

LECTURE 25 : THURSDAY APRIL 21

Part 1 : Finish §10.2, expectation & variance

Part 2 : Two special distributions

Part 3 : Bias and variance

Part 4 : The German Tank problem

We will do
two on
Tuesday

PART 1. $f(x) \geq 0$ s.t. $\int_{-\infty}^{\infty} f(x) dx = 1$ \rightarrow density function
under density curve areas describe prob for a continuous r.v.

Let X be a r.v. with the density function $f(x)$ described below:

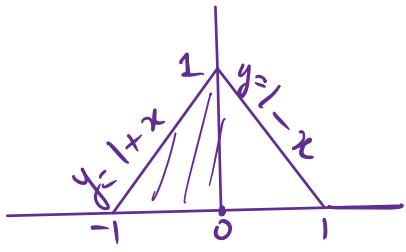
$$f(x) = \begin{cases} 0, & x < -1 \\ x+1, & -1 \leq x < 0 \\ 1-x, & 0 \leq x < 1 \\ 0, & x \geq 1 \end{cases}$$

Find ① $P\left(\frac{1}{2} < X < 1\right) = \frac{1}{2} b h = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$

② $P\left(0 < X < \frac{1}{2}\right) = \frac{1}{2} - \frac{1}{8} = \frac{4-1}{8} = \frac{3}{8}$
 $= \text{area of } \square + \text{area of } \triangle = \left(\frac{1}{2}\right)^2 + \frac{1}{8} = \frac{3}{8}$

$F(x) = \text{Area under the curve up to } x$

$$= \int_{-\infty}^x f(t) dt = P(X \leq x)$$



$$F(x) = \begin{cases} 0, & x < -1 \\ \int_{-1}^x (1+t) dt, & -1 \leq x < 0 \\ \frac{1}{2} + \int_0^x (1-t) dt, & 0 \leq x \leq 1 \\ 1, & x > 1 \end{cases}$$

$$\begin{aligned} \int_{-1}^x (1+t) dt &= t + \frac{t^2}{2} \Big|_{-1}^x \\ &= x + \frac{x^2}{2} - \left(-1 + \frac{1}{2} \right) = x + \frac{x^2}{2} + \frac{1}{2} = \frac{x^2+2x+1}{2} \end{aligned}$$

$$\int_0^x (1-t) dt = t - \frac{t^2}{2} \Big|_0^x = x - \frac{x^2}{2} = \frac{2x-x^2}{2}$$

$$F(x) = \begin{cases} 0, & x < -1 \\ \frac{x^2+2x+1}{2}, & -1 \leq x < 0 \\ \frac{1}{2} + \frac{2x-x^2}{2} = \frac{1+2x-x^2}{2}, & 0 \leq x \leq 1 \\ 1, & x > 1 \end{cases}$$

$$\begin{aligned} P\left(\frac{1}{2} < X < 1\right) &= F(1) - F\left(\frac{1}{2}\right) \\ &= \frac{1+2\cdot\frac{1}{2}}{2} - \frac{\frac{1+2\cdot\frac{1}{2}-\frac{1}{4}}{2}}{2} \Bigg\} \frac{7}{8} \\ &= 1 - \frac{7}{8} = \frac{1}{8} \end{aligned}$$

EXPECTED VALUE & VARIANCE

- $P(X \leq x) = \int_{-\infty}^x f(t) dt$ Note
 $f(x) \neq P(X=x)$
 for continuous r.v. X .
- $E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx = \mu$
- $E(g(X)) = \int_{-\infty}^{\infty} g(x) \cdot f(x) dx$

• $\text{Var}(X) = E[(X - \mu)^2]$
 $g(X) = (X - \mu)^2$, $\text{Var}(X) = E(g(X))$

$$\begin{aligned}\text{Var}(X) &= \int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx \\ &= \int_{-\infty}^{\infty} [x^2 - 2\mu x + \mu^2] \underline{f(x)} dx \\ &= \int_{-\infty}^{\infty} x^2 f(x) dx - 2\mu \underbrace{\int_{-\infty}^{\infty} x f(x) dx}_{=\mu} + \mu^2\end{aligned}$$

$$\text{Var}(X) = E(X^2) - \mu^2 = \int_{-\infty}^{\infty} x^2 \cdot f(x) dx - \left(\int_{-\infty}^{\infty} x f(x) dx \right)^2$$

PROPERTIES

$$\textcircled{1} \quad E(aX+b) = aE(X)+b$$

$$\int c f(x) dx = c \int f(x) dx$$

$$\textcircled{2} \quad Var(aX+b) = a^2 Var(X)$$

$$SD(aX+b) = |a| SD(X)$$

$$\textcircled{3} \quad E(aX+bY) = aE(X) + bE(Y)$$

\textcircled{4} If for every interval A, B

$$P(X \in A, Y \in B) = P(X \in A) P(Y \in B)$$

then we say that X & Y are independent r.v.

