Statistics

Probability

Set theory

random variable — discrete — binomial, hyper, poisson — continuous — Normal, Exponential
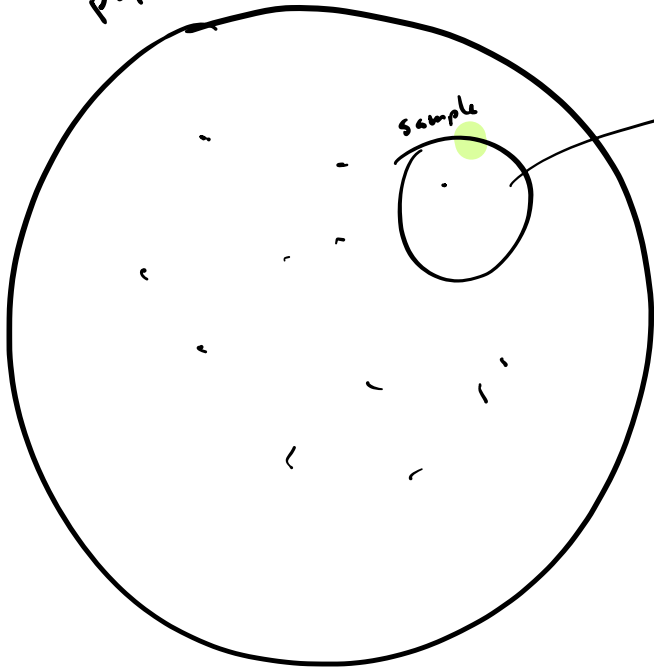
# STAT 515 fa 2023 Lec 09 slides

## Sampling distributions and the Central Limit Theorem

Karl B. Gregory

University of South Carolina

These slides are an instructional aid; their sole purpose is to display, during the lecture, definitions, plots, results, etc. which take too much time to write by hand on the blackboard. They are not intended to explain or expound on any material.

Population

$\mu$ = population mean

$\sigma^2$ = population variance

random sample

$X_1, X_2, ..., X_n$

$\underline{n}$ = sample size

sample

$$\bar{X}_n = \frac{1}{n}(X_1 + \cdots + X_n) = \frac{1}{n}\sum_{i=1}^{n} X_i$$

"X bar"

Sample mean

$$S_n^2 = \frac{1}{n-1}\sum_{i=1}^{n}(X_i - \bar{X}_n)^2$$

= avg. squared deviation from the mean in our sample

Sample variance

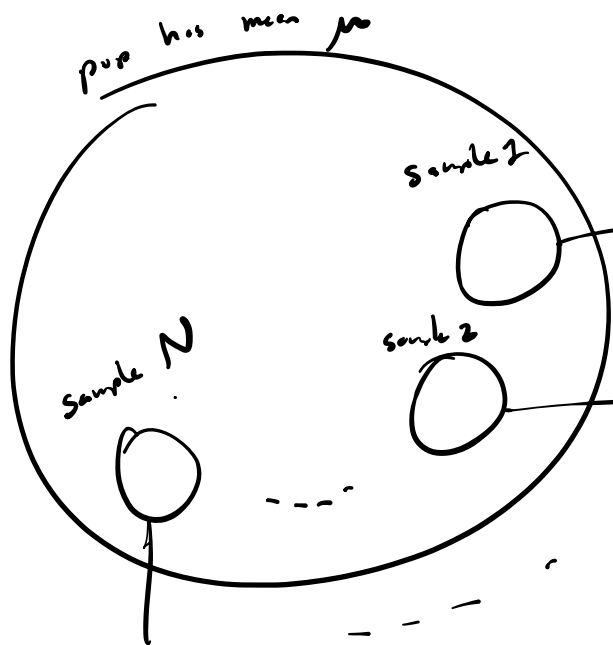$X \sim \text{Binom}(n, p)$

## Random sample

*random variable*

A collection of independent rvs with the same distribution is a *random sample*.

- Often denote by $X_1, \ldots, X_n$, where $n$ is the *sample size*.
- In random sample, $X_1, \ldots, X_n$ are *iid*: independent and identically distributed.
- Common distribution of $X_1, \ldots, X_n$ called the *population distribution*.
- Can write $X_1, \ldots, X_n \overset{\text{ind}}{\sim} F$ if a rs from a distribution $F$.

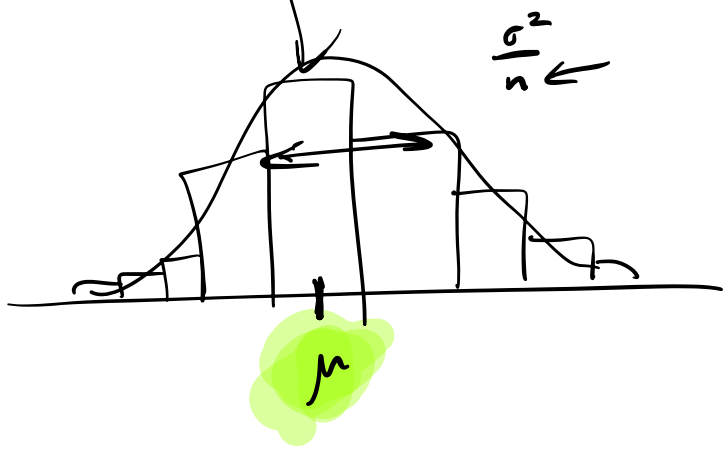$$X_1, \ldots, X_n \overset{\text{ind}}{\sim} \text{Binom}(n, p)$$

Goal is to learn from $X_1, \ldots, X_n$ about the population distribution.

pop has mean $\mu$

Sample 1

Sample N

Sample 2

$\overline{X}_1$

$\overline{X}_2$

$\overline{X}_N$

$\frac{\sigma^2}{n}$

$\mu$

$\overline{x}$

*random sample*

$$\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n}X_i$$

## Expected value and variance of the sample mean

Let $X_1, \ldots, X_n$ be a rs from a population with mean $\mu$ and $\sigma^2$. Then

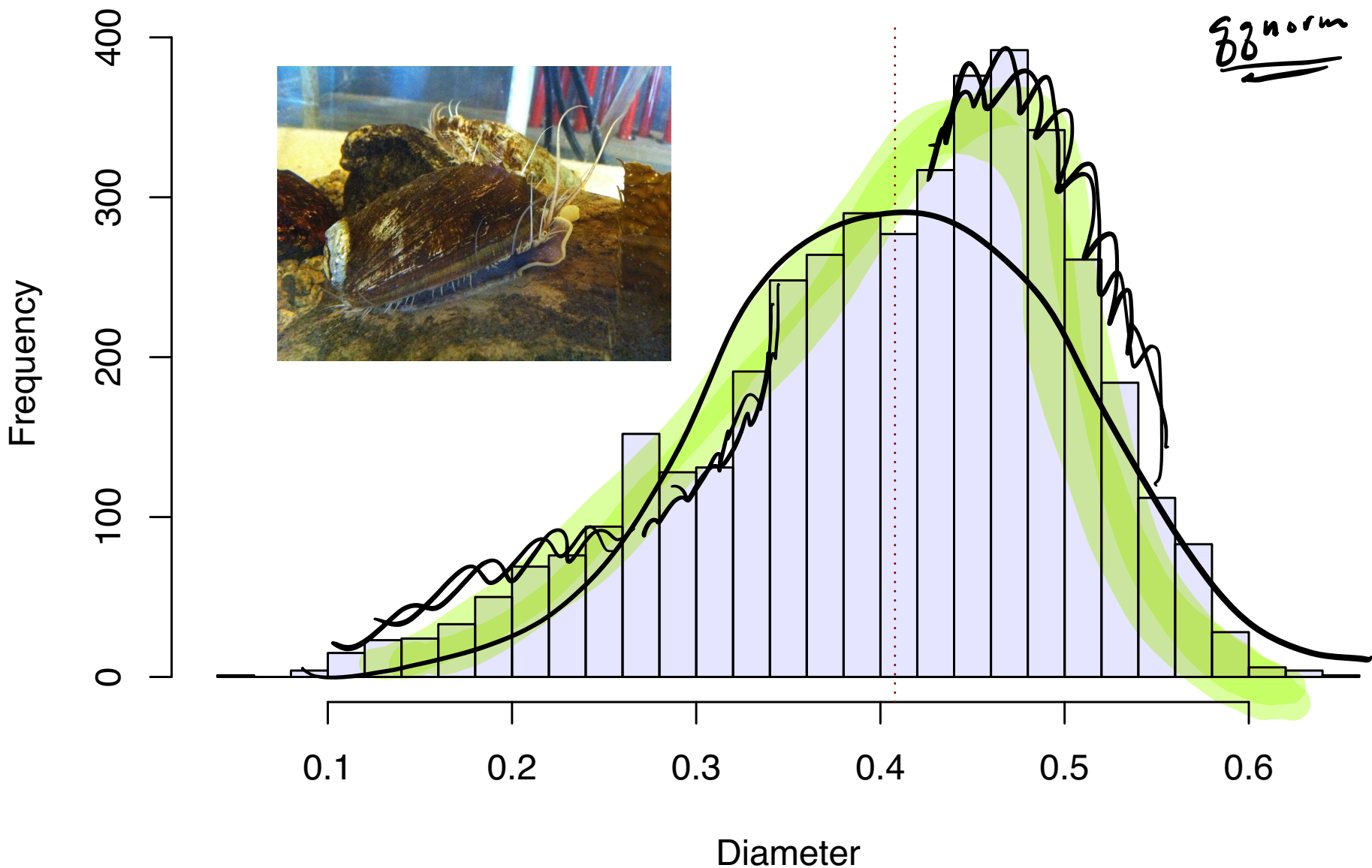$$\mathbb{E}\bar{X}_n = \mu \quad \text{and} \quad \text{Var}\,\bar{X}_n = \frac{\sigma^2}{n}.$$

## Examples:

1. If $X_1, \ldots, X_n \overset{\text{ind}}{\sim} \text{Normal}(\mu, \sigma^2)$, then $\mathbb{E}\bar{X}_n = \mu$ and $\text{Var}\,\bar{X}_n = \sigma^2/n$.

2. If $X_1, \ldots, X_n \overset{\text{ind}}{\sim} \text{Bernoulli}(p)$, then $\mathbb{E}\bar{X}_n = p$ and $\text{Var}\,\bar{X}_n = p(1-p)/n$.
$\mu = p$
$\sigma^2 = p(1-p)$

3. If $X_1, \ldots, X_n \overset{\text{ind}}{\sim} \text{Poisson}(\lambda)$, then $\mathbb{E}\bar{X}_n = \lambda$ and $\text{Var}\,\bar{X}_n = \lambda/n$.
$\mu = \lambda$
$\sigma^2 = \lambda$

