# SAMPLING DISTRIBUTIONS AND THE CENTRAL LIMIT THEOREM

Normal and chi-squared distributions

# Moving towards statistical inference

- Statistical inference refers to making conclusions about a *population* based on *sample data*.

- Definition 1: Random variables $Y_1, Y_2, \ldots, Y_n$ constitute a RANDOM SAMPLE if they are independent and identically distributed (iid).

- A very common statistical inference problem involves estimating a *parameter* with a (sample) *statistic*.

- Definition 2: A STATISTIC is a function of observed random variables $Y_1, Y_2, \ldots, Y_n$.

# Sampling distributions

- The value of a statistic will depend on the specific sample observed.
- In order to perform inference, we need to know the distribution of a statistic.   *Random variable*
- Definition 3: The SAMPLING DISTRIBUTION of a statistic is the distribution of possible values of the statistic across many repeated samples.
- The sampling distribution naturally depends on the distribution of the sample data $Y_1, Y_2, \ldots, Y_n$.

# Sampling distributions for normally distributed data

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

- We will start by considering the situation in which $Y_1, Y_2, \ldots, Y_n \sim iid\ N(\mu, \sigma^2)$.

- Theorem 1: If $Y_1, Y_2, \ldots, Y_n \sim iid\ N(\mu, \sigma^2)$, then $\bar{Y} = \frac{1}{n}\sum_{i=1}^{n} Y_i$ is normally distributed with mean $\mu$ and variance $\frac{\sigma^2}{n}$.

Proof    Before (method of mgfs) showed:    MGFs

$$Y_1 + Y_2 + \ldots + Y_n = \sum_{i=1}^{n} Y_i \sim N(n\mu, n\sigma^2)$$

$$\bar{Y} = \frac{1}{n}\sum Y_i$$

$$Y \sim N(\mu, \sigma^2) \quad \text{what is the dist. of } \frac{Y}{a}?$$

$$U = \frac{Y}{a} \quad \text{dist. of } U?$$

mgf of $U$  $E[e^{tu}] = E[e^{t\frac{Y}{a}}] = E[e^{\frac{t}{a}Y}]$

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

$$= e^{\mu\frac{t}{a}} e^{\frac{1}{2}\sigma^2\frac{t^2}{a^2}} = e^{\frac{\mu}{a}t} e^{\frac{1}{2}\frac{\sigma^2}{a^2}t^2}$$

mgf of $N(\frac{\mu}{a}, \frac{\sigma^2}{a^2})$

mgf of $N(\mu, \sigma^2)$

$$e^{\mu t} e^{\frac{1}{2}\sigma^2 t^2}$$

$$E[e^{tY}]$$

$$Y \sim N(\mu, \sigma^2) \implies \frac{Y}{a} \sim N(\frac{\mu}{a}, \frac{\sigma^2}{a^2})$$

$$\sum_{i=1}^{n} Y_i \sim N(n\mu, n\sigma^2) \implies \overline{Y} = \frac{\sum Y_i}{n} \sim N(\mu, \frac{\sigma^2}{n})$$

$\square$

# Normally distributed data: Example

- Example 1: Australian men aged 50-60 currently using antihypertensive drugs have mean diastolic blood pressure 94.9 mm Hg and standard deviation 11.5 Hg.  Assuming that these measures are normally distributed, what is the probability that the average DBP of a sample of 8 men from this population will be at least 100 mm Hg?

$n = 8$

$M$

$P(\bar{Y} > 100)$

# Normally distributed data: Another example

- Example 2: In the previous example how large must the sample be for the mean to be within 2 mm Hg from the mean with 95% probability?

$$\frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

- Example 3: If $Z \sim N(0,1)$ find the distribution of $Z^2$.

Method of mgf

mgf of $Y = Z^2$ : $E[e^{tY}] = E[e^{tZ^2}] = \int_{-\infty}^{\infty} e^{tz^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}(1-2t)} dz$$

$$K = (1-2t)^{-\frac{1}{2}}$$

Restrict to $t < \frac{1}{2}$

(so $1 - 2t > 0$)

$$= K \int_{-\infty}^{\infty} \underbrace{\frac{1}{\sqrt{2\pi} K} e^{-\frac{z^2}{2K^2}} dz}_{=1 \quad \text{pdf of } N(0, K^2)}$$

