*There will be an extra lecture next Wednesday evening through Webcast → more details soon

Def: rv $\underline{X}_1, ..., \underline{X}_n$ → sample mean of $(\underline{X}_1, ..., \underline{X}_n)$ is

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} \underline{X}_i$$

## Consequences continued

5) $\underline{X}_i \overset{iid}{\sim} N(\mu, \sigma^2)$ $(i = 1, ..., n)$     (All of them)

→ $\overline{\overline{X}}_n \sim N(\mu, \frac{\sigma^2}{n})$            (Each of them)

So $SD(\overline{X}_n) = \frac{\sigma}{\sqrt{n}}$

Because $E(\overline{X}_n) = \mu$, $\overline{\overline{X}}_n$ is an unbiased estimator of $\mu$.

Def: In frequentist statistics, the standard deviation (SD) of an estimator $\hat{\theta}(rv)$ of a parameter $\theta$ is called the standard error $SE(\hat{\theta})$ of $\hat{\theta}$

So if you use $\overline{X}_n$ as an estimate of $\mu$, $SE(\overline{X}_n) = \frac{\sigma}{\sqrt{n}}$
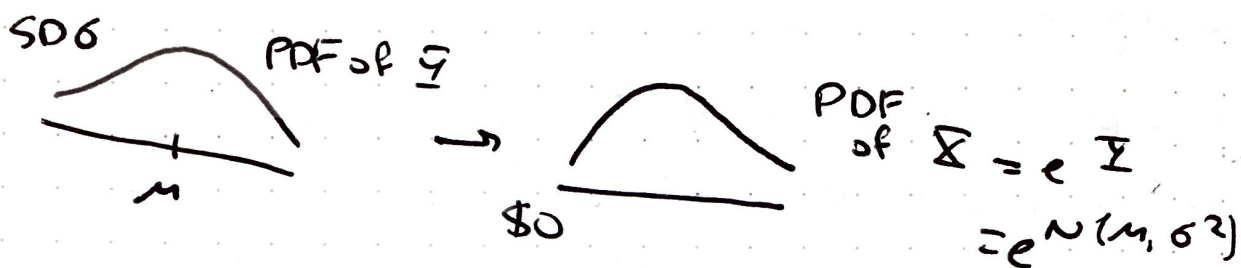
$SE \frac{\sigma}{\sqrt{n}}$

$SD = SE \sigma$


—PDF of $\overline{\overline{X}}_n$, $n>1$
—PDF of $\underline{X}_1$

This is the basis of most frequentist statistical inference

As $n \uparrow$, $\overline{\overline{X}}_n$ gets better as an estimate of $\mu$, at a $\frac{1}{\sqrt{n}}$ rate. This is called the Square Root Law.

This means that to cut the $SE(\bar{X}_n)$ in half, you have to quadruple the sample size.

## Lognormal Distribution

should be exponential normal

**Def:** If $X > 0$ and $Y = \log(X) \sim N(\mu, \sigma^2)$, then $X \sim$ Lognormal $(\mu, \sigma^2)$

SD $\sigma$



PDF of $Y$

$\mu$

$\longrightarrow$

\$0

PDF of $X = e^Y$
$= e^{N(\mu, \sigma^2)}$

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

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

Can get MGF of $X$ from MGF of $Y$

MGF of $Y$ is $\psi_Y(t) = \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$

But by def: $\psi_Y(t) = E(e^{tY}) = E(e^{t\log X}) = E(X^t)$

So we can read the moments of $X$ directly from the MGF of $Y$

$E(X) = \psi_Y(1) = \exp(\mu + \frac{\sigma^2}{2})$

$V(X) = \psi_Y(2) - [\psi_Y(1)]^2 = \exp(2\mu + \sigma^2)[e^{\sigma^2} - 1]$

# Gamma distribution

$(\alpha, \beta > 0)$ $X$ has the Gamma dist. with parameters $(\alpha, \beta)$, written $X \sim \Gamma(\alpha, \beta)$ or $X \sim \text{Gamma}(\alpha, \beta)$

$\rightarrow X$ continuous on $(0, \infty)$ with...