4. If $X_1, \ldots, X_n \overset{\text{ind}}{\sim} \text{Exponential}(\lambda)$, then $\mathbb{E}\bar{X}_n = 1/\lambda$ and $\text{Var}\,\bar{X}_n = 1/(n\lambda^2)$.
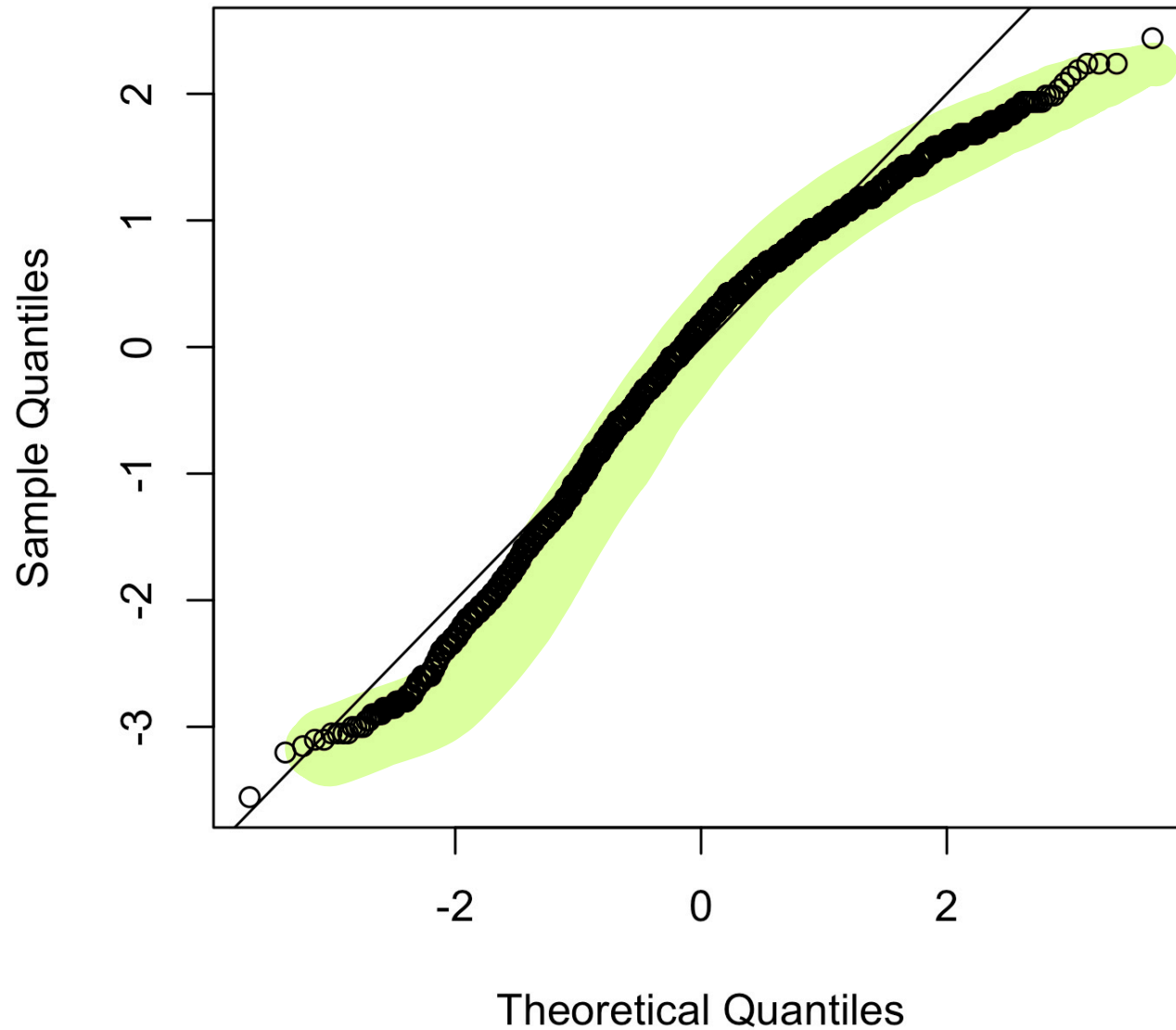$\mu = \frac{1}{\lambda}$
$\sigma^2 = \frac{1}{\lambda^2}$

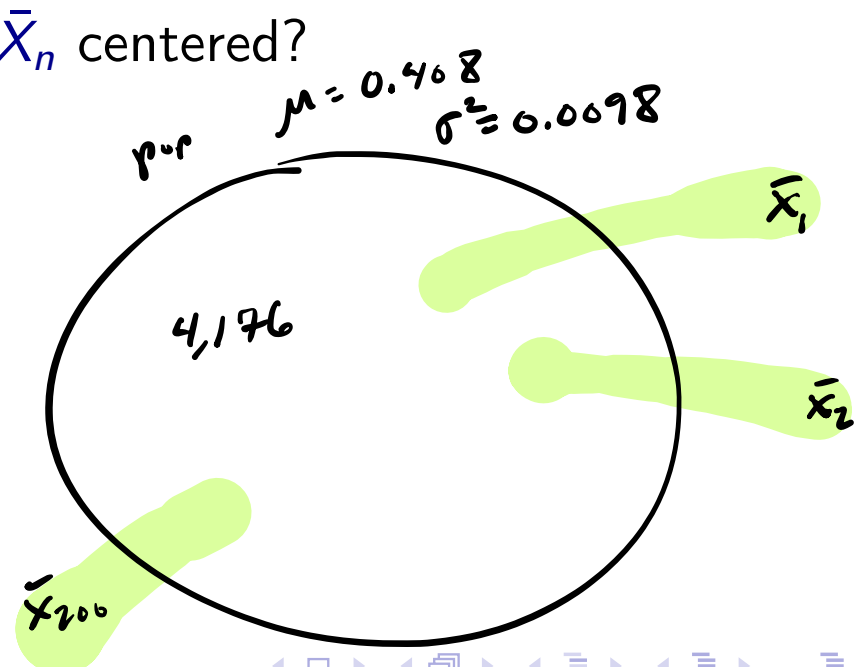Consider the diameters of 4,176 abalones with mean 0.4078915.   link to data

← population

← μ

ggnorm

Normal Q-Q plot of abalone diameters

**Exercise:** Treat the 4,176 abalone as a population. The mean diameter is $\mu = 0.408$. Let $\bar{X}_n$ be the mean diameter from a sample of abalone.
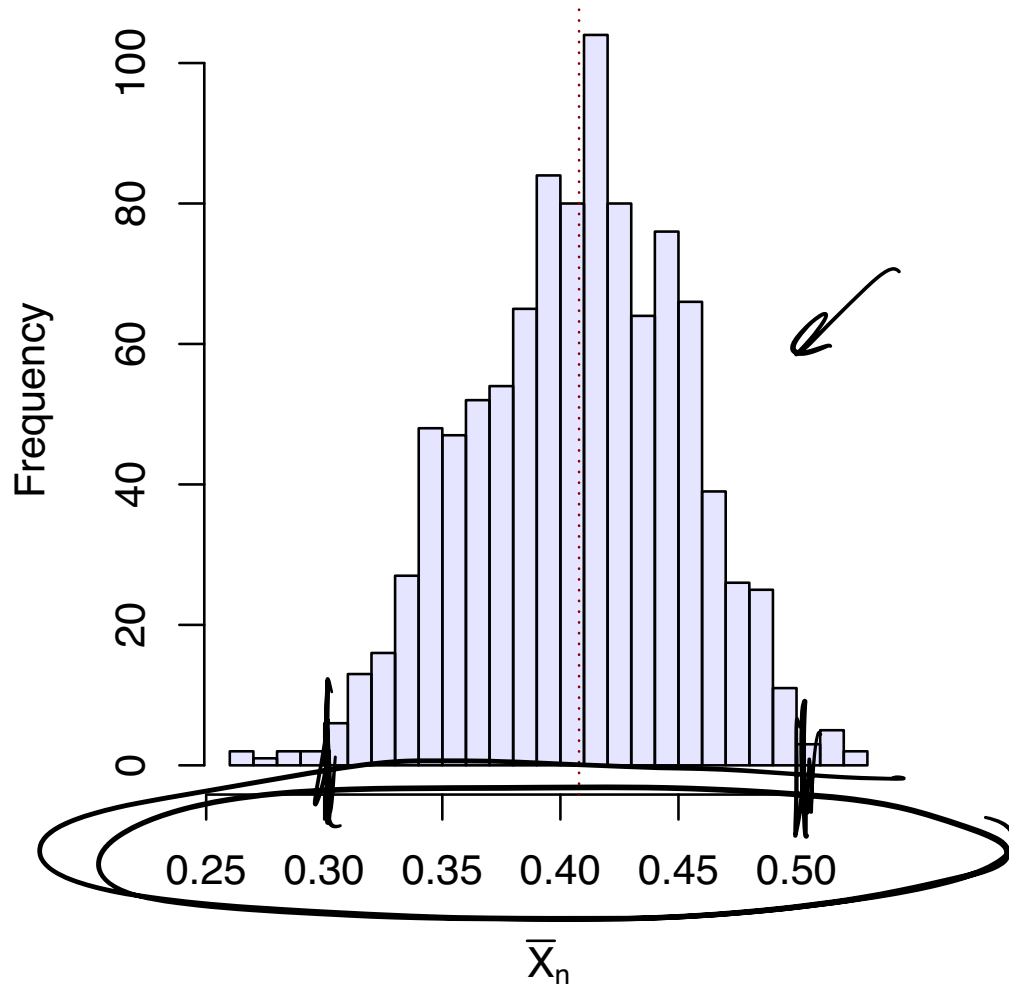
1. For the sample sizes $n = 5, 25, 100$, draw 1,000 samples and
   1. Make a histogram of the $\bar{X}_n$ values.
   2. Make a Normal Q-Q plot of the $\bar{X}_n$ values.
2. Around what value are the values of $\bar{X}_n$ centered?
3. What changes as $n$ changes?

$\mu = 0.408$

$\sigma^2 \approx 0.0098$

pop

4,176

$\bar{X}_1$

$\bar{X}_2$

$\bar{X}_{200}$

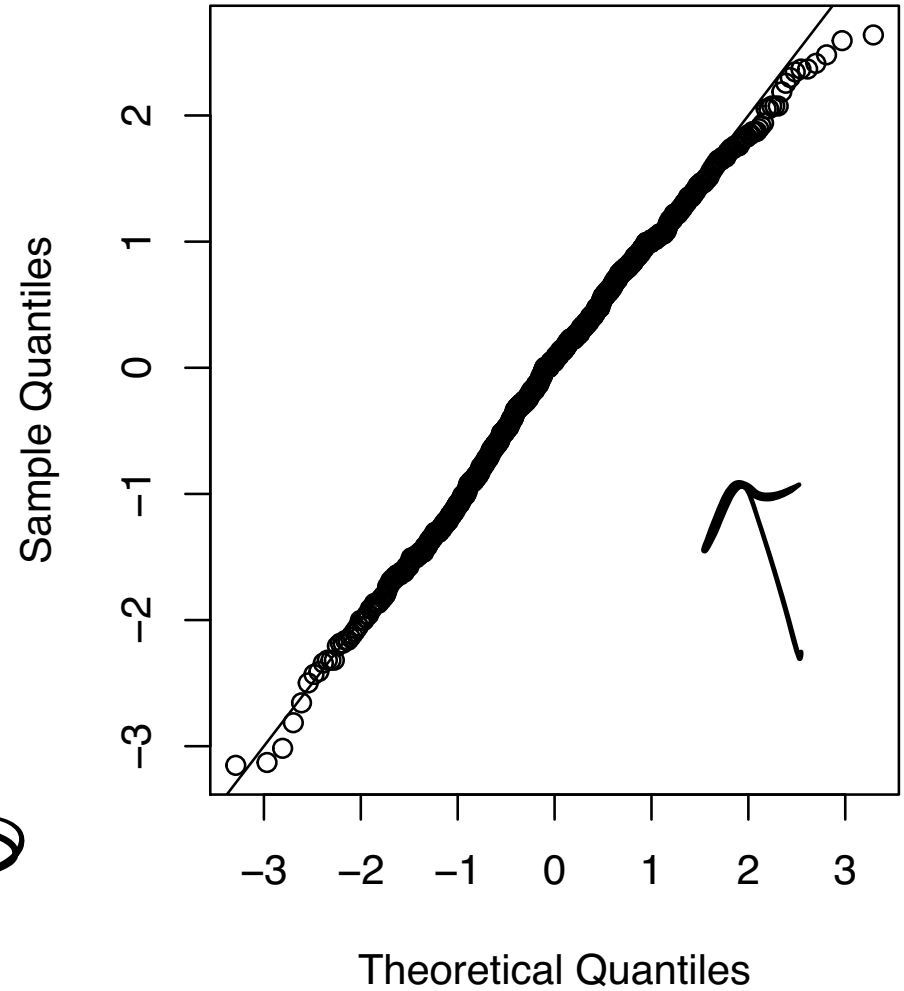| $n$ | $\mathrm{Var}\,\bar{X}_n$ | $\dfrac{\sigma^2}{n}$ |
|-----|------|------|
| 5 | 0.002 | 0.002 |
| 25 | 0.00036 | 0.00039 |
| ↓ | | |

```r
abalones <- read.csv("/Users/karlgregory/Desktop/abalone/abalone.data")

diam <- abalones$V3

mean(diam)
var(diam)

hist(diam)

n <- 5
xbar <- numeric(200)
for(i in 1:200){

    # draw a random sample of size n from the population of diameters:
    X <- sample(diam,n,replace = FALSE)
    xbar[i] <- mean(X)

}

mean(xbar)
var(xbar)
var(diam)/n
```
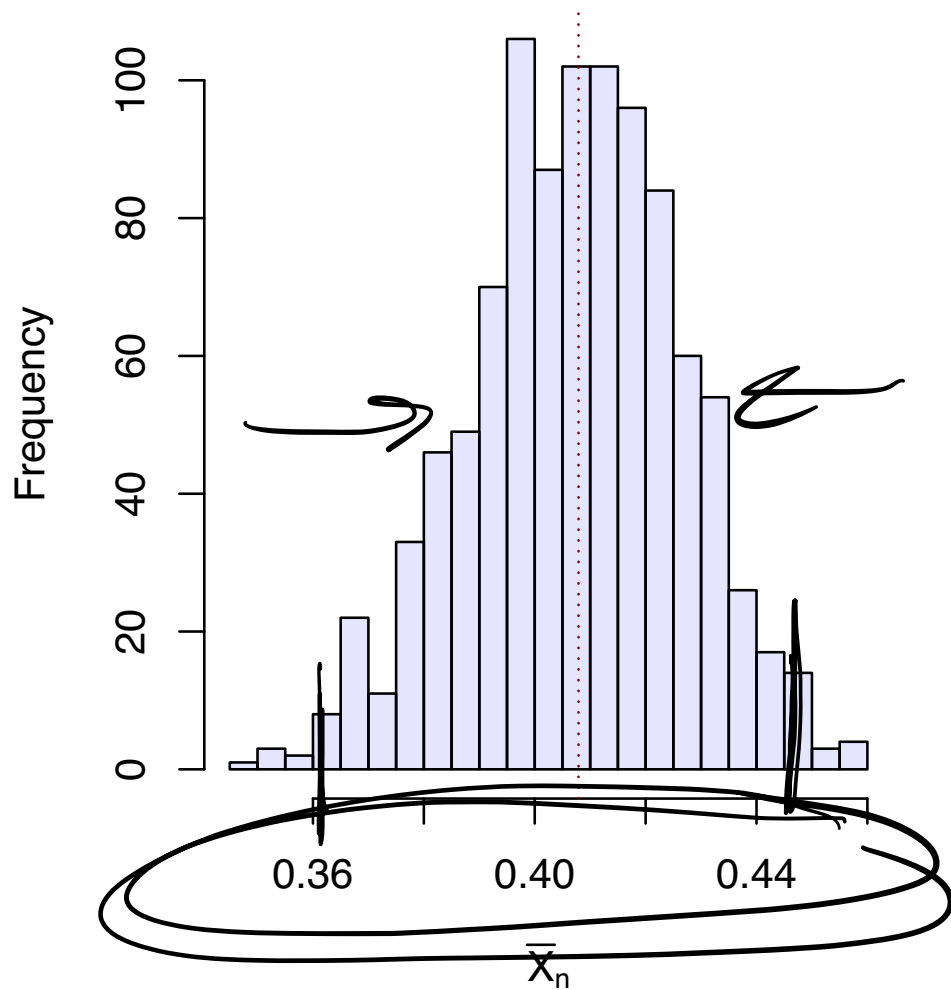
Histogram of $\overline{X}_n$ with n = 5

