and the second one goes to $-\infty$. Similarly, if $\mu > \mu_0$, then $T \to +\infty$.

Idea: base the tests on this statistics → so called t-tests

① $(H_0: \mu = \mu_0)$ The indication that $H_0$ is not true would be if $|T|$ becomes too large, i.e. $T \to \pm\infty$.

So, let's devise the following test: $\delta = \begin{cases} H_0, & \text{if } |T| \leq c \\ H_1, & \text{if } |T| > c \end{cases}$

What is $c$? It depends on the level of significance $\alpha$. We want.

$$\alpha \geq \mathbb{P}(\delta = H_1 | H_0) = \mathbb{P}(|T| > c | H_0) = t_{n-1}(|T| > c) = 2 t_{n-1}(c, \infty) = \alpha$$

given that $H_0$ holds, we have $T \sim t_{n-1}$

So, from $2 t_{n-1}(c, \infty) = \alpha$, you find $c$.

② $(H_0: \mu \geq \mu_0)$ The indication that $H_0$ is not true would be if $T \to -\infty$

So, $\delta = \begin{cases} H_0, & \text{if } T \geq c \\ H_1, & \text{if } T < c \end{cases}$  What is $c$? Depends on $\alpha$ again.

$$\alpha \geq P(\delta = H_1 | H_0) = \mathbb{P}(T < c | H_0) =$$

$$= \mathbb{P}\left(T - \sqrt{n-1} \cdot \frac{\mu - \mu_0}{\hat{\sigma}} < c - \sqrt{n-1} \cdot \frac{\mu - \mu_0}{\hat{\sigma}} \Big| H_0\right) =$$

$$= \sup_{\mu \geq \mu_0} \mathbb{P}\left(T - \sqrt{n-1} \cdot \frac{\mu - \mu_0}{\hat{\sigma}} < c - \sqrt{n-1} \cdot \frac{\mu - \mu_0}{\hat{\sigma}}\right) =$$

this is $t_{n-1}$-distributed          this is maximized when $\mu = \mu_0$

$$= t_{n-1}((-\infty, c]) = \alpha.$$

So, all you need to do is find $c$ so that $t_{n-1}(-\infty, c) = \alpha$, and you have the t-test to test ②
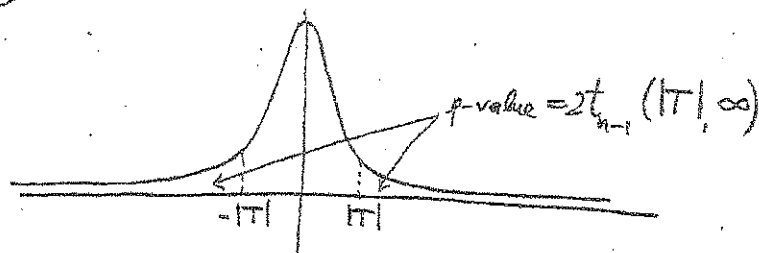
For ③ similar!

<u>p-value</u>: Rather than specifying $\alpha$ and deciding whether to accept or reject $H_0$ at level $\alpha$, we can ask "for what values of $\alpha$ do we reject $H_0$?"

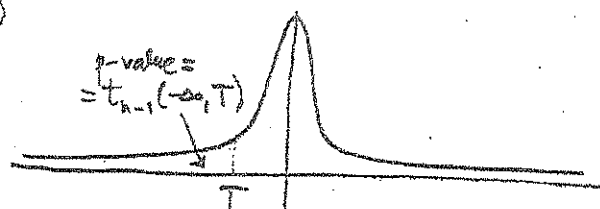p-value: the smallest value of $\alpha$ for which $H_0$ is rejected.

p-value can be understood as a probability, given that $H_0$ is true, to observe a value of $t$-statistic equally or less likely than the one that was observed. So, the small p-value means that the observed $t$-statistic is very unlikely under the null hypothesis, which in turn provides strong evidence against $H_0$.