PDF $X \sim \Gamma(\alpha, \beta)$



PDF $f_X(x \mid \alpha, \beta) = \dfrac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \, \underbrace{I(x > 0)}_{\text{Support of } X}$

$\alpha$ is called a shape parameter in the $\Gamma(\alpha, \beta)$ family because it governs things like skewness of the dist.

$\beta$ is related to the scale of the distribution, which measures how spread out the distribution is

$\Gamma(\alpha)$ is the Gamma function, invented to deal with integrals of functions like ✴ above:

$$\Gamma(\alpha) \overset{\triangle}{=} \int_0^\infty x^{\alpha-1} e^{-x} \, dx$$

$\underbrace{\qquad}_{\uparrow}$

has no anti-derivative in closed form

$\Gamma(\alpha)$ turns out to be a continuous generalization of the factorial function, because (n positive integer)

$\rightarrow \Gamma(n) = (n-1)!$

$\Gamma(\alpha) \to \infty$ really quickly as $\alpha \to \infty$, so it's better to evaluate the Gamma PDF on the log scale and then exponentiate:

$$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} = \exp\left[\begin{array}{l} \alpha \ln(\beta) - \ln\Gamma(\alpha) \\ + (\alpha-1)\ln(x) - \beta x \end{array}\right]$$

Another way to tame $\Gamma(x)$ is with a Stirling's approximation

$$\Gamma(x) \doteq \sqrt{2\pi} x^{x-\frac{1}{2}} e^{-x} \text{ for large } x$$

so that $\ln\Gamma(x) \doteq \frac{1}{2}\ln(2\pi) + (x-\frac{1}{2})\ln x - x$

$X \sim \Gamma(\alpha, \beta)$

$$\Psi_X(t) = \left(1 - \frac{t}{\beta}\right)^{-\alpha} \text{ for } t < \beta$$

so $E(X) = \frac{\alpha}{\beta}$ and $V(X) = \frac{\alpha}{\beta^2}$   $SD(X) = \frac{\sqrt{\alpha}}{\beta}$

## Alternate expression

$$\Psi_X(t) = \left(\frac{\beta}{\beta-t}\right)^\alpha \text{ for } t < \beta$$

Special case of $\Gamma(\alpha, \beta)$
with $\alpha = 1$ the PDF is $f_X(x|\beta) = \beta e^{-\beta x} I(x > 0)$

This is just the exponential distribution

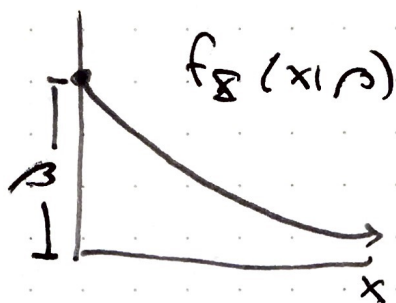$\underline{X} \sim \text{Exponential} (\beta)$

$\Psi_{\underline{X}}(t) = \dfrac{\beta}{\beta - t} , \quad 1 + t\beta$

$E(\underline{X}) = \dfrac{1}{\beta}$

$V(\underline{X}) = \dfrac{1}{\beta^2}$

$SD(\underline{X}) = \dfrac{1}{\beta}$

$f_{\underline{X}}(x|\beta)$



__Thm:__ Suppose that arrivals (events) occur according to a Poisson process with rate $\beta$ per unit time. and define $\underline{Z}_1 = Z_1 - 0$

$$\underline{Z}_2 = Z_2 - Z_1$$
$$\cdots \quad \underline{Z}_k = Z_k - Z_{k-1} \quad \text{for } k=2,3,\ldots$$

Set $Z_k = $ time until $k$th arrival $k=1,2,\ldots$

The $\underline{Z}_i$ are the inter-arrival times.

$\underline{Z}_i \overset{IID}{\sim} \text{Exponential} (\beta)$ as a result.

The exponential dist. is also related to the Geometric dist. in that they both have a <u>memoryless property</u>.