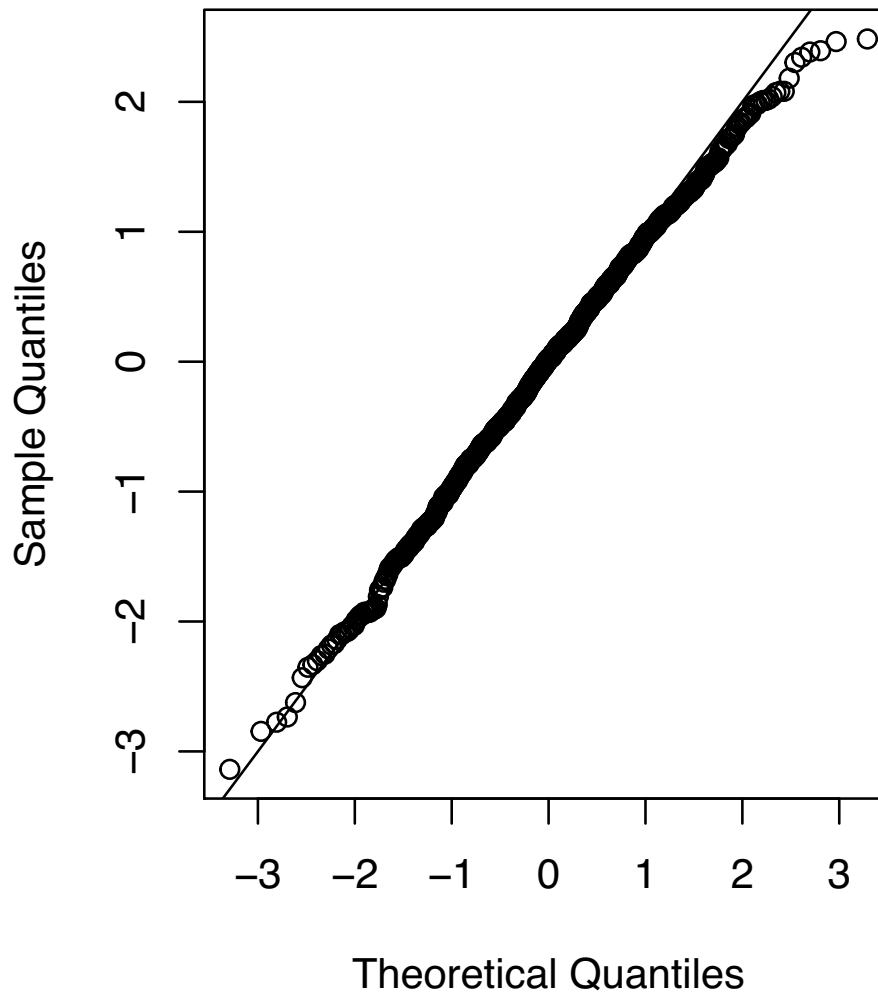Normal Q–Q plot of $\sqrt{n}(\overline{X}_n - \mu)/\sigma$

## Histogram of $\overline{X}_n$ with n = 25



## Normal Q–Q plot of $\sqrt{n}(\overline{X}_n - \mu)/\sigma$

Histogram of $\overline{X}_n$ with n = 100

Normal Q–Q plot of $\sqrt{n}(\overline{X}_n - \mu)/\sigma$



| x | | 0 | 1 | $\overline{X}_n$ |
|---|---|---|---|---|
| $P(X=x)$ | | 1-p | p | |

$EX = 1 \cdot p + 0(1-p) = p$

$$X_1, \ldots, X_n \overset{ind}{\sim} Bernoulli(p)$$
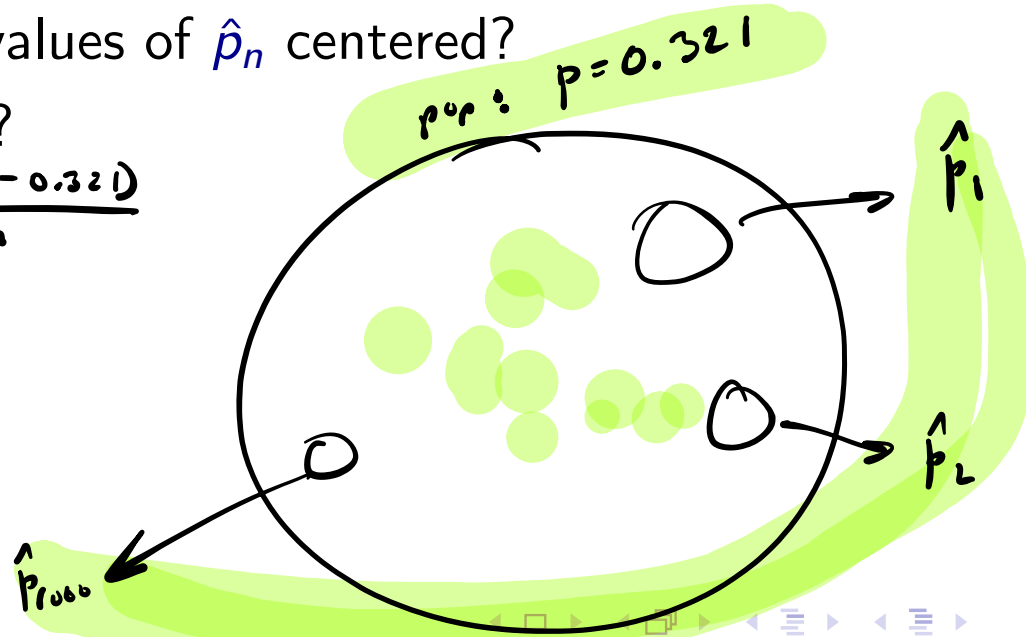
$\mu = p$

$\sigma^2 = p(1-p)$

$$\hat{p}_n = \bar{X}_n = \frac{1}{n}(X_1 + \cdots + X_n)$$

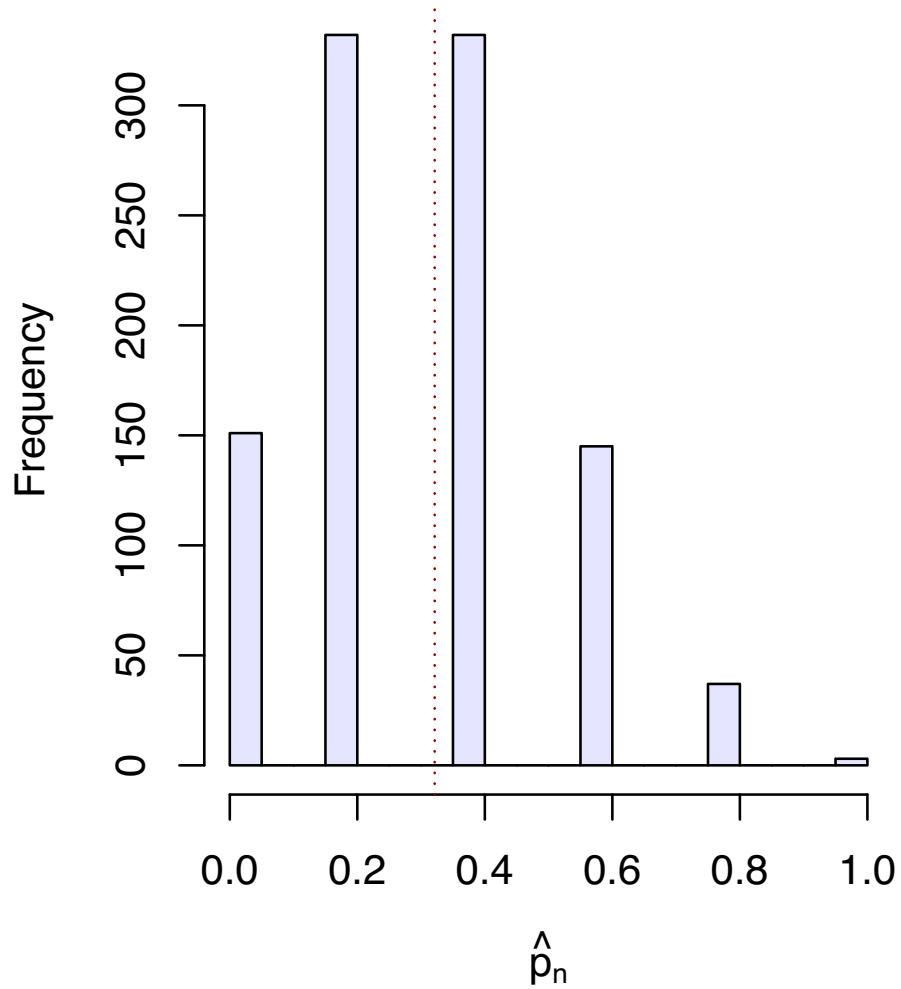$$= \frac{\#\{1s\}}{n} = \text{proportion of } 1s.$$

**Exercise:** Treat the 4,176 abalone as a population. The proportion classified as infants among the abalone is $p = 0.321$; let $\hat{p}_n$ represent the proportion of infants in a random sample of abalone.

1. For the sample sizes $n = 5, 25, 100$, draw 1,000 samples and
   1. Make a histogram of the $\hat{p}_n$ values.
   2. Make a Normal Q-Q plot of the $\hat{p}_n$.

