

## Lecture 12

AMS 131

Hypergeometric: a finite population has  $A$  elements of type 1 and  $B$  of type 2, total size  $(A+B)$

you choose  $n$  elements random without replacement (SRS)

$X = \left( \begin{array}{l} \text{\# of type 1} \\ \text{in sample} \end{array} \right) \sim \text{hypergeometric}(A, B, n)$

$$f_X(X|A, B, n) = \begin{cases} \frac{\binom{A}{x} \binom{B}{n-x}}{\binom{A+B}{n}} & \text{support} \\ I[\max\{0, n-B\} \leq x \leq \min\{n, A\}] \end{cases}$$

PMF

$(A, B, n)$  non negative  $A+B \geq n$

1.  $E(X) = n \cdot \frac{A}{A+B}$  2.  $V(X) = n \cdot \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right) \left( \frac{A+B-n}{A+B-1} \right)$

If taken with replacement (IID),  $X$  would be Binomial with same  $n$  and:

$$p = \frac{A}{A+B} \quad E(X) = np = n \cdot \frac{A}{A+B} \quad V(X) = np(1-p) = n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right)$$

If  $T = A+B$ , total # of elements in population

	mean	variance (of $X$ )
with replacement (IID)	$n \left( \frac{A}{A+B} \right)$	$n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right)$
without replacement (SRS)	$n \left( \frac{A}{A+B} \right)$	$n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right) \left( \frac{T-n}{T-1} \right)$

SRS more information, no repetition, smaller variance

$$0 \leq \alpha = \frac{T-n}{T-1} \leq 1 \quad \text{called the finite population correction.}$$

a.)  $n=1$ ,  $\alpha=1 \iff$  SRS = IID (only 1 element)

b.)  $n=T$ ,  $\alpha=0 \iff$  exhaust the entire population with SRS, no uncertainty left

c.) fix  $n$ , allow  $T \uparrow$ ,  $\alpha \rightarrow 1 \iff$  small sample from a large population, SRS = IID

## Lecture 12 (cont.)

 $\mathcal{X} \sim \text{Poisson}(\lambda) \quad \lambda > 0$ 

$$\text{PMF } f_{\mathcal{X}}(x) = \frac{\lambda^x e^{-\lambda}}{x!} \mathbb{I}_{\{0,1,\dots\}}(x) \quad \leftarrow \text{support of } x$$

$$\left. \begin{array}{l} E(\mathcal{X}) = \lambda \\ V(\mathcal{X}) = \lambda \end{array} \right\} \frac{V(\mathcal{X})}{E(\mathcal{X})} = 1 \quad \begin{array}{l} \text{variance - mean} \\ \text{ratio (VTMR)} \end{array}$$

$\Psi_{\mathcal{X}}(t) = e^{\lambda(e^t - 1)}$  unrealistic, most R.V.s that represent counts of occurrences of events in fixed time interval lengths = VTMR  $\gg 1$ )

The Poisson and Binomial both count the number of "successes" as time goes on, so:

When  $\left( \begin{array}{l} n \text{ is large} \\ p \text{ is close to } 0 \end{array} \right) \rightarrow \text{Binomial}(n, p) \doteq \text{Poisson}(n \cdot p)$

$n$  pos. integer  $\mathcal{X} \sim \text{Binomial}(n, p)$

$0 < p < 1$   $\mathcal{X} \sim \text{Poisson}(\lambda) \quad \lambda > 0$

any sequence of numbers  $\left\{ p_n \right\}_{n=1}^{\infty}$  of values between 0 and 1 with  $\lim_{n \rightarrow \infty} n \cdot p_n = \lambda$

$n = \#$  of successes.

PMF Binomial

PMF Poisson

$$f_{\mathcal{X}}(x | n, p_n) \rightarrow n \rightarrow \infty \rightarrow f_{\mathcal{X}}(y | \lambda)$$

A Poisson Process with rate  $\lambda$  per unit time is a stochastic process with two properties

1. # of arrivals in every interval of time length  $t \sim \text{Poisson}(\lambda t)$
2. # of arrivals in all disjoint (non overlapping) time intervals are independent.

Cryptosporidium case study

people get sick when

1 oocyst

organisms, at one stage in life

Concentration is

5 liters ← 1.3 gal of water

Suppose oocysts occur according to a Poisson process with rate  $\lambda$  oocysts per liter. A filtering system can capture all of the oocysts in a sample, but only has a probability  $p$  of detecting if they are actually there.