① 



$$p\text{-value} = 2 t_{n-1}(|T|, \infty)$$

② 



$$p\text{-value} = t_{n-1}(-\infty, T)$$

Stated differently, to perform a test using a given sample, we first find the p-value of the sample and then $H_0$ is rejected if we decide to use $\alpha$ larger than the p-value; and accept otherwise

(in this case, under $H_0$, $T \to +\infty$)

p-value also tells us whether the decision to accept or reject $H_0$ is a close call.

## 4.2. Hypothesis about variance of one normal sample

We know that $\dfrac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$. So, similarly to 4.1, we base our tests on the following statistic:
$$Q = \frac{n\hat{\sigma}^2}{\sigma_0^2}$$

Since, $Q = \dfrac{n\hat{\sigma}^2}{\sigma^2} \cdot \dfrac{\sigma^2}{\sigma_0^2} \sim \dfrac{\sigma^2}{\sigma_0^2}\chi^2_{n-1}$, then $Q$ behaves differently depending on whether $\sigma = \sigma_0$, $\sigma > \sigma_0$ or $\sigma < \sigma_0$ (exactly what we need!)

① $(H_0 : \sigma = \sigma_0)$  The decision rule will be
$$\delta = \begin{cases} H_0, & \text{if } c_1 \leq Q \leq c_2 \\ H_1, & \text{if } Q < c_1 \text{ or } c_2 < Q. \end{cases}$$

Thresholds $c_1, c_2$ should satisfy the condition

$$\alpha = \mathbb{P}(\delta = H_1 | H_0) = \mathbb{P}(Q < c_1 | \sigma = \sigma_0) + \mathbb{P}(Q > c_2 | \sigma = \sigma_0) =$$

$$= \chi^2_{n-1}(0, c_1) + \chi^2_{n-1}(c_2, \infty)$$

So, for example, you can set $\chi^2_{n-1}(0, c_1) = \chi^2_{n-1}(c_2, \infty) = \alpha/2$.

② $(H_0 : \sigma \leq \sigma_0)$ In this case, the decision rule will be

$$\delta = \begin{cases} H_0, & \text{if } Q \leq c \\ H_1, & \text{if } Q > c. \end{cases}$$

Threshold $c$ should satisfy the condition

$$\alpha = \mathbb{P}(\delta = H_1 | H_0) = \sup_{\sigma \leq \sigma_0} \mathbb{P}(Q > c) = \sup_{\sigma \leq \sigma_0} \mathbb{P}\left( \frac{n \hat{\sigma}^2}{\sigma_0^2} > c \right) =$$

$$= \sup_{\sigma \leq \sigma_0} \mathbb{P}\left( \frac{n \hat{\sigma}^2}{\sigma^2} > \underbrace{\frac{\sigma_0^2}{\sigma^2} c}_{\substack{\text{make this as} \\ \text{small as possible} \\ \text{under } \sigma \leq \sigma_0}} \right) = \mathbb{P}\left( \underbrace{\frac{n \hat{\sigma}^2}{\sigma^2}}_{\chi^2_{n-1}\text{-distributed}} > c \right) = \chi^2_{n-1}(c, \infty)$$

So, find $c$ from $\alpha = \chi^2_{n-1}(c, \infty)$

③ $(H_0 : \sigma \geq \sigma_0)$ is similar.

## 4.3 Two-sample t-tests

### 4.3.1 Paired samples

Suppose we wish to compare returns on small-cap vs. large-cap stocks.
For each of $n$ yrs we have the returns on a portfolio of small-cap stocks $(x_1, \ldots, x_n)$
and on a portfolio of large-cap stocks $(y_1, \ldots, y_n)$. Form $z_i = x_i - y_i$, $i = \overline{1, n}$.
We want to test $H_0 : \mu_x = \mu_y$ for the means of the two samples.
Assuming that $x_i, y_i$ are normal and independent, so then $z_i$ will be normal.
So, $H_0$ is equivalent to $H_0 : \mu_z = 0$. Now, do the usual t-test from 4.1 ① for $z_1, \ldots, z_n$