2. Around what value are the values of $\hat{p}_n$ centered?
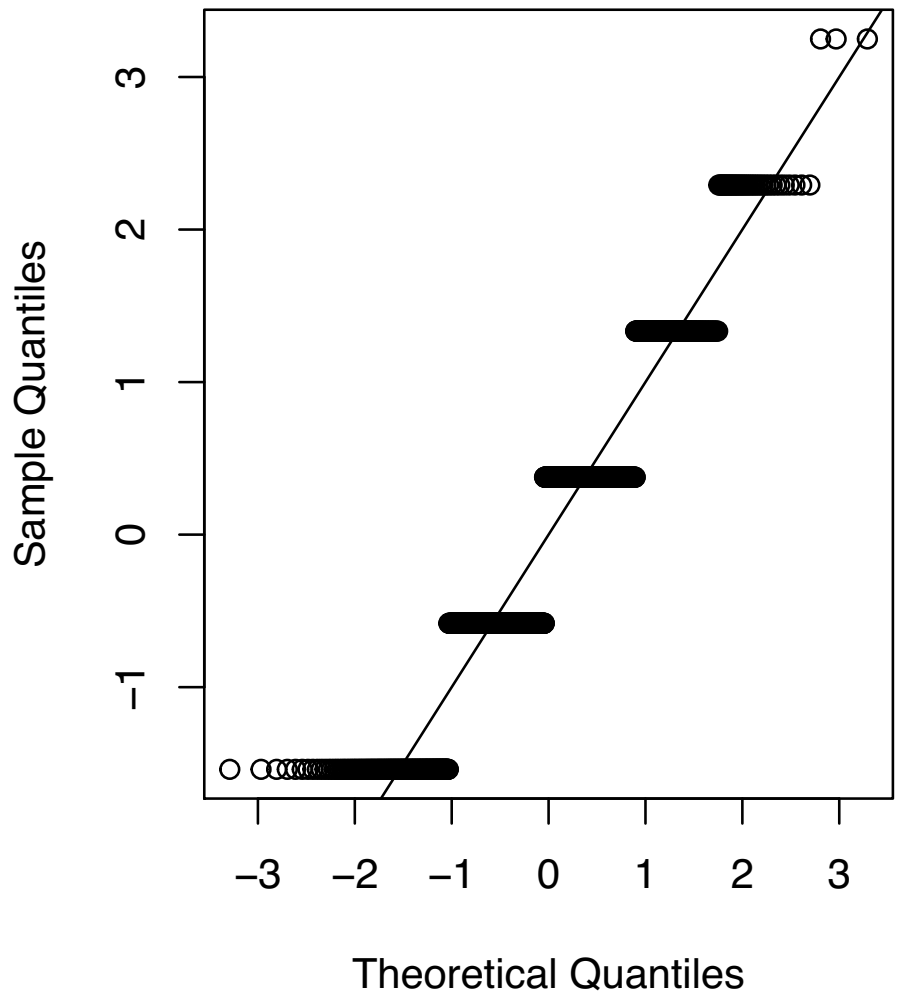
3. What changes as $n$ changes?

pop: $p = 0.321$

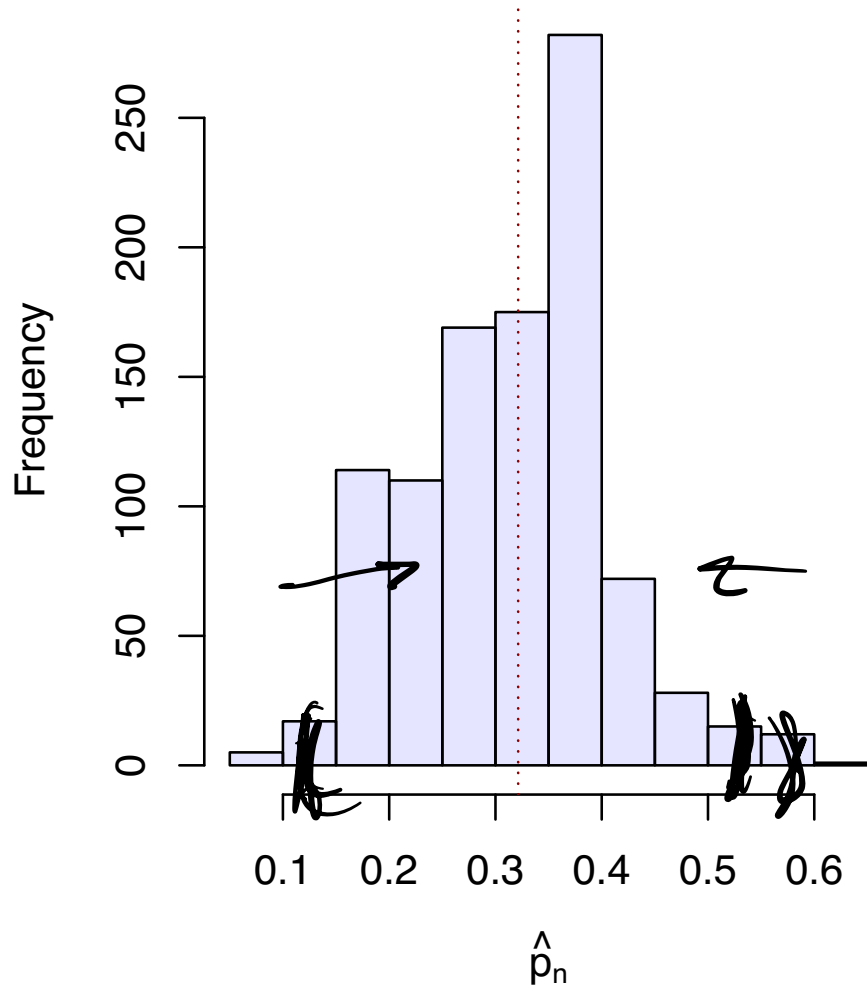| $n$ | Var $\hat{p}_n$ | $\frac{p(1-p)}{n} = \frac{0.321(1-0.321)}{n}$ |
|-----|-----------------|-----------------------------------------------|
| 5   | 0.0416          | 0.0435                                          |
| 25  | 0.0083          | 0.0087                                          |

$\hat{p}_1$

$\hat{p}_2$

$\hat{p}_{1000}$
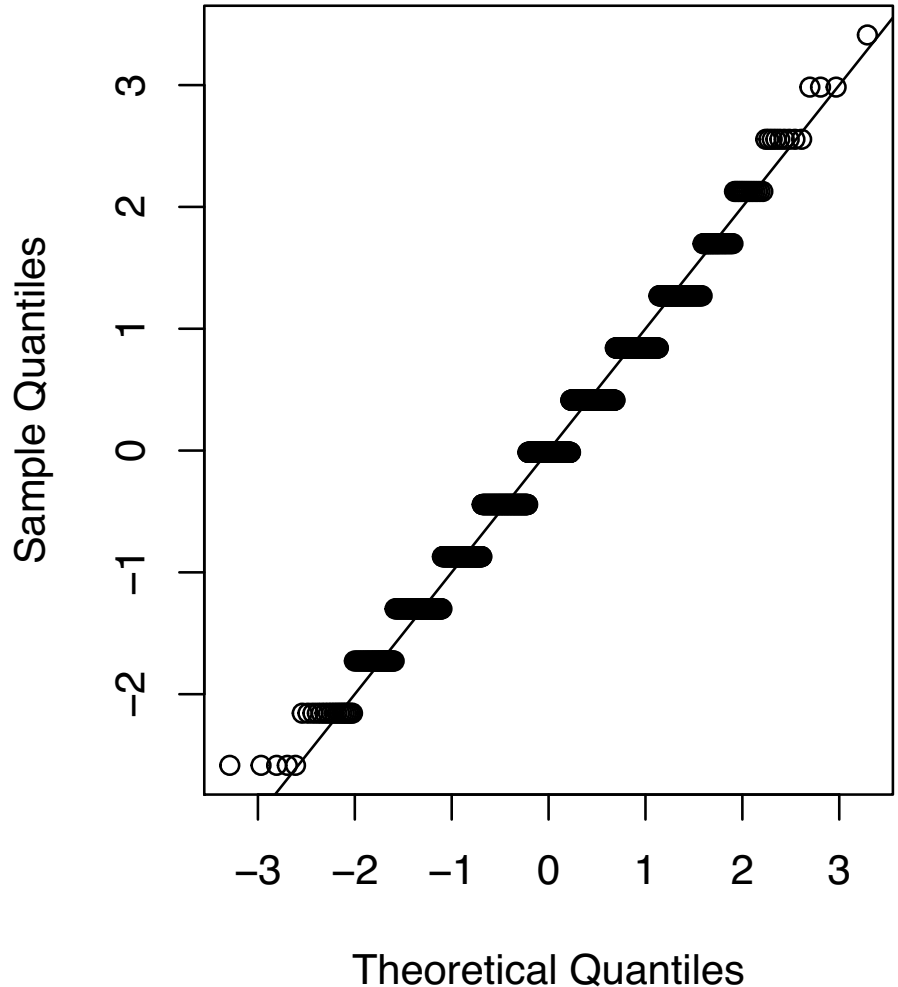
$n = 5$

## Histogram of $\hat{p}_n$ with $n = 5$

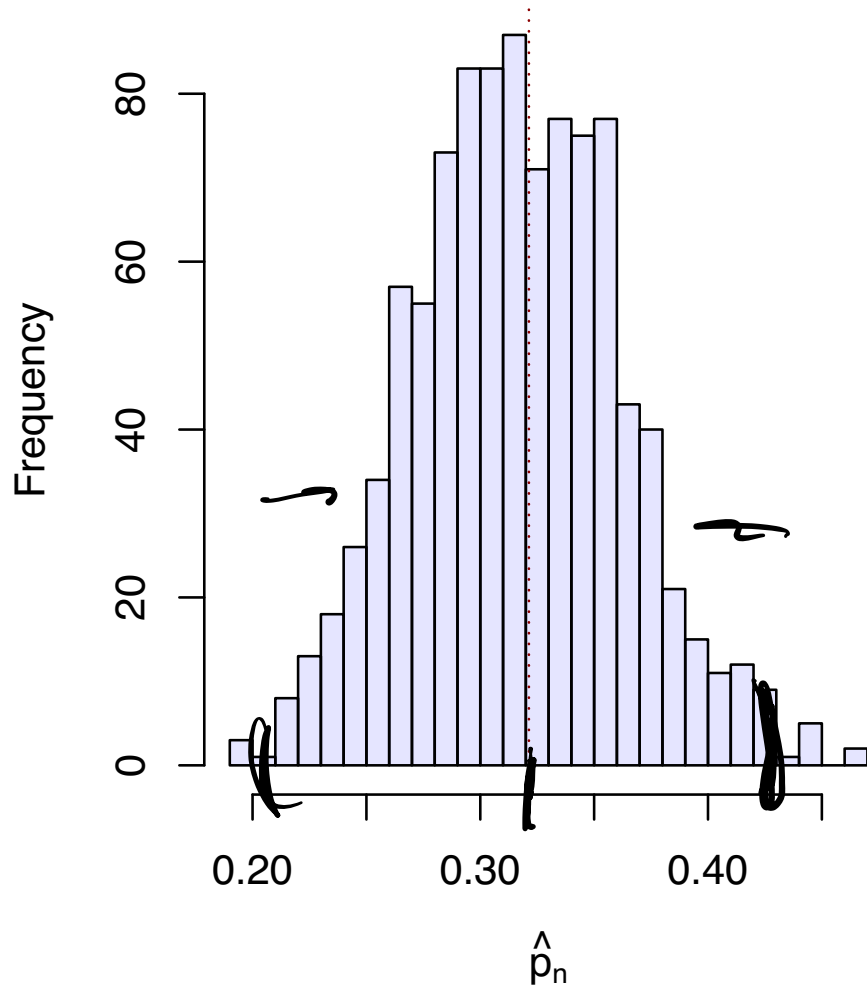## Normal Q–Q plot of $\sqrt{n}(\hat{p}_n - p)/\sqrt{p(1-p)}$
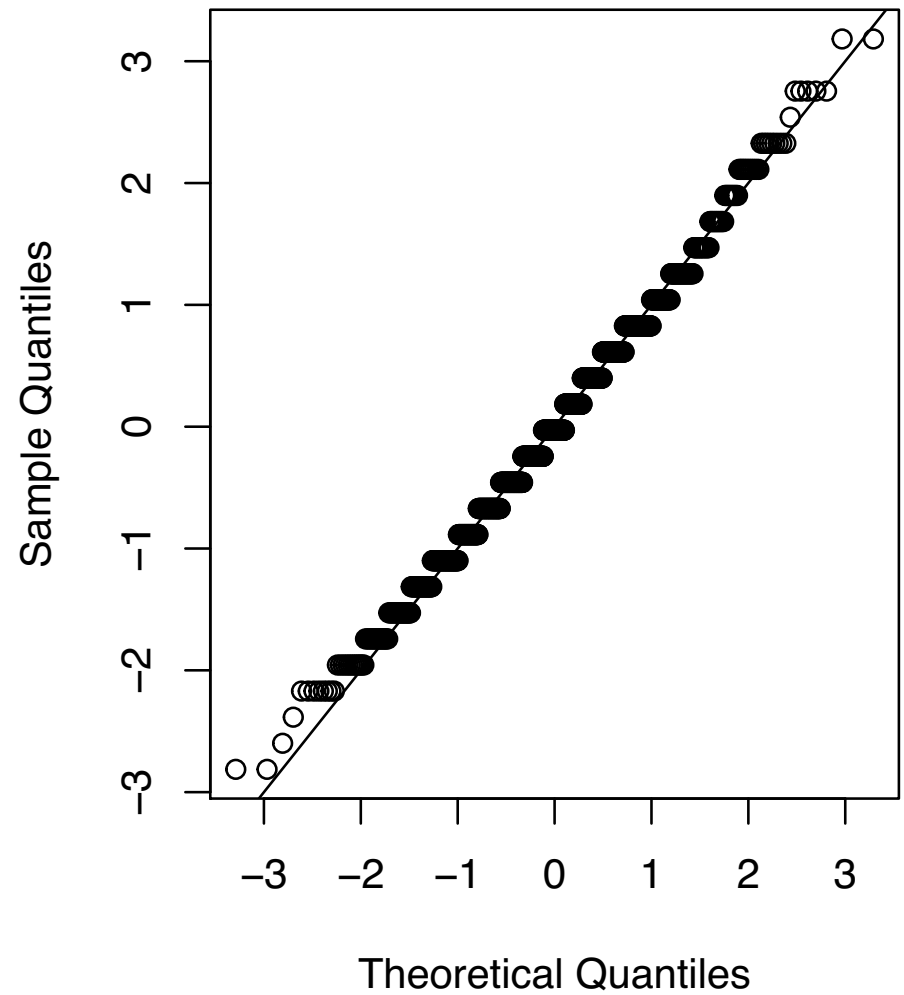
Histogram of $\hat{p}_n$ with n = 25