\textcircled{5} If X, Y are independent r.v.

$$\text{then } Var(X+Y) = Var(X) + Var(Y)$$

$$Var(X-Y) = Var(X) + Var(Y)$$

PART 2 : TWO SPECIAL DISTRIBUTIONS

• The Exponential distribution

$$f(x) = e^{-x}, x > 0 \quad (0 \text{ elsewhere})$$

• The r.v. X that has the pdf $f(x) = e^{-x}, x > 0$

is called an EXPONENTIAL r.v. WITH RATE 1

$$X \sim \exp(1)$$

• In general, we say $X \sim \exp(\lambda)$

$$\text{If } f(x) = \lambda e^{-\lambda x}, x > 0, \lambda > 0 \quad (\lambda \text{ is some constant})$$

$$F(x) = \int_{-\infty}^x f(t) dt = \int_0^x \lambda e^{-\lambda t} dt$$

$$= \lambda \int_0^x e^{-\lambda t} dt = \lambda \cdot \frac{-e^{-\lambda t}}{\lambda} \Big|_0^x = 1 - e^{-\lambda x}$$

$$F(x) = P(X \leq x) = 1 - e^{-\lambda x}, x > 0$$

$$(F(x) = 0, x \leq 0)$$

We often denote the exponential r.v. by T .

$$T \sim \exp(\lambda), \quad P(T \leq x) = 1 - e^{-\lambda x}, \quad x > 0, \lambda > 0$$

$$\begin{aligned} P(T > x) &= 1 - P(T \leq x) \\ &= 1 - (1 - e^{-\lambda x}) = e^{-\lambda x} \end{aligned}$$

MEMORYLESS PROPERTY OF $T \sim \exp(\lambda)$

$$P(T > s+t | T > t) = P(T > s)$$

$$\mathbb{E}(T) = \frac{1}{\lambda} \leftarrow \text{use integration by parts}$$

$$\mathbb{E}(T) = \int_{-\infty}^{\infty} x \cdot f(x) dx = \int_0^{\infty} x \cdot \lambda e^{-\lambda x} dx$$

$$\mathbb{E}(T) = \underbrace{\cancel{x \cdot \lambda e^{-\lambda x}}_0}_{0} + \int_0^{\infty} \frac{x e^{-\lambda x}}{f(\lambda)} dx$$

$$\mathbb{E}(T) = \int_0^{\infty} e^{-\lambda x} dx = \frac{e^{-\lambda x}}{-\lambda} \Big|_0^{\infty} = \frac{1}{\lambda}$$

$$\text{Var}(T) = \mathbb{E}(T^2) - (\mathbb{E}(T))^2$$

$$= \underbrace{\mathbb{E}(T^2)}_{\frac{2}{\lambda^2}} - \frac{1}{\lambda^2}$$

$$\text{Var}(T) = \frac{1}{\lambda^2}, \quad \text{SD}(T) = \frac{1}{\sqrt{\lambda}}$$

$$T \sim \exp(\lambda) , f(t) = \lambda e^{-\lambda t}, t > 0, \lambda > 0$$

$$F(t) = \begin{cases} 1 - e^{-\lambda t}, & t > 0 \\ 0, & t \leq 0 \end{cases}$$

$$\mathbb{E}(T) = \frac{1}{\lambda}$$

$$\text{Var}(T) = \frac{1}{\lambda^2}, \text{SD}(T) = \frac{1}{\sqrt{\lambda}}$$

$P(T > t) = S(t)$, S is the survival function

T models (randomly distributed) lifetime

S is the chance of lasting past time t

$$P(T > t+s | T > t) = P(\underline{T > s}) \quad \text{the object}$$

This means that the chance λ will get past s more time units if it has already lasted t time units is the same as the object lasting s time units when new.