__Thm:__ $\underline{X} \sim \text{Exponential} (\beta)$ ; $t > 0$; $h > 0$

$\rightsquigarrow P(\underline{X} \geq t + h \mid \underline{X} \geq t) = P(\underline{X} > h)$

__Ex:__ $\underline{X} = $ time from initial use until a manufactured product fails
(e.g. light bulb)

$$F_{\underline{X}}(x) = P(\underline{X} \leq x)$$

$$1 - F_{\underline{X}}(x) = P(\underline{X} > x)$$

$$= P(\text{"system survives" at least to time } x)$$

For this reason, $1 - F_{\underline{X}}(x)$ is the survival function

$S_{\underline{X}}(x) = 1 - F_{\underline{X}}(x)$ in medicine and the reliability

function $R_{\underline{X}}(x) = 1 - F_{\underline{X}}(x)$ in engineering

For $\underline{X} \sim$ Exponential $(\beta) \rightarrow F_{\underline{X}}(x) = 1 - e^{-\beta x}$ for $x > 0$

So $S_{\underline{X}}(x) = R_{\underline{X}}(x) = e^{-\beta x}$ for this dist.

The instantaneous failure rate or hazard rate function
is defined to be $H_{\underline{X}}(x) = \dfrac{f_{\underline{X}}(x)}{S_{\underline{X}}(x)} = \dfrac{f_{\underline{X}}(x)}{R_{\underline{X}}(x)}$

This gives $P\left(\text{failure in interval } (x, x+\epsilon) \,\middle|\, \begin{array}{l}\text{survival to} \\ \text{time } x\end{array}\right)$
for small $\epsilon$

Notice that if $\underline{X} \sim$ Exponential $(\beta)$ then $H_{\underline{X}}(x) = \dfrac{\beta e^{-\beta x}}{e^{-\beta x}}$
$= \beta$ (constant in $x$)

The exponential is the only failure rate distribution with
constant hazard.

Returning to the earlier result that $\underline{X} \sim$ Exponential $(\beta)$

$\rightarrow P(\underline{X} \geq t + h \mid \underline{X} \geq t) = P(\underline{X} \geq h)$

For all $t > 0$  $h > 0$

This says that if the product has survived to time $t$, the chance it will survive to time $(t+h)$ is the same as the original chance of surviving from time 0 to time $h$

"the system doesn't remember how long it survived" (This often makes the Exponential unrealistic in practice)

## Consequences

1) $X_i \overset{IID}{\sim}$ Exponential $(\beta)$ $(i = 1, \ldots, n)$ then
$$Y_1 = \min(X_1, \ldots, X_n) \sim \text{Exponential}(n\beta)$$

## Beta Distribution

$\alpha, \beta > 0$

$$X \sim \text{Beta}(\alpha, \beta) \iff f_X(x) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$$

$\underbrace{(0 < x < 1)}_{\text{support of } X}$

The name comes from the normalizing constant: the function $x^{\alpha-1}(1-x)^{\beta-1}$ has no closed-form anti-derivative, so people just made a definition...

Def: for all $\alpha > 0$ $\beta > 0$
$$B(\alpha, \beta) \overset{\Delta}{=} \int_0^1 x^{\alpha-1}(1-x)^{\beta-1} dx$$
$\uparrow$
beta function

Can show that $B(\alpha, \beta) = \dfrac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$

$(\alpha, \beta)$ jointly control the shape of the Beta$(\alpha, \beta)$ dist.

$\underline{X} \sim$ Beta$(\alpha, \beta)$

$$\Psi_{\underline{X}}(t) = 1 + \sum_{k=1}^{\infty} \left( \prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r} \right) \cdot \frac{t^k}{k!}$$

$E(\underline{X}) = \dfrac{\alpha}{\alpha+\beta}$

$V(\underline{X}) = \left( \dfrac{\alpha}{\alpha+\beta} \right) \left( \dfrac{\beta}{\alpha+\beta} \right) \left( \dfrac{1}{\alpha+\beta+1} \right)$

## Multinomial Distributions (back to discrete)

You're contemplating a population that contains
elements of $k \geq 2$ types

(e.g. $\{$ Democrat, Republican, Libertarian, Independent, Green$\}$)