Normal Q–Q plot of $\sqrt{n}(\hat{p}_n - p)/\sqrt{p(1-p)}$

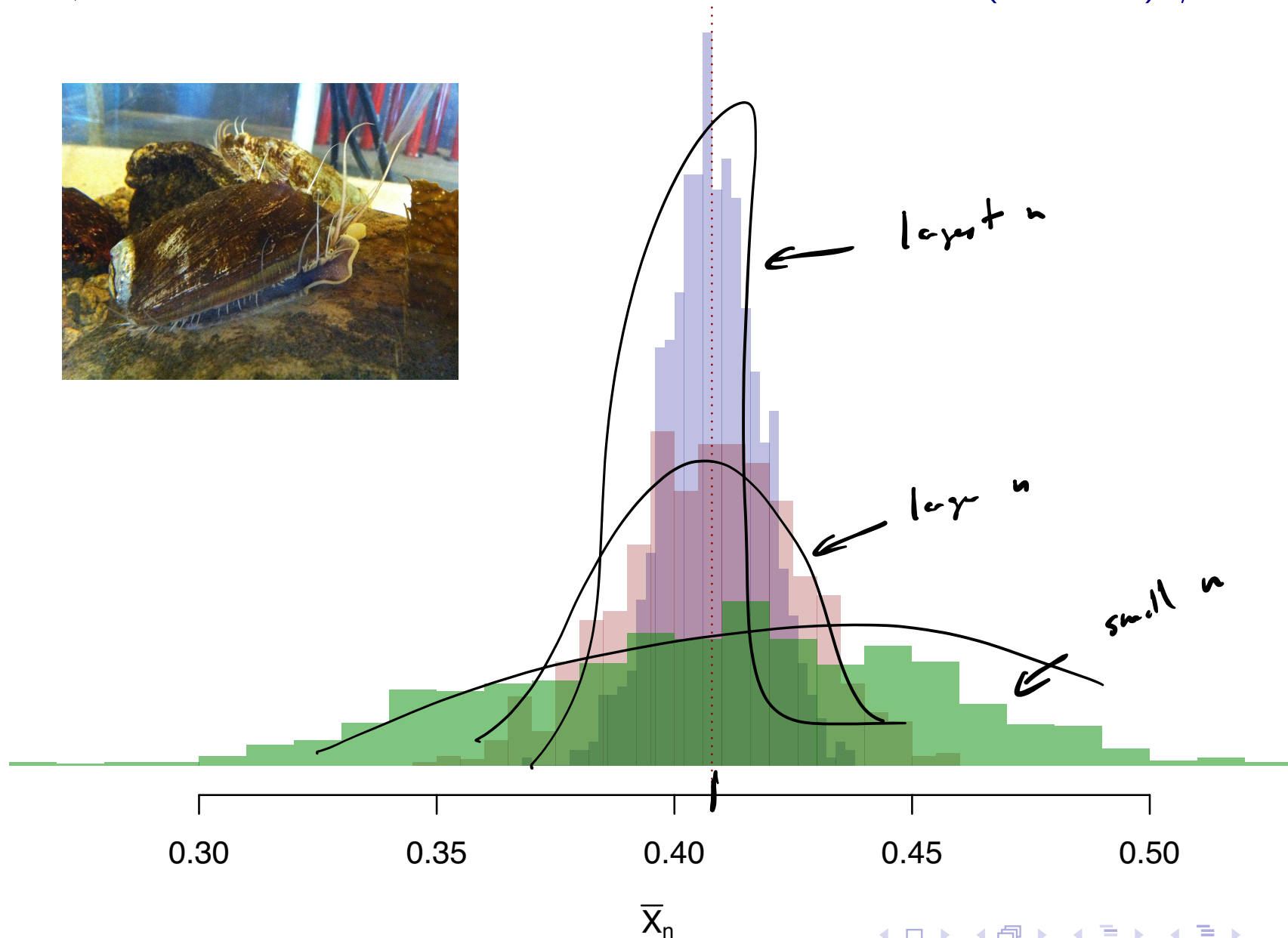Histogram of $\hat{p}_n$ with n = 100

Normal Q–Q plot of $\sqrt{n}(\hat{p}_n - p)/\sqrt{p(1-p)}$

If $X_1, \ldots, X_n$ a rs of abalone, $\mathbb{E}\bar{X}_n = 0.4079$ and $\mathrm{Var}\,\bar{X}_n = (0.09924)^2/n$.



$\bar{X}_n$

$\bar{X}_n$ has standard deviation $\frac{\sigma}{\sqrt{n}}$.

## Distribution of sample mean when population is Normal

Let $X_1, \ldots, X_n \overset{ind}{\sim} \text{Normal}(\mu, \sigma^2)$. Then $\bar{X}_n \sim \text{Normal}(\mu, \sigma^2/n)$.

Can use this to get probabilities like $P(a < \bar{X}_n < b)$ as follows:

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

$Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$

1. Transform $a$ and $b$ to the $Z$-world (# of standard deviations world):

$$a \mapsto \frac{a - \mu}{\sigma/\sqrt{n}} \quad \text{and} \quad b \mapsto \frac{b - \mu}{\sigma/\sqrt{n}},$$

2. Find

$$P\left( \frac{a - \mu}{\sigma/\sqrt{n}} < Z < \frac{b - \mu}{\sigma/\sqrt{n}} \right).$$

$$\sigma^2 = 2500$$

**Exercise:** Let $X$ = minutes talking on phone in last month of a randomly selected USC student. Assume $X \sim \text{Normal}(\mu = 450, \sigma^2 = 50^2)$.

1. Find $P(|X - 450| > 50)$. $= 0.3174$
2. Find $P(X < 425)$. $= 0.3085$

Now let $\bar{X}_n$ be the mean talk time from $n = 9$ randomly selected students.

$$\bar{X}_n \sim \text{Normal}\left(\mu = 450, \frac{50^2}{9}\right)$$

1. Find $P(|\bar{X}_n - 450| > 50)$.
2. Find $P(\bar{X}_n < 425)$.

① $P\left( \underbrace{|X - 450|}_{\text{distance of X from 450}} > 50 \right)$ = $P\Big( X$ is more than 50 minutes from 450

$X \sim \text{Normal} \left( \mu = 450, \ \sigma^2 = 50^2 \right)$

$= 1 - P\left( |X - 450| \leq 50 \right)$

$= 1 - P\left( -50 \leq X - 450 \leq 50 \right)$

$= 1 - P\left( 400 \leq X \leq 500 \right)$

$\underbrace{\hspace{4cm}}_{0.6826}$

$= .3174$



$\sigma = 50$

$\mu = 450$

400    500

0.3413

$Z \sim N(0,1)$

0

$\dfrac{400 - 450}{50} = -1$      $\dfrac{500 - 450}{50} = 1$

$P\left( |X - 450| > 50 \right) = 2(.3413) = .6826.$



$P(X < 425) = 0.3085$

?

425    450

0.3085

$Z \sim N(0,1)$

0.1915 from z-table

.3085

0

$\dfrac{425 - 450}{50} = -\frac{1}{2}$     $\frac{1}{2}$

Now let $\bar{X}_n$ be the mean talk time from $n = 9$ randomly selected students.

① Find $P(|\bar{X}_n - 450| > 50)$.

② Find $P(\bar{X}_n < 425)$.

$$\bar{X}_n \sim \text{Normal}\left(\mu = 450, \frac{50^2}{9}\right)$$

② $P\left(|\bar{X}_n - 450| > 50\right) = P\left(\bar{X}_n \text{ is more than } 50 \text{ minutes from the mean}\right)$
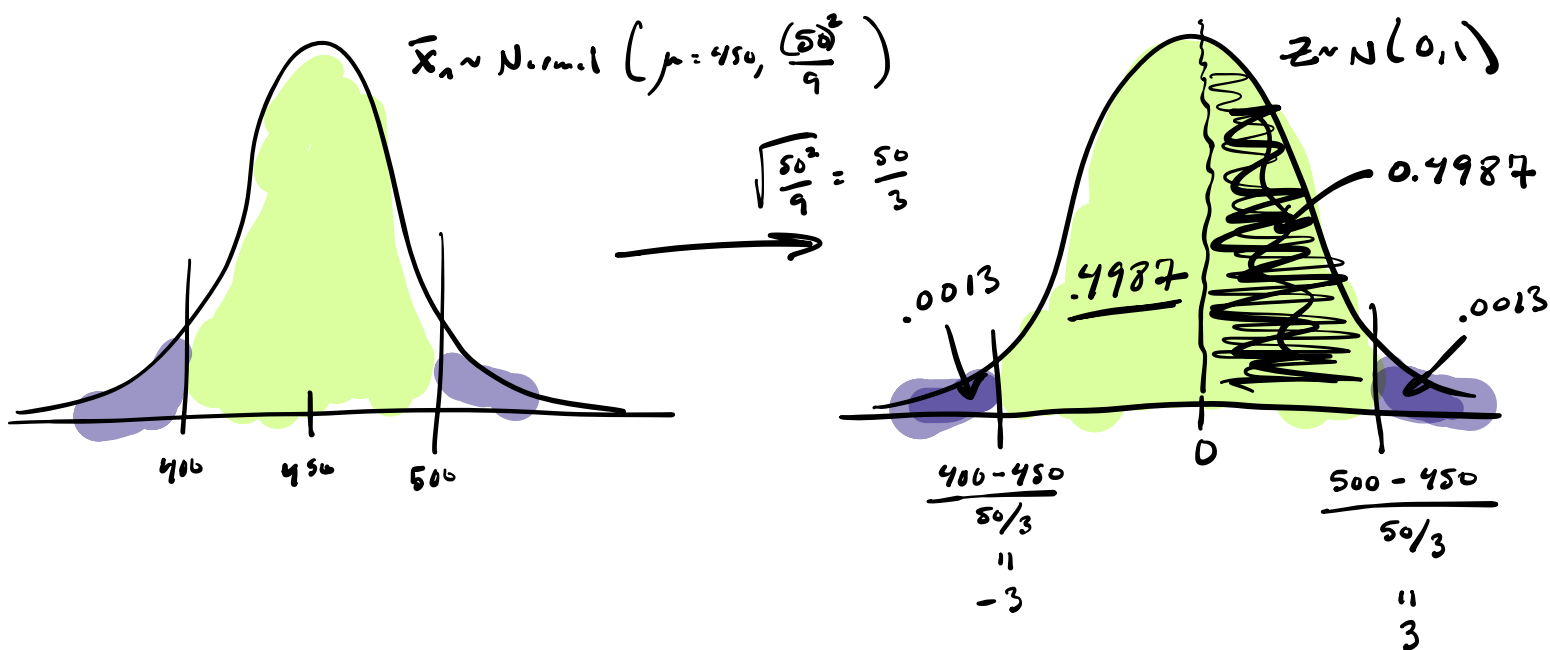
$$= 1 - P\left(|\bar{X}_n - 450| \leq 50\right)$$