Object forgets that it has lasted for t time units.

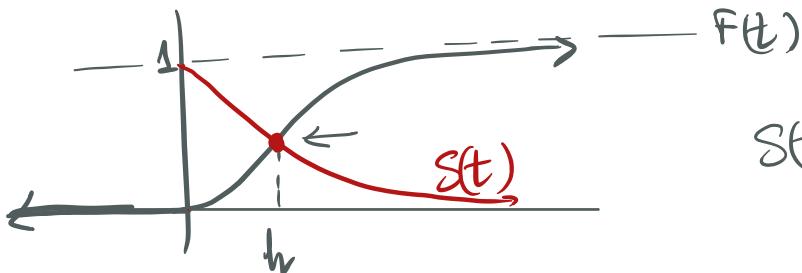
Median of the exponential distribution:

Median is the point where half goes to left

& half to the right

$$P(T \leq t) = P(T < t) = \underbrace{P(T > t)}_{e^{-\lambda t}}$$

$$F(t) = P(T \leq t) = 1 - e^{-\lambda t}, T \sim \exp(\lambda)$$



$$\begin{aligned}S(t) &= P(T > t) \\&= 1 - F(t)\end{aligned}$$

$$\text{at } h : S(h) = F(h)$$

$$e^{-\lambda h} = 1 - e^{-\lambda h}$$

$$2e^{-\lambda h} = 1$$

$$e^{-\lambda h} = \frac{1}{2}$$

h is the median or halfway point of the dsn.

$$e^{-\lambda h} = \frac{1}{2}$$

$$e^{\lambda h} = 2$$

$$\lambda h = \ln(2) \Leftrightarrow \log(2)$$

$$h = \frac{\log(2)}{\lambda} = \left(\frac{\ln(2)}{\lambda} \right) = \ln(2) \cdot E(T)$$

$$\ln(2) \approx 0.69 < 1$$

Median < $E(T)$ → Dsn is skewed right



HALF-LIFE

The half-life of a radioactive isotope is the time until half the atoms have decayed.

λ is called the decay rate & we use $T \sim \exp(\lambda)$ to model lifetime of radioactive atom

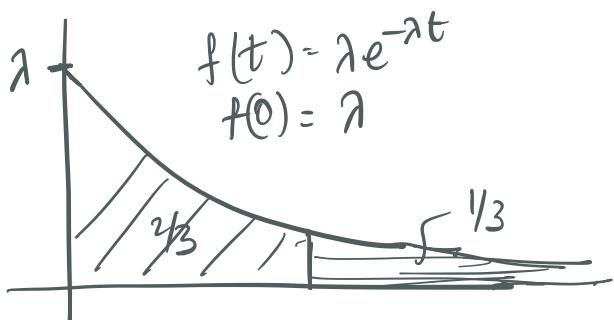
Ex 10.5.4 Strontium-90 which has a half life of 28.8 years.

If we assume exponential decay, how long will it take for $\frac{2}{3}$ of a lump of Strontium-90 to decay?

Half-life h is value s.t. $P(T < h) = P(T > h)$

$$h \text{ is given, } h = 28.8 \text{ years} = \frac{\log(2)}{\lambda}$$

$$\lambda = \frac{\log(2)}{h} \propto \frac{0.69}{28.8} \approx 0.024$$



want t s.t.
 $P(T > t) = \frac{1}{3}$

$$F(t) = \frac{2}{3}$$

$$1 - e^{-\lambda t} = \frac{2}{3}, \lambda = 0.024$$

$$\Rightarrow t = \frac{\log(3)}{\lambda} \approx 45.78 \text{ years.}$$

$$e^{-\lambda t} = \frac{1}{3} \Rightarrow \lambda t = 3, \lambda = \frac{\log(3)}{t}$$

This is the kind of computation in Carbon-dating which uses the proportion of Carbon 14 in organic matter to determine its age.