Suppose the proportion of elements of type $i$ is $0 \leq p_i \leq 1$
with $\sum_{i=1}^{k} p_i = 1$ $\underset{\sim}{p} = (p_1, \ldots, p_k)$

You take an IID sample of size $n$ from this pop.
$\underline{X}_i = \#$ elements of type $i$ in your sample

$\sum_{i=1}^{k} \underline{X}_i = n$

Can show that the vector $\underset{\sim}{\underline{X}} = (\underline{X}_1, \ldots, \underline{X}_k)$ has pmf:

$$f_{\underset{\sim}{\underline{X}} | n, p}(\underline{x} | n, p) = \begin{cases} \binom{n}{x_1, \ldots, x_k} p_1^{x_1} \cdots p_k^{x_k} & \text{if } \sum_{i=1}^{k} x_i = n \\ 0 & \text{else} \end{cases}$$

where $\binom{n}{x_1, \ldots, x_k} \triangleq \dfrac{n!}{x_1! \, x_2! \ldots x_k!}$ is the multinomial coefficient

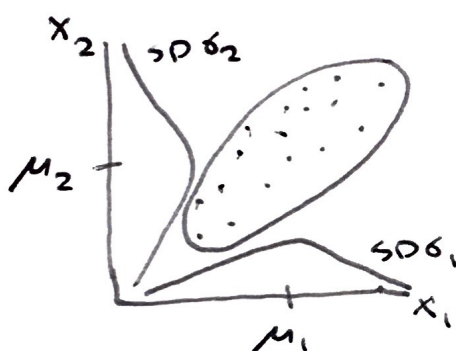This is called the Multinomial $(n, p)$ distribution

$E(X_i) = n p_i$ \quad $\left.\begin{array}{l} \\ \\ \end{array}\right\}$ just like binomial

$V(X_i) = n p_i (1 - p_i)$

$C(X_i, X_j) = - n p_i p_j$

Negatively correlated because $\sum\limits_{i=1}^{k} X_i = n$

## Bivariate Normal Dist.

Can build a 2-dimensional (bivariate) version of the Normal dist. as follows:



$Z_1, Z_2 \overset{IID}{\sim} N(0, 1)$

specify 5 parameters

1) $-\infty < \mu_1 < \infty$

2) $-\infty < \mu_2 < \infty$

3) $0 < \sigma_1 < \infty$

4) $0 < \sigma_2 < \infty$

5) $-1 < \rho < 1$

Now build $(X_1, X_2)$ with the transformation

$X_1 = \mu_1 + \sigma_1 Z_1$

$X_2 = \sigma_2 \left[ \rho Z_1 + \sqrt{1 - \rho^2} \, Z_2 \right] + \mu_2$

The joint PDF of $\underset{\sim}{X} = (X_1, X_2)$ is then

$f_{X_1, X_2}(x_1, x_2) = \dfrac{1}{2 \pi \sqrt{1 - \rho^2} \, \sigma_1 \sigma_2} \cdot \exp \left\{ \dfrac{1}{2(1 - \rho^2)} \left[ \left( \dfrac{x_1 - \mu_1}{\sigma_1} \right)^2 - \right.\right.$

$\left.\left. 2 \rho \left( \dfrac{x_1 - \mu_1}{\sigma_1} \right) \left( \dfrac{x_2 - \mu_2}{\sigma_2} \right) + \left( \dfrac{x_2 - \mu_2}{\sigma_2} \right)^2 \right] \right\}$ \quad standard units

This is the Bivariate Normal $(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$ dist

Easy to show that $E(X_1) = \mu_1$, $E(X_2) = \mu_2$

$V(X_1) = \sigma_1^2$, $V(X_2) = \sigma_2^2$, $\rho(X_1, X_2) = \rho$

## Consequences of this def:

1) $(X_1, X_2) \sim$ Bivariate Normal $\rightarrow \binom{X_1, X_2}{\text{independent}} \leftrightarrow \binom{X_1, X_2}{\text{uncorrelated}}$

We already knew the $\rightarrow$ direction in general; what's new here is that correlation 0 implies independence if $(X_1, X_2) \sim$ Bivariate Normal