$$= 1 - P\left(400 \leq \bar{X}_n \leq 500\right)$$

$$\bar{X}_n \sim \text{Normal}\left(\mu = 450, \frac{(50)^2}{9}\right)$$

$$\sqrt{\frac{50^2}{9}} = \frac{50}{3}$$

$Z \sim N(0,1)$

$0.4987$

$.0013$   $.4987$   $.0013$

$\frac{400-450}{50/3} = -3$

$0$

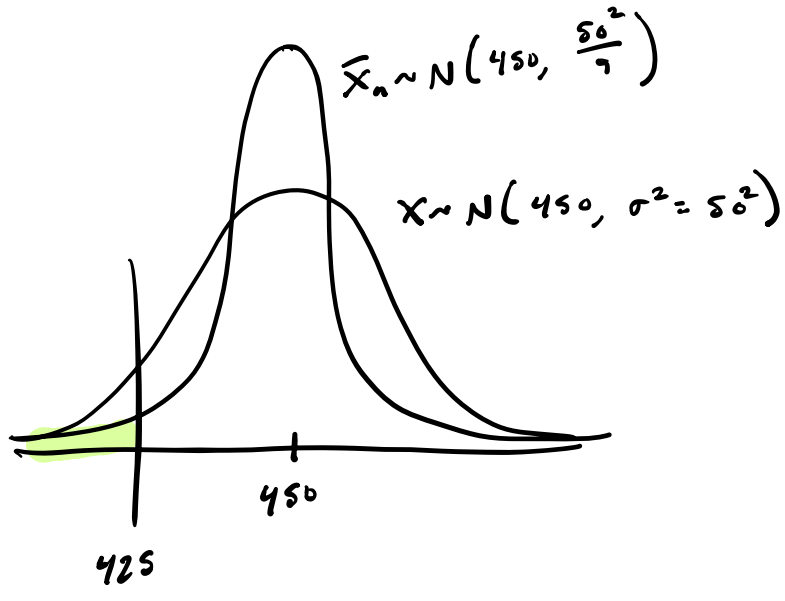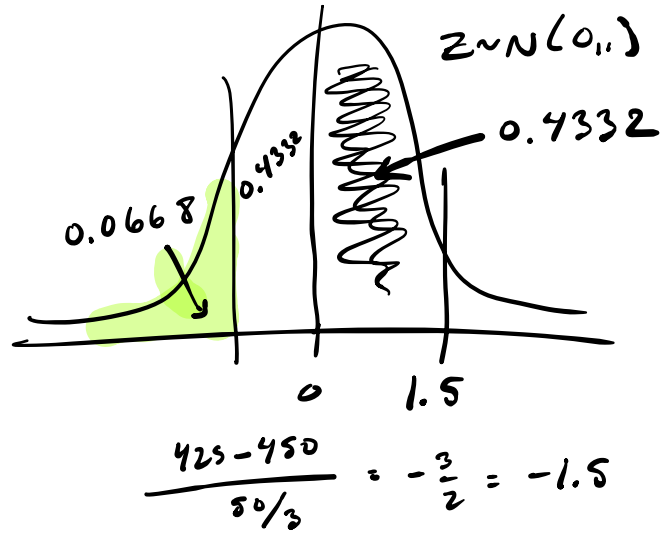$\frac{500-450}{50/3} = 3$

400   450   500

$$Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

$$P\left(|\bar{X}_n - 450| > 50\right) = 2(0.0013) = 0.0026.$$

$\mu = 450$

② $P\left( \bar{X}_n < 425 \right) = 0.0668.$

$\bar{X}_n \sim N\left( 450, \frac{50^2}{9} \right)$



425    450

$Z \sim N(0,1)$

0.4332

0.0668    0.4332

0    1.5

$\dfrac{425 - 450}{50/3} = -\dfrac{3}{2} = -1.5$



$\bar{X}_n \sim N\left( 450, \frac{50^2}{9} \right)$

$X \sim N\left( 450, \sigma^2 = 50^2 \right)$

450

425

$$P\left( |X - 450| > 50 \right) = 0.3174$$



Normal $\left( \mu = 450, \quad \dfrac{50^2}{9} \right)$

$N \left( \mu = 450, \quad \sigma^2 = 50^2 \right)$

400   450   500

**Exercise:** You sell jars of baby food labelled as weighing 4oz ≈ 113g. Suppose your process results in jar weights with the Normal($\m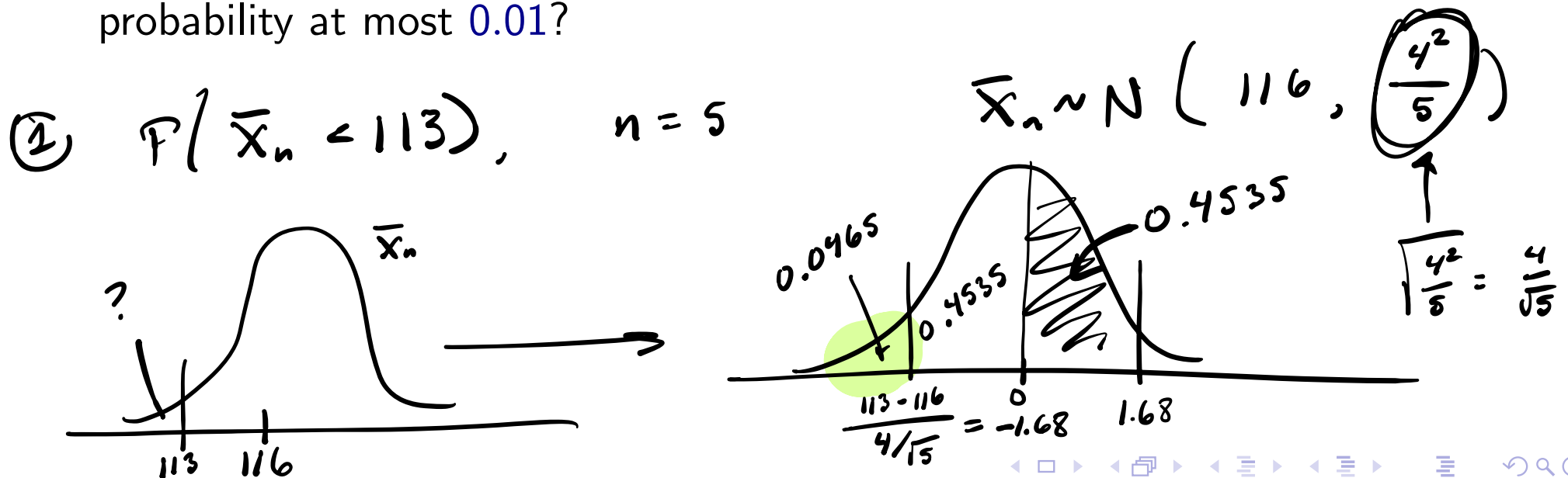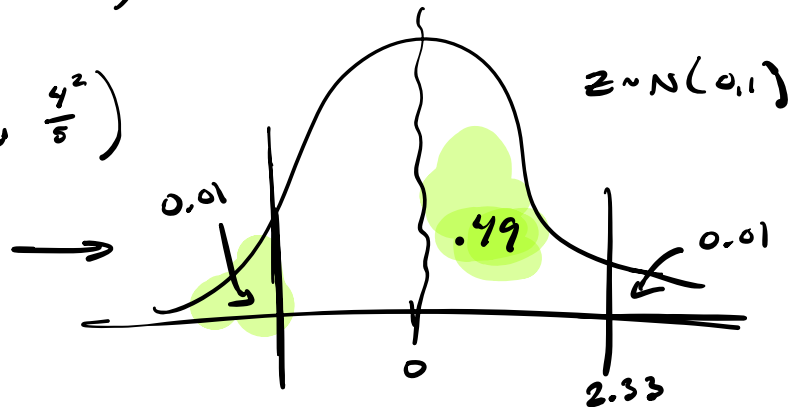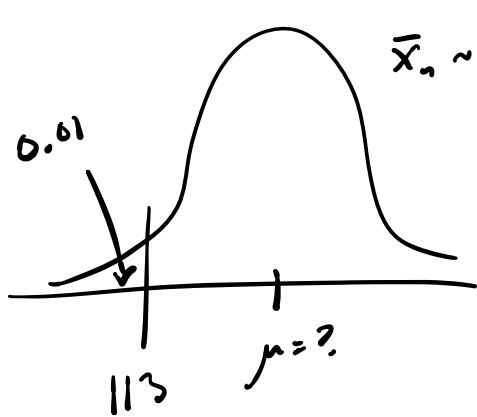u = 116, \sigma^2 = 4^2$) distribution. A regulator will sample 5 jars and fine you if the average weight is less than 113g.

1. With what probability will you get fined?

2. To what must you increase $\mu$ so that you are fined with prob. at most 0.01?

3. Keeping $\mu = 116$g, to what must you reduce $\sigma$ so that you are fined with probability at most 0.01?

① $P\left(\bar{X}_n < 113\right),$  $n = 5$  $\bar{X}_n \sim N\left(116, \left(\dfrac{4^2}{5}\right)\right)$

$\bar{X}_n$

?

$0.0965$  $0.4535$  $0.4535$

$\dfrac{113-116}{4/\sqrt{5}} = -1.68$  $0$  $1.68$

$\sqrt{\dfrac{4^2}{5}} = \dfrac{4}{\sqrt{5}}$

113  116

$$P\left(\bar{X}_n < 113\right) = 0.0465.$$

② $P\left(\bar{X}_n < 113\right) \overset{set}{=} 0.01$, $\mu$ is unknown.

$$\bar{X}_n \sim N\left(\mu = ?, \frac{4^2}{5}\right)$$

$z \sim N(0,1)$

0.01

0.01    .49    0.01

113    $\mu = ?$

0    2.33

$$\frac{113 - \mu}{4/\sqrt{5}} = z_{0.01}^2 = -2.33$$

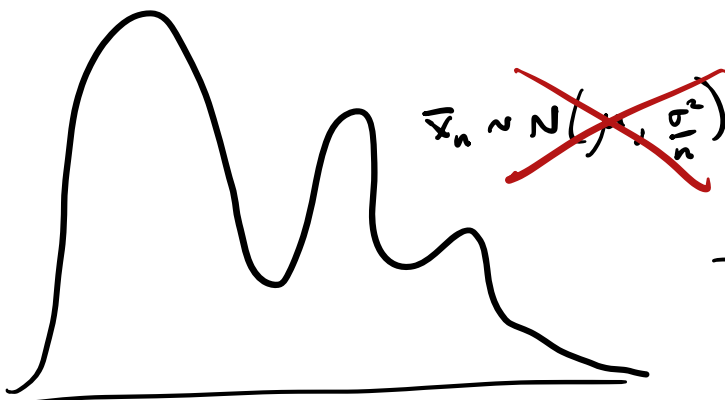$$Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

$$\frac{113 - \mu}{4/\sqrt{5}} = -2.33$$

$\Longleftrightarrow$

$$113 - \mu = (-2.33)\frac{4}{\sqrt{5}}$$

$\Longrightarrow$ $113 + 2.33\frac{4}{\sqrt{5}} = \mu$

$\Longrightarrow$ $\mu = 117.17.$

$$\bar{X}_n \sim \cancel{N\left(\phantom{x}, \frac{\sigma^2}{n}\right)}$$

$Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$

$z \sim N(0,1)$

↑ If not normal, but n is large, can still do the Z stuff.