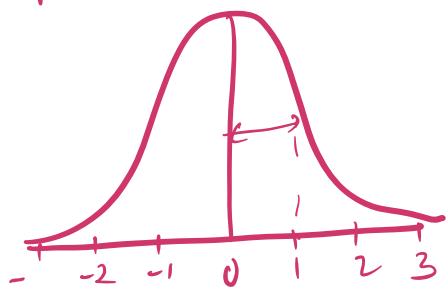
Look at $P(T > t) = \text{prop. of carbon}/4 \text{ left.}$

$$P = e^{-\lambda t}$$

$$t = \frac{1}{\lambda} \log\left(\frac{1}{P}\right)$$

THE NORMAL DISTRIBUTION 8.10.4

Well-known density curve $f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$, $x \in \mathbb{R}$



We can define a more general density function centered at $\mu \in \mathbb{R}$, spread $\sigma > 0$

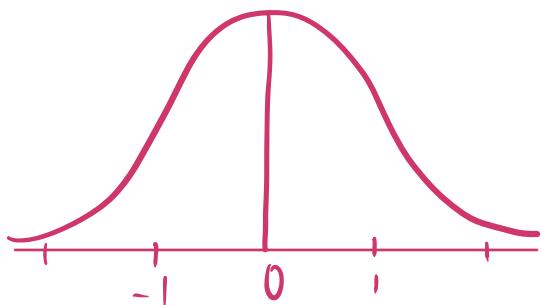
We say that X has the Normal dsn $N(\mu, \sigma^2)$

if the density of X is given by

$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}, -\infty < x < \infty$$

If $\mu=0$, $\sigma^2=1$ ($\sigma=1$), then we denote the associated r.v. by Z and call this the Standard Normal dsn.

$$F(x) = \Phi(x) = P(Z \leq x)$$

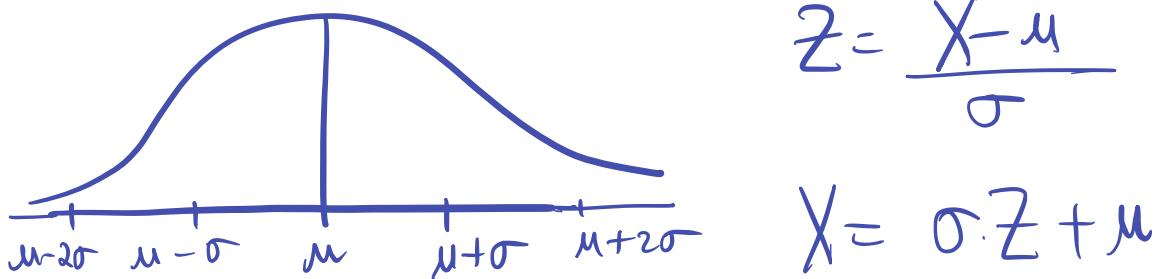


$$Z \sim N(0, 1)$$

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

If we put X in standard units

$$Z = \frac{X - \mu}{\sigma}$$



$$X = \sigma Z + \mu$$

$$\mathbb{E}(Z) = 0, \text{Var}(Z) = 1, SD(Z) = 1$$

$$\mathbb{E}(X) = \mu, \text{Var}(X) = \sigma^2, SD(X) = \sigma$$

IMPORTANT FACT

If X, Y are independent & normally distributed
 $X \sim N(\mu_X, \sigma_X^2)$, $Y \sim N(\mu_Y, \sigma_Y^2)$

$$X+Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

For instance

$$X + 3Y - 2 \sim N(\mu_X + 3\mu_Y - 2, \sigma_X^2 + 9\sigma_Y^2)$$

↑ Why?

$$\underline{X - Y} \sim N(\mu_X - \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

This r.v. is what we use to compute $\hat{\mu}$ confidence intervals for the difference in means.

Suppose we have $X_1, X_2, \dots, X_n \sim \text{iid } \mu_X, \sigma_X^2$

$Y_1, Y_2, \dots, Y_m, \mu_Y, \sigma_Y^2$

by CLT for n, m large enough.

$$\bar{X} \sim N\left(\mu_X, \frac{\sigma_X^2}{n}\right), \quad \bar{Y} \sim N\left(\mu_Y, \frac{\sigma_Y^2}{m}\right)$$

$$\bar{X} - \bar{Y} \sim N\left(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right)$$

Approx 95% CI for $\mu_X - \mu_Y$

$$(\bar{X} - \bar{Y}) \pm 2 \times SD(\bar{X} - \bar{Y})$$

$$: (\bar{X} - \bar{Y}) \pm 2 \times \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}$$