$$= K = (1-2t)^{-\frac{1}{2}}$$
$\uparrow$ mgf of $Z^2$

check book
mgf of Gamma $(\alpha = \frac{1}{2}, \beta = 2)$

$Z^2 \sim$ Gamma $(\frac{1}{2}, 2)$

$$\boxed{Z^2 \sim \chi^2_1}$$

# The $\chi^2$ distribution

- We saw before that the $\chi^2$ distribution with $\nu$ degrees of freedom is simply the $gamma\left(\frac{\nu}{2}, 2\right)$ distribution.
- In Example 3 we saw that if $Z \sim N(0,1)$ then $Z^2 \sim \chi_1^2$.
- Theorem 2: If $Y_1, Y_2, \ldots, Y_n$ are iid $N(\mu, \sigma^2)$ random variables, then

$$\sum_{i=1}^{n} \left(\frac{Y_i - \mu}{\sigma}\right)^2 \sim \chi_n^2$$

$$\frac{Y_i - \mu}{\sigma} \sim N(0,1)$$

$$\text{Let } Z_i = \frac{Y_i - \mu}{\sigma}$$

$$\sum_{i=1}^{n} \left(\frac{Y_i - \mu}{\sigma}\right)^2 = \sum_{i=1}^{n} Z_i^2$$

By ② $Z_i^2 \sim \chi_1^2$     by ① $\sim gamma(\frac{1}{2}, 2)$

we did this for exp but same method show ~~$\beta$~~ (mgf)?

③ $W_1, W_2, \ldots, W_n \sim Gamma(\alpha_i, \beta)$ all indep.

then $\sum_{i=1}^{n} W_i \sim Gamma\left(\sum_{i=1}^{n} \alpha_i, \beta\right)$ (check)

$$\sum_{i=1}^{n}\left(\frac{Y_i - \mu}{\sigma}\right)^2 = \sum_{i=1}^{n} Z_i^2 \overset{③}{\sim} Gamma\left(\sum_{i=1}^{n} \frac{1}{2}, 2\right)$$

$$\sim Gamma\left(\frac{n}{2}, 2\right)$$

$$\overset{①}{\sim} \chi_n^2$$

$\square$

# The distribution of the sample variance

- <u>Definition 4</u>: Given data $Y_1, Y_2, \ldots, Y_n$, the <u>SAMPLE VARIANCE</u> is defined to be

$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i - \bar{Y})^2$$

where $\bar{Y} = \frac{1}{n} \sum_{i=1}^{n} Y_i$ is the sample mean.

- Theorem 3: If $Y_1, Y_2, \ldots, Y_n$ are iid $N(\mu, \sigma^2)$ random variables, then

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$$

Furthermore, $\bar{Y}$ and $s^2$ are independent.

Look at $n=2$ case $\quad Y_1, Y_2$ iid $N(\mu, \sigma^2)$

$$S^2 = \frac{1}{2-1}\left[(Y_1 - \bar{Y})^2 + (Y_2 - \bar{Y})^2\right] \qquad \bar{Y} = \frac{1}{2}(Y_1 + Y_2) = \frac{1}{2}Y_1 + \frac{1}{2}Y_2$$

$$S^2 = \left(Y_1 - \frac{1}{2}Y_1 - \frac{1}{2}Y_2\right)^2 + \left(Y_2 - \frac{1}{2}Y_1 - \frac{1}{2}Y_2\right)^2$$

$$= \left(\frac{1}{2}Y_1 - \frac{1}{2}Y_2\right)^2 + \left(\frac{1}{2}Y_2 - \frac{1}{2}Y_1\right)^2$$

$$= \frac{1}{4}(Y_1 - Y_2)^2 + \frac{1}{4}(Y_1 - Y_2)^2 = \frac{1}{2}(Y_1 - Y_2)^2$$

What is the dist of $Y_1 - Y_2$? $\qquad a_1 Y_1 + a_2 Y_2$

$$Y_1 - Y_2 \sim N(\mu - \mu, \sigma^2 + \sigma^2) \qquad \begin{array}{c} a_1 = 1 \\ a_2 = 1 \end{array}$$

$$\sim N(0, 2\sigma^2)$$

$$\left(\frac{Y_1 - Y_2}{\sqrt{2}\,\sigma}\right)^2 = \frac{(Y_1 - Y_2)^2}{2\sigma^2} \sim \chi_1^2 \qquad \frac{(2-1)S^2}{\sigma^2} \sim \chi_1^2$$

$$2 \rightsquigarrow$$