## Central Limit Theorem

Let $X_1, \ldots, X_n$ be a rs from a dist. with mean $\mu$ and variance $\sigma^2 < \infty$. Then
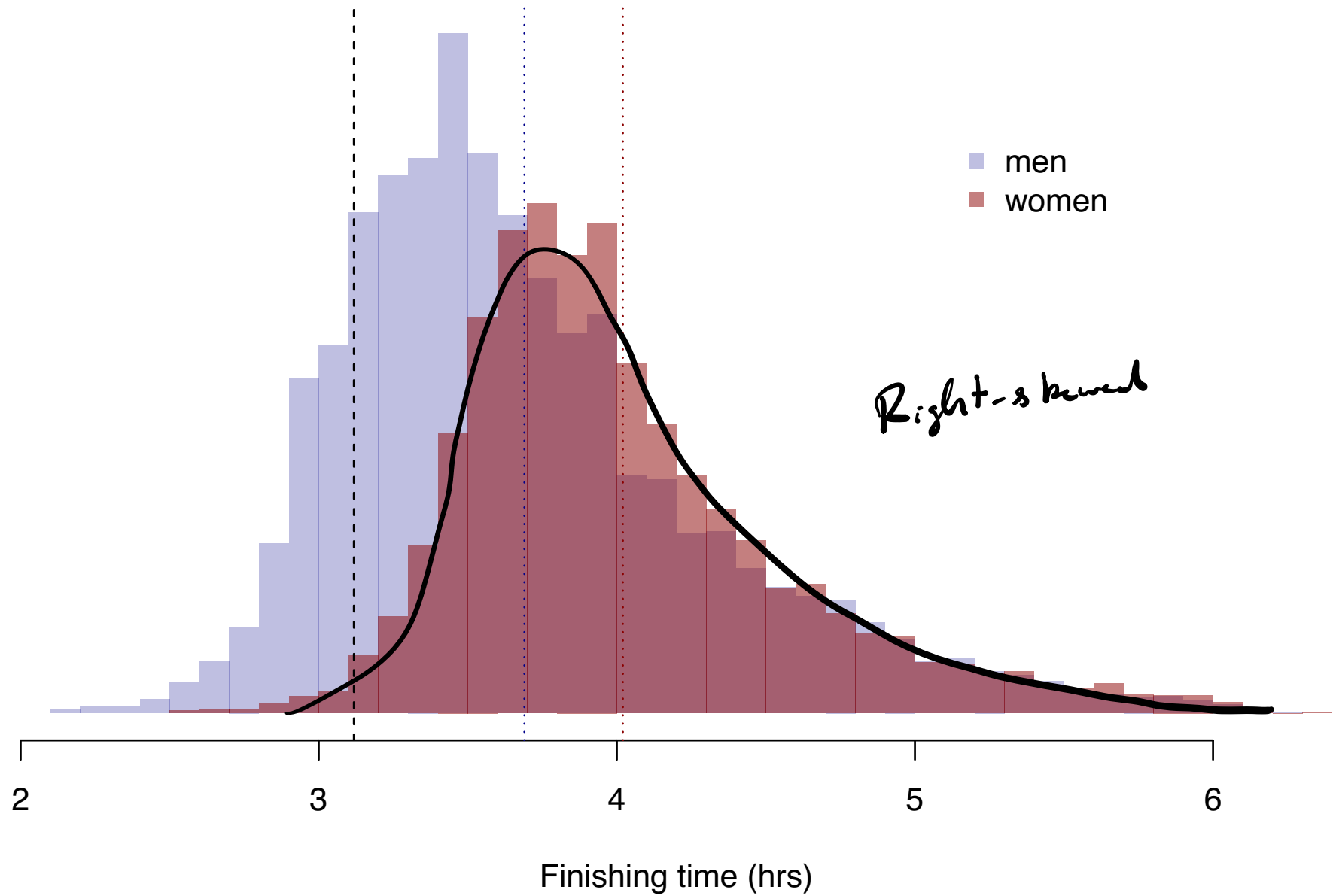
$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$ behaves more and more like $Z \sim \text{Normal}(0, 1)$
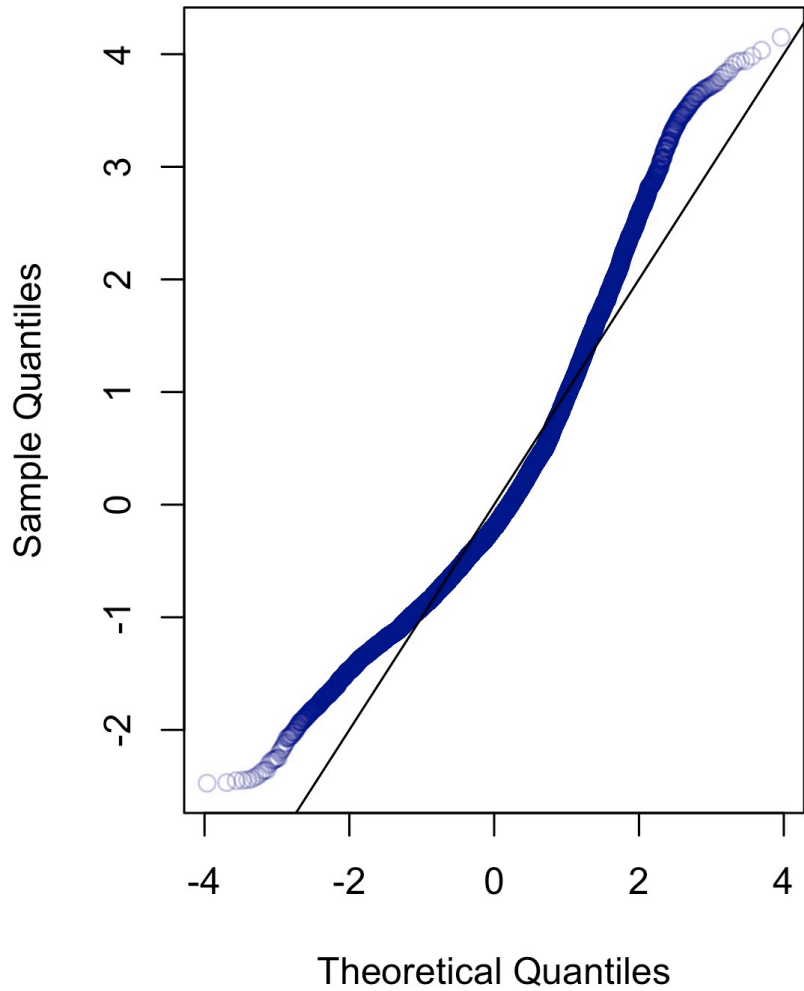
for larger and larger $n$.

This means that for large $n$ (say $n \geq 30$), we have

$$\bar{X}_n \overset{\text{approx}}{\sim} \text{Normal}\left(\mu, \frac{\sigma^2}{n}\right).$$

2009 Boston Marathon finishing times (hrs)
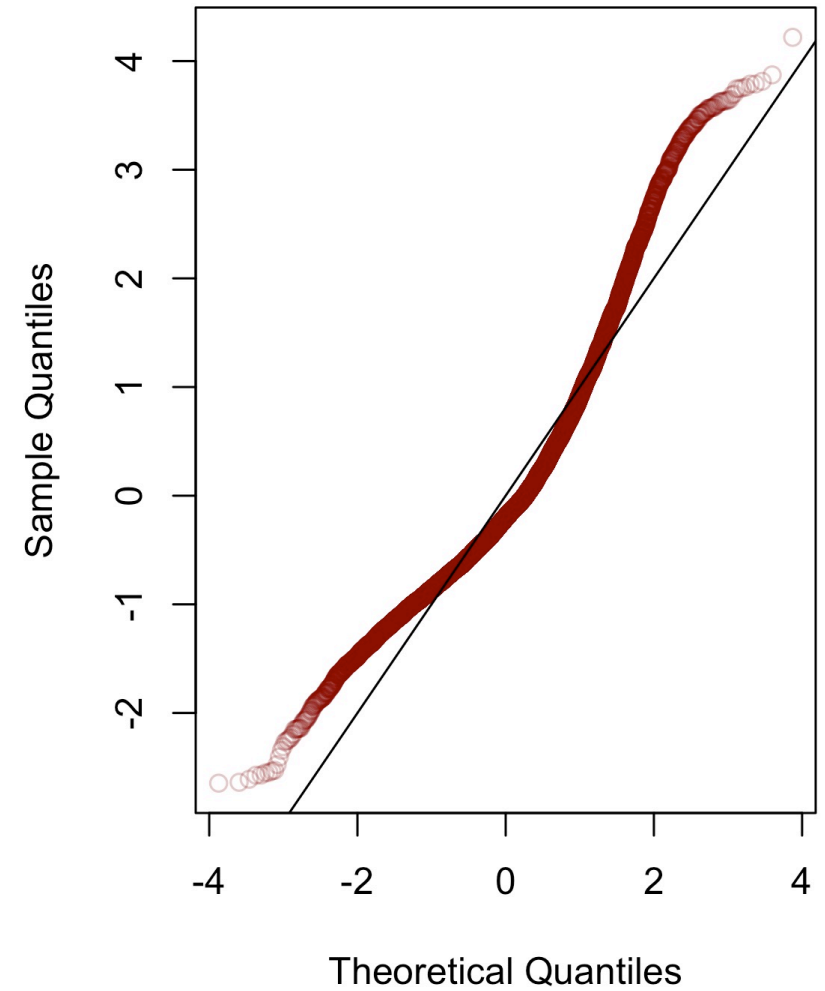
Right-skewed

Finishing time (hrs)

Not Normal.



Normal Q-Q plot for men

Normal Q-Q plot for women

**Exercise:** Women's finishing times for the 2009 Boston Marathon had mean $4.02$ hours and standard deviation $0.555$ hours.
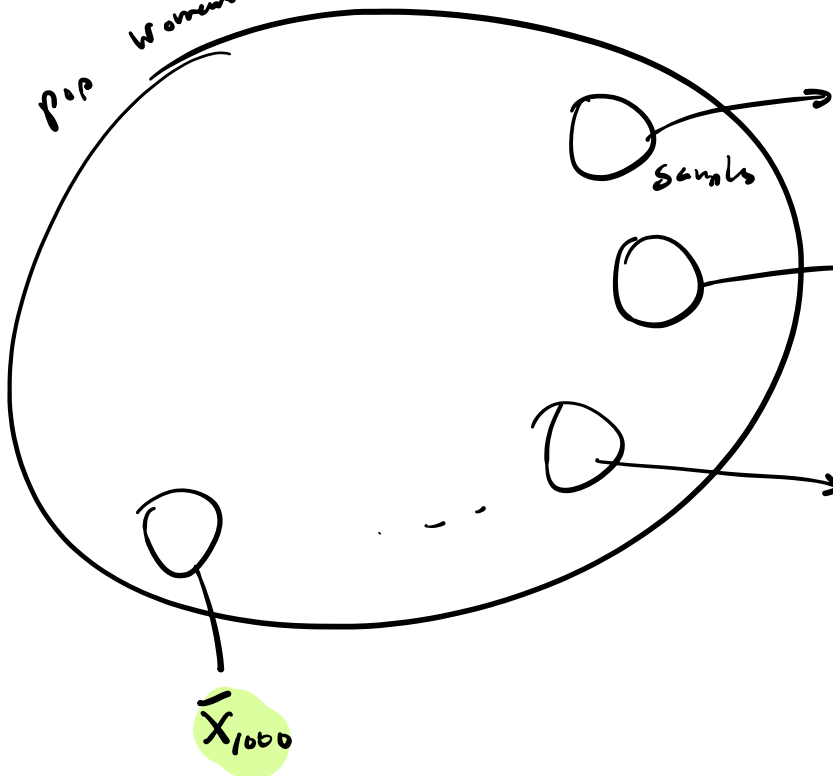
Consider sampling $n = 30$ women and let $\bar{X}_n$ be the mean of their finishing times.

1. Find an approximation to $P(\bar{X}_n < 3.90)$.
2. Find an approximation to $P(\bar{X}_n > 4.25)$.
3. Find an approximation to $P(|\bar{X}_n - 4.02| < 0.2)$.
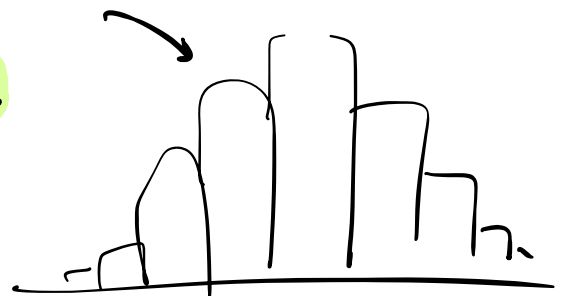
Now use R to draw $1{,}000$ samples of size $n = 30$. <u>`link to women's data`</u>.