$Y =$  # of oocysts in water (actual)  $t$  liters

$X_i = \begin{cases} 1 & \text{if oocyst } i \text{ gets counted} \\ 0 & \text{else} \end{cases}$  (counting events are independent)

$X =$  # of counted oocysts  $(X|Y=y) = \sum_{i=1}^y X_i \sim \text{Binomial}(y, p)$

$X$  only becomes binomial given  $Y=y$

The Law of Total Probability says:  $f_X(x) = P(X=x) = \sum_{y=0}^{\infty} P(Y=y)P(X=x|Y=y)$  for all  $x=0,1,\dots$

$\rightarrow P(Y=y) = \frac{(\lambda t)^y e^{-\lambda t}}{y!}$  for  $y=0,1,\dots$  and  $P(X=x|Y=y) = \binom{y}{x} p^x (1-p)^{y-x}$

if  $X=x$ , then  $Y \geq x$  - at least as large as the number of oocysts detected

$$f_X(x) = \sum_{y=x}^{\infty} \binom{y}{x} p^x (1-p)^{y-x} \frac{(\lambda t)^y e^{-\lambda t}}{y!} = \frac{e^{-p\lambda t} (p\lambda t)^x}{x!}$$

$X \sim \text{Poisson}(p\lambda t)$ : losing  $(1-p)$  proportion to faulty counting lowers the rate of the Poisson process from  $\lambda$ /liter to  $\lambda \cdot p$ /liter

Poisson mixture of Binomials (= another Poisson)

In practice,  $p$  is small (close to zero). How much water ( $t$  liters) do you need to filter to achieve:

$P(\text{at least 1 oocyst detected}) \geq 1 - \alpha$  for small  $\alpha$ ?

## Lecture 12 (cont.)

$P(\text{at least 1 detected}) = 1 - P(\text{none detected})$

$$= 1 - P(X=0) = 1 - e^{-p\lambda t} \geq 1 - \alpha$$

$$\Leftrightarrow \alpha \geq e^{-p\lambda t} \Leftrightarrow \ln(\alpha) \geq -p\lambda t \quad \rightarrow \quad t \geq \frac{-\ln(\alpha)}{p\lambda}$$

$\alpha = 0.01$  (failure rate)       $p = 0.1$

$\lambda = 0.2$  / liter = 1 / 5 liters

to achieve 99%,  $t$  has to be at least 230.3 liters

**Negative Binomial Distribution:** you're watching a potentially endless sequence of Bernoulli trials with constant success probability  $p$ .

$X = \#$  of failures before  $r^{\text{th}}$  success       $r \geq 1$  integer

$X \sim$  Negative Binomial dist. PMF:  $f_X(x|r, p) = \binom{r+x-1}{x} p^r (1-p)^x$

$$\cdot \mathbb{I}_{\{0, 1, 2, \dots\}}(x)$$

now we're counting failures instead of successes, there are two ways to estimate (unknown)  $p$ .

1. decide ahead of time to sample  $n$  success/failures and record the (random)  $\#$   $S$  of success you see, a good estimate would be  $\hat{p}_B = \frac{S}{n}$  (binomial)
2. Sample until you see  $s$  successes, then record the  $\#$  of trials  $N$  needed to accumulate that many successes, estimate  $\hat{p}_{NB} = \frac{s}{N}$  (neg. binomial)

Binomial: numerator random, denominator fixed

N-Binomial: numerator fixed, denominator random.

Set  $r=1$  and record the number  $X$  of failures until the first success,  $X$  follows the Geometric ( $p$ ) dist.

$$\text{PMF } f_X(x|p) = p(1-p)^x I_{\{0,1,\dots\}}(x) \quad \leftarrow \begin{array}{l} \text{parameter } p \\ \text{support of } X \end{array}$$

$X_1, \dots, X_n$  IID Geometric ( $p$ )

$$\rightarrow \sum_{i=1}^n X_i \sim \text{Negative Binomial}(n, p)$$

$X_1, \dots, X_n$  IID Bernoulli ( $p$ )

$$\rightarrow \sum_{i=1}^n X_i \sim \text{Binomial}(n, p)$$

$X \sim \text{Negative Binomial}(r, p)$

$$Y_X(t) = \left[ \frac{p}{1 - (1-p)e^t} \right]^r \quad \text{for } t < \log\left(\frac{1}{1-p}\right)$$

$$E(X) = \frac{r(1-p)}{p} \quad V(X) = \frac{r(1-p)}{p^2}$$

$X \sim \text{Geometric}(p)$

$$P(X = k+t | X \geq k) = P(X = t) \quad \left. \begin{array}{l} k \\ t \end{array} \right\} \begin{array}{l} \text{non-negative} \\ \text{integers} \end{array}$$