2) $(X_1, X_2) \sim$ Bivariate Normal $(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) \rightarrow$ conditional distribution of $X_2$ given that $X_1 = x_1$ is (univariate) normal with mean $E(X_2 | x_1) =$

$\mu_2 + \rho \frac{\sigma_2}{\sigma_1}(x_1 - \mu_1)$ and variance $V(X_2 | x_1) = (1 - \rho^2)\sigma_2^2$



Galton revisited

Result 2 says that if $(X_1, X_2)$ are Bivariate Normal then the conditional distributions of $X_2$ given $X = x_1^*$ in all of the vertical strips are also normal

And the means of all these normal distributions in the vertical strips are connected together by Galton's regression line

$\hat{X}_2 = \mu_2 + \rho \frac{\sigma_2}{\sigma_1}(x_1 - \mu_1)$          $\hat{X}_2 = \beta_0 + \beta_1 x_1$

This line has slope $\beta_1 = \rho \frac{\sigma_2}{\sigma_1}$ and "y"-intercept $\beta_0 = \mu_2 - \beta_1 \mu_1$

Moreover, we can now quantify an earlier insight

$$\boxed{\text{ignore } x}, \quad \text{predict } (\hat{x}_2)_{\substack{no \\ x_1}} = \mu_2 = E(\mathcal{X}_2)$$

(root mean squared error) (RMSE) of this prediction is

$$\sqrt{V(\mathcal{X}_2)} = \sigma_2$$

$$\boxed{\text{use } x_1 \text{ to predict } \hat{x}_2} \quad \text{predict } (\hat{x}_2)_{\substack{use \\ x_1}} = E(\mathcal{X}_2 \mid \mathcal{X}_1 = x_1)$$

$$= \mu_2 + \rho \frac{\sigma_2}{\sigma_1}(x_1 - \mu_1)$$

RMSE of this prediction is $\sqrt{V(\mathcal{X}_2 \mid x_1)} = \sigma_2 \sqrt{1-\rho^2}$

Since $-1 < \rho < 1$, $\sigma_2\sqrt{1-\rho^2} \leq \sigma_2$ with equality only when $\rho=0$

3) $(\mathcal{X}_1, \mathcal{X}_2) \sim \text{Bivariate Normal}(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$,

$$Y = a_1 \mathcal{X}_1 + a_2 \mathcal{X}_2 + b, \quad (a_1, a_2, b) \text{ arbitrary constants}$$

$$\rightarrow Y \sim N(a_1\mu_1 + a_2\mu_2 + b, \; a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + 2a_1 a_2 \rho\sigma_1\sigma_2)$$

## Large random samples

You draw an IID random sample $\mathcal{X}_1, \dots, \mathcal{X}_n$ from a population, with the goal of estimating the population mean $\mu = E(\mathcal{X}_i)$

We've already seen that from a root mean squared error point of view, the sample mean $\bar{\mathcal{X}}_n = \frac{1}{n}\sum_{i=1}^{n} \mathcal{X}_i$ is the best you can do (in the absence of prior information)

It would be nice if $\bar{\mathcal{X}}_n$ approached the right answer $\mu$ as $n$ increases; how to quantify that idea?

# Two inequalities that help

## Markov inequality

Suppose $\underline{X}$ is a non-negative rv, ie $P(\underline{X} \geq 0) = 1$

Then for all real $t > 0$, $P(\underline{X} \geq t) \leq \dfrac{E(\underline{X})}{t}$ says that if

$E(\underline{X})$ is fixed, you can't move more and more probability out into the right tail beyond a certain point

Ex: $E(\underline{X}) = 1$, $\underline{X}$ non-negative $\rightarrow P(\underline{X} \geq 100) \leq \dfrac{1}{100}$

The inequality is sharp, meaning that the upper bound $\dfrac{E(\underline{X})}{t}$ on $P(\underline{X} \geq t)$ is attainable but most of the time its a crude upper bound

Ex: $E(\underline{X}) = 1$, $\underline{X}$ —non-negative $\rightarrow$ put probability $0.99$ on $\underline{X} = 0$ and probability $0.01$ on $\underline{X} = 100$