1. Make histogram and Normal Q-Q plot of $\bar{X}_n$.
2. Get the probabilities above using the output of the simulation.

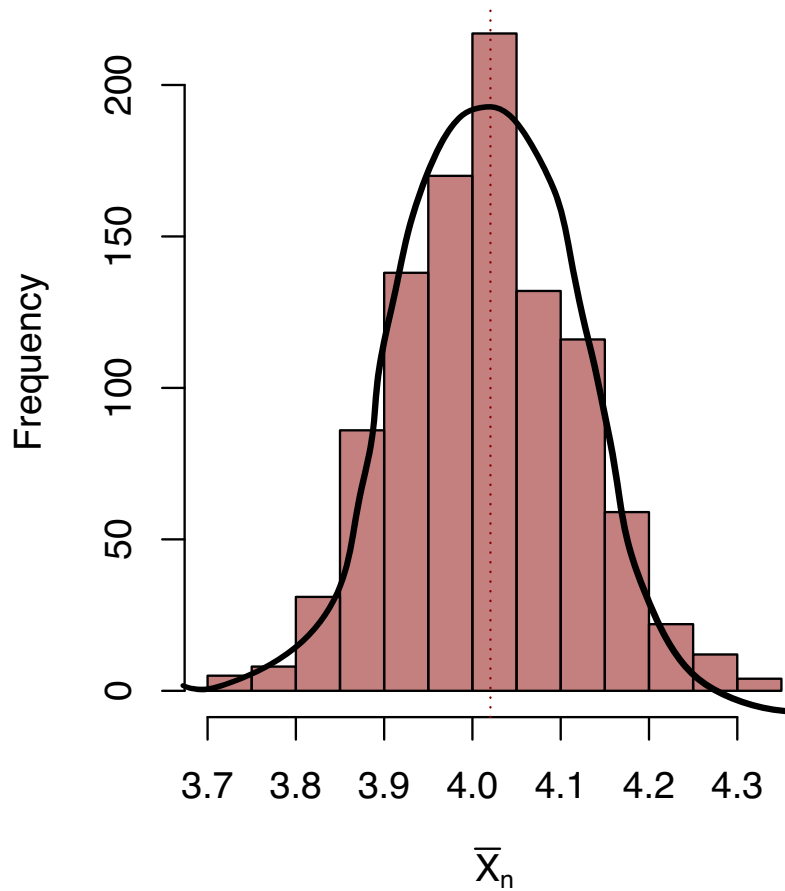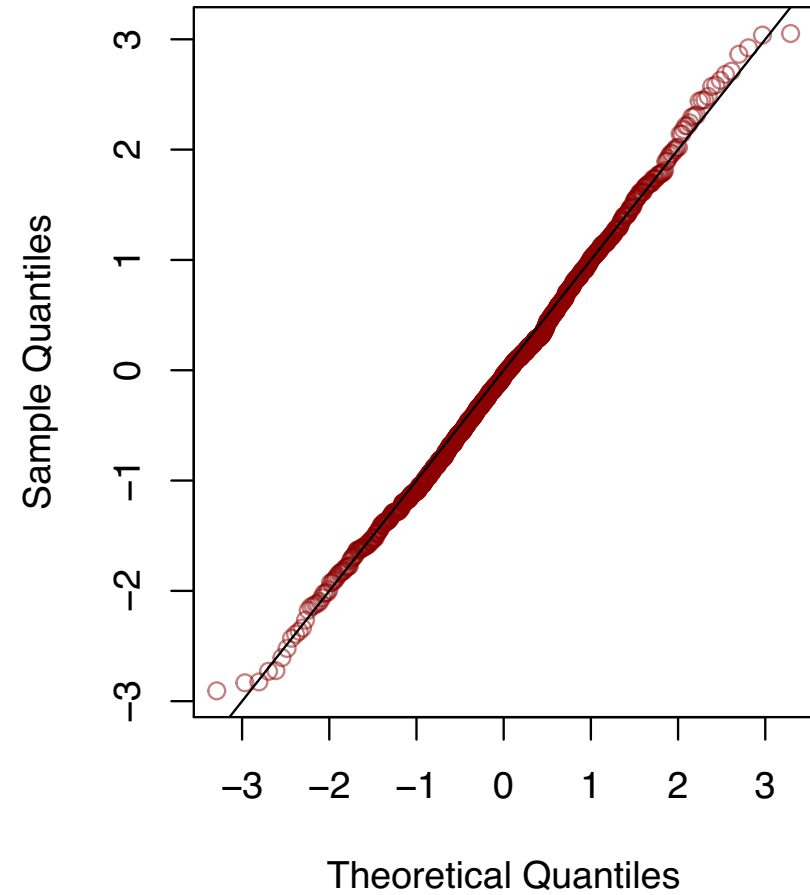Pop women's finishing times   NOT NORMAL

$\overline{X}_1$

sample

$\overline{X}_2$

$n = 30$

$\overline{X}_3$

$\overline{X}_{1000}$

## Histogram of $\overline{X}_n$ with n = 30

## Normal Q–Q plot of $\sqrt{n}(\overline{X}_n - \mu)/\sigma$

$$\frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} \overset{approx}{\sim} Z \sim N(0,1) \quad \text{and} \quad \bar{x}_n \overset{approx}{\sim} N\left(\mu, \frac{\sigma^2}{n}\right).$$

$$X \sim \text{Bernoulli}(p) \quad \Rightarrow \quad EX = p \quad, \quad \text{Var } X = p(1-p).$$

$\underset{\mu}{\uparrow} \qquad \underset{\sigma^2}{\uparrow}$

Central Lim. Thm

$$\frac{\hat{p}_n - p}{\sqrt{\frac{p(1-p)}{n}}} \overset{approx}{\sim} Z \sim N(0,1) \qquad \text{for large } n.$$

Also

$$\hat{p}_n \overset{approx}{\sim} N\left(p, \frac{p(1-p)}{n}\right) \qquad \text{for large } n.$$

p is proportion of 1s
in the population.

population of 0s and 1s.



sample $\rightarrow X_1, \dots X_n$

$$\bar{x}_n = \frac{x_1 + \dots + x_n}{n}$$

$$= \text{proportion of 1s}$$

$$= \hat{p}_n$$

We can apply the Central Limit theorem to proportions.

$$\text{Central Lim Thm}$$
$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \longrightarrow Z \sim N(0,1)$$
$$\text{as} \quad n \to \infty.$$

## Central Limit Theorem for the sample proportion

Let $X_1, \ldots, X_n \overset{\text{ind}}{\sim}$ Bernoulli$(p)$ and let $\hat{p}_n = \bar{X}_n$. Then

For Bernoulli$(p)$,
$$\mu = p$$
$$\text{and} \quad \sigma^2 = p(1-p)$$

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{\hat{p}_n - p}{\sqrt{\underbrace{p(1-p)}_{\sigma^2}/n}}$$

behaves more and more like $Z \sim \text{Normal}(0,1)$

for larger and larger $n$.

This means that for large $n$ (say $np \geq 5$ and $n(1-p) \geq 5$), we have

$$\bar{X}_n \overset{\text{approx}}{\sim} N\left(\mu, \frac{\sigma^2}{n}\right)$$
for large $n$.

$$\hat{p}_n \overset{\text{approx}}{\sim} \text{Normal}\left(p, \frac{p(1-p)}{n}\right).$$

Also: $\sum_{i=1}^{n} X_i = n\hat{p}_n \overset{\text{approx}}{\sim} \text{Normal}\left(p, np(1-p)\right)$ for large $n$.

$$n = 15 \qquad X_1, \ldots, X_n \overset{ind}{\sim} \text{Bernoulli} (p = 0.60)$$

population $\qquad p = 0.60$

sample $\hat{p}_n$ $\qquad n = 15$

**Exercise:** Suppose 60% of USC undergraduates are registered to vote. Consider taking a sample of size $n = 15$. Let $\hat{p}_n$ be the number in your sample who are registered to vote.

1. Find the approximate value of $P(\hat{p}_n > 0.70)$ using the Normal distribution.
2. Find the exact value of $P(\hat{p}_n > 0.70)$ using the Binomial distribution.
3. Find the approximate value of $P(0.30 < \hat{p}_n < 0.80)$ using the Normal dist.
4. Find the exact value of $P(0.30 < \hat{p}_n < 0.80)$ using the Binomial dist.
5. Repeat the above for a sample of size $n = 100$.

(1) $\quad P\left( \hat{p}_n > .70 \right) \approx .2148$

$\hat{p}_n$



$Z \sim N(0,1)$

$.2852$

$.5 - .2852 = .2148$

$p = 0.60$

$0.70$

$0$

$$\frac{.70 - 0.60}{\sqrt{\frac{.6(1-.6)}{15}}} = 0.79$$

$$Z = \frac{X - \mu}{\sigma} \qquad Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \qquad Z = \frac{\hat{p}_n - p}{\sqrt{\frac{p(1-p)}{n}}}$$

② Find $P\left(\hat{p}_n > 0.70\right)$ using Binomial.

$$\hat{p}_n = \frac{X_1 + \cdots + X_n}{n} = \frac{\#\{successes\}}{n}$$

$\iff \quad n\hat{p}_n = \#\{successes\} \sim \text{Binomial}\left(n = 15, \; p = 0.60\right)$

$$P\left(\hat{p}_n > 0.70\right) = P\left(n\hat{p}_n > n \, 0.70\right)$$

$15(.7) \quad 40 \, rts$
$10.5$

$$= P\left(15 \cdot \hat{p}_n > 15(.7)\right)$$

$$= P\left(Y > 10.5\right) \qquad Y \sim \text{Binomial}\left(n = 15, \; p = 0.60\right)$$

$$= P\left(Y > 10\right)$$

$$= 1 - P\left(Y \leq 10\right)$$

$$= 1 - pbinom\left(\underline{10}, \; 15, \; 0.60\right)$$

$$= 0.2173$$

$\hat{p}_n \overset{approx}{\sim} \text{Normal}\left(0.60, \frac{0.6(1-0.6)}{15}\right)$

$Z \sim N(0,1)$

0.2852

0.2148

$\frac{0.70 - 0.60}{\sqrt{\frac{0.60(1-0.60)}{15}}} = 0.79$

0.60

0.70

0

$Z = \dfrac{\hat{p}_n - p}{\sqrt{\dfrac{p(1-p)}{n}}}$

$$P\left(\hat{p}_n > 0.70\right) \approx 0.2148.$$

approx.

② Find $P\left(\hat{p}_n > 0.7\right)$ using **binomial**.

$P\left(\hat{p}_n > 0.7\right) = P\left(\bar{X}_n > 0.7\right)$

$= P\left(\dfrac{X_1 + \cdots + X_n}{n} > 0.7\right)$

$= P\left(X_1 + \cdots + X_n > n(0.7)\right)$

# successes in $n$ trials $\sim \text{Binomial}(n, 0.7)$

## GRADE DISTRIBUTION

| | |
|---|---|
| Greater than 100 | 4 |
| 90 - 100 | 4 |
| 80 - 89 | 8 |
| 70 - 79 | 4 |
| 60 - 69 | 3 |
| 50 - 59 | 4 |
| 40 - 49 | 3 |
| 30 - 39 | 2 |
| 20 - 29 | 2 |
| 10 - 19 | 0 |
| 0 - 9 | 0 |
| Less than 0 | 0 |

14 - 15

Workflow for success
2.5

M - 30 minutes - redo all examples from class.

W - 30 min.

If needed read notes.

Summary of sampling distribution results for $\bar{X}_n$:



$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \overset{\texttt{approx}}{\sim} \texttt{Normal}(0,1)$$

$n \geq 30$

Population
X non - Normal

$n < 30$

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \text{ non-Normal}$$

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

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \texttt{Normal}(0,1)$$

Summary of sampling distribution results for $\hat{p}_n$:

# ds in the samples

$\boxed{n\hat{p}_n} \sim \text{Binomial}(n, p)$

if $n$ is small

$\min\{np, n(1-p)\} < 5$

$\min\{np, n(1-p)\} \geq 5$

if $n$ is large

From central limit theorem.

$$\frac{\hat{p}_n - p}{\sqrt{\frac{p(1-p)}{n}}} \overset{\text{approx}}{\sim} \text{Normal}(0, 1)$$

$$n\hat{p}_n \overset{\text{approx}}{\sim} \text{Normal}(np, np(1-p))$$

$$n\hat{p}_n \sim \text{Binomial}(n, p)$$