Wait  $k$  trials, what's the chance you succeed on the next  $k+t$  trial, just  $P(X=t)$  waiting  $t$  trials - memoryless property of the Geometric distribution  
forget the condition of waiting  $k$  trials  
ONLY discrete distribution with this property.

$\underbrace{F \ F \ F \ F \ F \ S}_{X}$   $X = \#$  of failures until first success  
 $\underbrace{\hspace{10em}}_{Y}$   $Y = \#$  of failures, starting at trial  $k+1$  until next success

$Y$  has the same dist. as  $X$  and is independent of what happened in the first  $k$  trials  
 $\hookrightarrow$  "the process has no memory"

Continuous distributions

Normal (Gaussian) dist  $X \sim \text{Normal}(\mu, \sigma^2)$  $\mu = \text{mean}$ ,  $\sigma^2 = \text{variance}$   $0 < \sigma^2 < \infty$ 

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \quad f_X(x|\mu, \sigma^2)$$

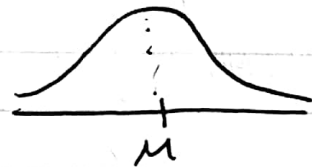
$\sigma \geq 0$   $\sigma = 0$   $\frac{1}{\mu}$  , use  $\sigma > 0$  converting to standard units

$-\infty < \mu < \infty$

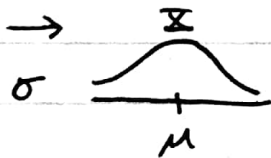
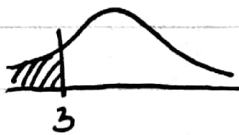
1. many random processes have this dist. shape
2. The Central Limit Theorem

 $X \sim \text{Normal}(\mu, \sigma^2) = N(\mu, \sigma^2)$  $E(X) = \mu$ ,  $V(X) = \sigma^2$ ,  $SD(X) = \sigma$ 

$$\psi_X(t) = \exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right)$$

Center of symmetry = median = mean = mode =  $\mu$  $(X|\mu, \sigma^2) \sim N(\mu, \sigma^2)$  $Y = aX + b$ ,  $a \neq 0$ ,  $b$  constants

Normality is preserved under

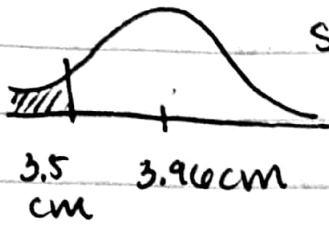
|a| $\sigma$  linear transformation $X \sim N(\mu, \sigma^2)$  $P(X \leq 3)$  CDF of Normal has no closed form.

SD=1

Standard Normal Dist.  $\mu=0$ ,  $SD=1$ 

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{t^2}{2}\right] dt$$

$n=103$  monarch butterflies, wing length.



SD = 0.29 cm

$$\begin{bmatrix} y_1 = 4.1 \\ y_2 = 3.3 \\ \dots \\ y_n = 4.7 \end{bmatrix} \quad n=103$$

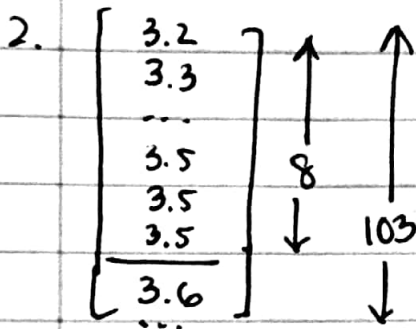
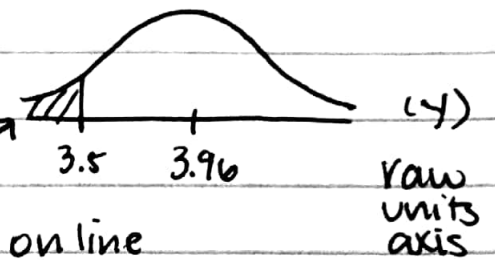
$$\bar{y} = 3.96 \text{ cm}$$

$$\sigma_y = 0.29 \text{ cm}$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} \quad P(X \leq 3.5)?$$