{ list of
"plausible" values
for $\mu_X - \mu_Y$

Setting up hyp. tests for equality of means

Test

$$H_0: \mu_x = \mu_y$$

$$H_A: \mu_x \neq \mu_y$$

$$\mu_x < \mu_y$$

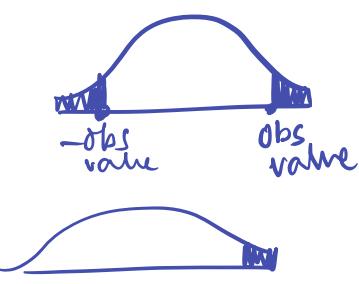
$$\mu_x > \mu_y$$

Statistic

$$|T| = |\bar{X} - \bar{Y}|$$

$$T = \bar{Y} - \bar{X}$$

$$T = \bar{X} - \bar{Y}$$



You are making a statement about $\mu_x - \mu_y$

Proportions C.I. for true difference in proportion

$$X_1, - - - X_n \sim \text{Bernoulli}(p_x)$$

$$Y_1, - - - Y_m \sim \text{Bernoulli}(p_y)$$

$$\bar{X} - \bar{Y} \underset{\text{CLT}}{\approx} N \left(p_x - p_y, \frac{p_x q_x}{n} + \frac{p_y q_y}{m} \right)$$

Difference for hypothesis tests.

$$H_0: p_x = p_y$$

H_A: p_x > p_y (Suppose we are doing a one-sided test)

$$T = \bar{X} - \bar{Y}$$

Under assumption H₀ is true, p_x = p_y = p \leftarrow common value

To estimate p , we combine or pool our samples.

X : $\boxed{0} \quad \boxed{0} \quad \boxed{0} \quad \boxed{1} \quad \dots \quad \boxed{1} \quad \dots \quad \dots$

Y : $\boxed{0} \quad \boxed{0} \quad \boxed{1} \quad \boxed{0} \quad \boxed{1} \quad \dots \quad \dots$

of successes in X sample = avg (total size)

$$\bar{X} = \frac{\text{# of successes}}{n}, \quad \sum X_i, \quad \bar{Y} = \frac{\sum Y_i}{m}$$

X_i, Y_i are 0-1 tkts values

$$\sum X_i = n\bar{X}, \quad \sum Y_i = m\bar{Y}$$

Total # of $\boxed{1}$'s in combined sample = $n\bar{X} + m\bar{Y}$

\hat{p} = estimate of p from the combined sample

$$= \frac{n\bar{X} + m\bar{Y}}{m+n}$$

total # of successes
or tkts marked
 $\boxed{1}$

sample size

Example

2 indep samples from Fresno ($n=400$)
& Irvine ($m=500$)
proportion in Fresno sample that agree with
lifiting mask mandate on airplanes: 57%
" Irvine " " : 61%

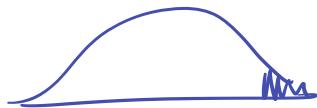
$$\bar{X} = 0.57, \bar{Y} = 0.61 \quad \sigma_{\bar{X}}^2 \approx \frac{\hat{p}_x(1-\hat{p}_x)}{n} = \frac{(0.57)(0.43)}{400}$$

$$SD(\bar{Y} - \bar{X}) = \sqrt{\sigma_{\bar{X}}^2 + \sigma_{\bar{Y}}^2} \quad \sigma_{\bar{Y}}^2 = \frac{\hat{p}_y(1-\hat{p}_y)}{m} = \frac{(0.61)(0.39)}{500}$$

$$(\bar{Y} - \bar{X}) \pm 2 \cdot SD(\bar{Y} - \bar{X})$$

$$H_0: p_x = p_y$$

$$H_1: p_x < p_y \quad T = \bar{Y} - \bar{X}$$



$$\bar{X} = 0.57$$

$$n\bar{X} = (0.57)(400), \quad (n\bar{Y}) = (0.61)(500)$$

$$\hat{p} = \frac{(0.57)(400) + (0.61)(500)}{400 + 500} \approx 0.5922$$

$$\sigma_{\bar{Y} - \bar{X}} = \sigma = \sqrt{\sigma_{\bar{X}}^2 + \sigma_{\bar{Y}}^2} = \sqrt{\frac{\hat{p}(1-\hat{p})}{400} + \frac{\hat{p}(1-\hat{p})}{500}}$$

$$Z = \frac{(\bar{Y} - \bar{X}) - 0}{\sigma}$$