1.  $\begin{bmatrix} y_1 = 4.1 \\ y_2 = 3.3 \\ \dots \\ y_n = 4.7 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ \dots \\ 0 \end{bmatrix}$   $X = \begin{cases} 1 & \text{if } y \leq 3.5 \\ 0 & \text{else} \end{cases}$

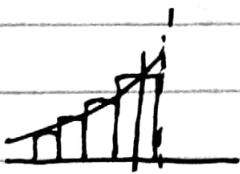
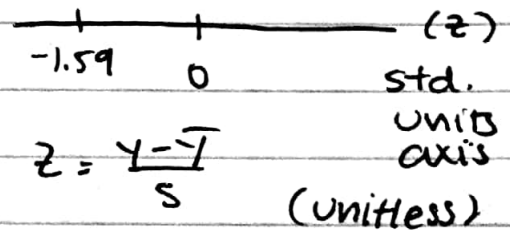
$$\mu = \frac{\text{sum}}{103}$$



(exact) =  $\frac{8}{103} \approx 7.8\%$

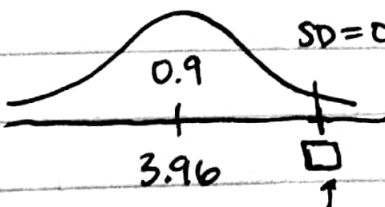
online  
calc = 5.63%

Sort, then count



continuity correction

3.5  
3.55 ← 0.0787 = 7.87%



What is the 90th percentile of the wing length dist?  
- Inverse CDF problem

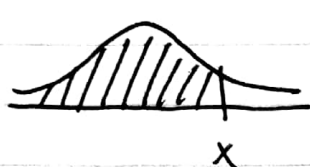
online calc = 4.332

PDF of  $Z \sim \text{Normal}(0,1)$

CDF

lower case  $\phi_Z(x) \triangleq \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \rightarrow \Phi_Z(x) \triangleq \int_{-\infty}^x \phi_Z(t) dt$

1.  $e^{-cx^2}$   $c > 0$  has no anti-derivative in closed form, we have approximations in a table using numerical integration
2. The Normal PDF is symmetric (for all  $x \in \mathbb{R}$ )  
 $\Phi(-x) = 1 - \Phi(x)$        $\phi_X(x) = \phi_X(-x)$  ← same height  
 $\Phi^{-1}(p) = -\Phi^{-1}(1-p)$  for all  $0 < p < 1$



x

 $\Phi_X(x)$ 

-x

 $\Phi_X(-x)$ 

Symmetric on both sides

$$\Phi_X(x) + \Phi_X(-x) = 1$$

$X \sim \text{Normal}(\mu, \sigma^2) \rightarrow Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$  standard

$$F_X(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

$$F_X^{-1}(p) = \mu + \sigma \Phi^{-1}(p)$$

Empirical Rule  $(\mu \pm \sigma) = 68\%$  of data. } exact for all normal distributions  
 $(\mu \pm 2\sigma) = 95\%$   
 $(\mu \pm 3\sigma) = 99.7\%$

$X_1, \dots, X_k$  independent  $X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$

$$\rightarrow \sum_{i=1}^k X_i \sim \text{Normal}\left(\sum_{i=1}^k \mu_i, \sum_{i=1}^k \sigma_i^2\right)$$

additive property, why Normals are indexed by  $V(X)$

women height

$$N(65 \text{ in}, (3.2)^2 \text{ in}^2)$$

$$\mu \quad \sigma = 3.2 \text{ in}$$

men height

$$N(69.5 \text{ in}, (3.3)^2 \text{ in}^2)$$

$$\mu \quad \sigma = 3.3 \text{ in}$$



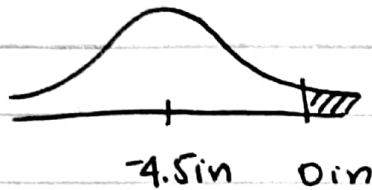
1 woman with height  $\underline{W}$  } chosen randomly,  
 1 man with height  $\underline{M}$  } independent

$P(\text{woman taller than man}) = P(\underline{W} \geq \underline{M})$  due to D:

$D = \underline{W} - \underline{M}$  looking for positive D ← linear

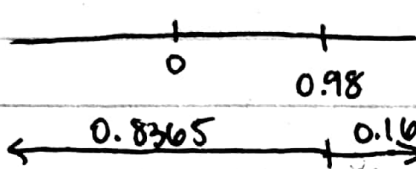
$D \sim N(65 - 69.5 = -4.5 \text{ in}, (3.2)^2 + (3.3)^2 = 21.1 \text{ in}^2)$  operation

$SD = \sqrt{21.1} = 4.6 \text{ in}$



convert to standard units:

$\frac{0 - (-4.5)}{4.6} = +0.98$

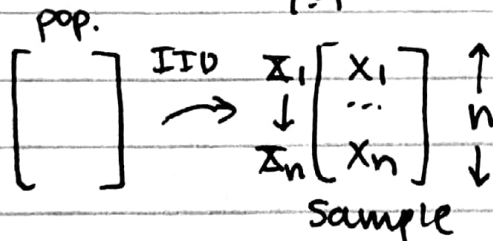


using online calc.

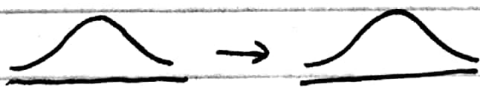
→ so  $P(\underline{W} \geq \underline{M}) \approx 16\%$   
 about 1 in 6

R.V.  $\underline{X}_1, \dots, \underline{X}_n \rightarrow$  sample mean of  $(\underline{X}_1, \dots, \underline{X}_n)$

is  $\bar{\underline{X}}_n = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$   $\left\{ \begin{array}{l} \underline{X}_i \stackrel{\text{IID}}{\rightarrow} N(\mu, \sigma^2) \\ (i=1, \dots, n) \end{array} \right.$



$\bar{\underline{X}}_n = \frac{1}{n} \sum_{i=1}^n \underline{X}_i = (\bar{X})$  linear transformation



$E(\bar{\underline{X}}_n) = E\left(\frac{1}{n} \sum_{i=1}^n \underline{X}_i\right)$

$E(\bar{\underline{X}}_n) = \mu \leftrightarrow \bar{\underline{X}}_n$  is unbiased

for  $\mu$ .

$\bar{\underline{X}}_n$  is an unbiased estimator

$= \frac{1}{n} E\left(\sum_{i=1}^n \underline{X}_i\right) = \frac{1}{n} \sum_{i=1}^n E(\underline{X}_i)$

$= \frac{1}{n} \sum_{i=1}^n \mu = \mu \checkmark$

In frequentist statistics, the standard deviation on an estimator  $\hat{\theta}$  (R.V.) of a parameter  $\theta$  is called the standard error  $SE(\hat{\theta})$  of  $\hat{\theta}_k$

AMS 131

Lecture 12 (cont.)

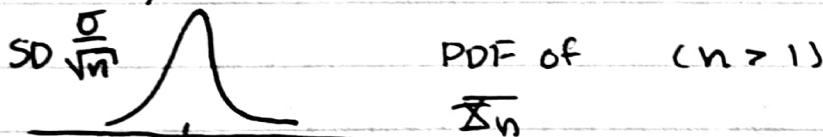
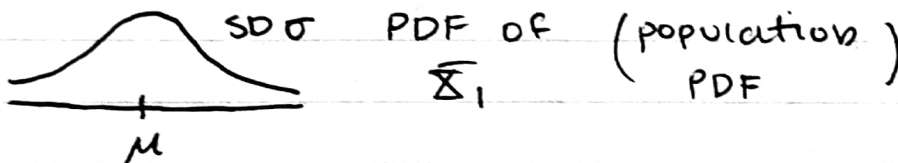
$$SD(\bar{X}_n) \triangleq SE(\bar{X}_n) = \frac{c}{\sqrt{n}}$$

$$SD\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} SD\left(\sum_{i=1}^n X_i\right) = \text{hard?}$$

$$V(\bar{X}_n) = V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} V\left(\sum_{i=1}^n X_i\right) \stackrel{\text{⊕}}{=} \frac{1}{n^2} \sum_{i=1}^n V(X_i)$$

$$\text{⊕} = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$SD(\bar{X}_n) = SE(\bar{X}_n) = \sqrt{V(\bar{X}_n)} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}} \quad \star \text{ square root law}$$



get closer with  $n \rightarrow \infty$

