

Topic 4 – Normal and sampling distributions

ENVX1002 Introduction to Statistical Methods

Dr. Floris van Ogtrop
The University of Sydney

Jan 2024



Topic 4 Outline – Normal and sampling distributions

- Example
- Normal distribution
- Other continuous distributions
- Sampling distribution

Learning Outcomes

- Understand what a (probability) distribution is:
 - the properties of a continuous distribution.
- Use Normal Distribution to understand/describe data
 - Be able to standardise a Normal;
 - Calculate probabilities based on Normal Distribution using R.
- Know that there are other continuous distributions useful in hypothesis testing.
- Distinguish between population, sample and sampling distributions;
- Distinguish between a standard deviation and standard error of the mean;
- Describe the Central Limit Theorem;
- Use R and Excel to calculate the standard error and probabilities associated with sampling distributions;

Types of data

- Numerical
 - Continuous: yield, weight
 - Discrete: weeds per m^2
- Categorical
 - Binary: 2 mutually exclusive categories
 - Ordinal: categories ranked in order
 - Nominal: qualitative data

Example

- The gestation period (in days) for American Simmental cattle is distributed with mean 284.3 and standard deviation 5.52. How often is a calf born a week early?

Wray et al. 1987

► J Anim Sci. 1987 Oct;65(4):970-4. doi: 10.2527/jas1987.654970x.

Analysis of gestation length in American Simmental cattle

N R Wray ¹, R L Quaas, E J Pollak

Affiliations + expand

PMID: 3667470 DOI: 10.2527/jas1987.654970x

Abstract

Records of gestation length (71,461) for Simmental cattle were distributed with mean 284.3 d and standard deviation 5.52 d. Gestation length was found to increase with percent Simmental and was 1.9 d longer for calves born to mature dams than for those born to heifer dams. Bull calves experienced gestation lengths 1.5 d longer than heifer calves. Sire, maternal grandsire, residual and total variances were estimated to be 2.42, .58, 22.78 and 25.78 d², respectively, by Henderson's Method III. Heritability of gestation length was calculated to be .374 from the sire variance and .09 from the maternal grandsire variance. Direct additive genetic variance was considered to be of greater importance than maternal additive genetic variance. Correlations between the evaluations of sires for gestation length and heifer calving ease, birth weight and weaning weight were .26, .26 and .13, respectively.

FULL TEXT LINKS

OXFORD
ACADEMIC

ACTIONS

“ Cite

🔖 Collections

SHARE



PAGE NAVIGATION

◀ Title & authors

Abstract



What is a distribution

- In our case we are generally referring to a distribution function
 - This is a function (or model) that describes the probability that a system will take on value or set of values x
- For any variable X , we describe probabilities by
 - Discrete variables: probability distribution function $P(X = x)$
 - Continuous variables: probability density function $f(x)$
 - Discrete and Continuous variables: cumulative density function $F(x) = P(X \leq x)$

Properties of a continuous distribution

- For any continuous distribution
 - There is an infinite number of possible values;
 - These values may be within a fixed interval. For example, male human heights (in cm) belong to $[54.6, 272]$.

Human height

- Specific values in a continuous distribution have probability 0. For example, the likelihood of measuring a Simmental cow at exactly 450kg is zero. This is because there are potentially an infinite number of other weights that are higher or lower than 450 kg so we say that measuring exactly 450 has a very very small probability which is equivalent to zero
- The total of all the probabilities = must be 1. (Total area under the pdf)

The Normal Distribution

- The Normal Distribution is **super important** because it occurs everywhere! It naturally describes many natural phenomenon and is a great for modelling the sample mean.
- It is a symmetric bell-shaped variable with two parameters μ and σ^2 such that:

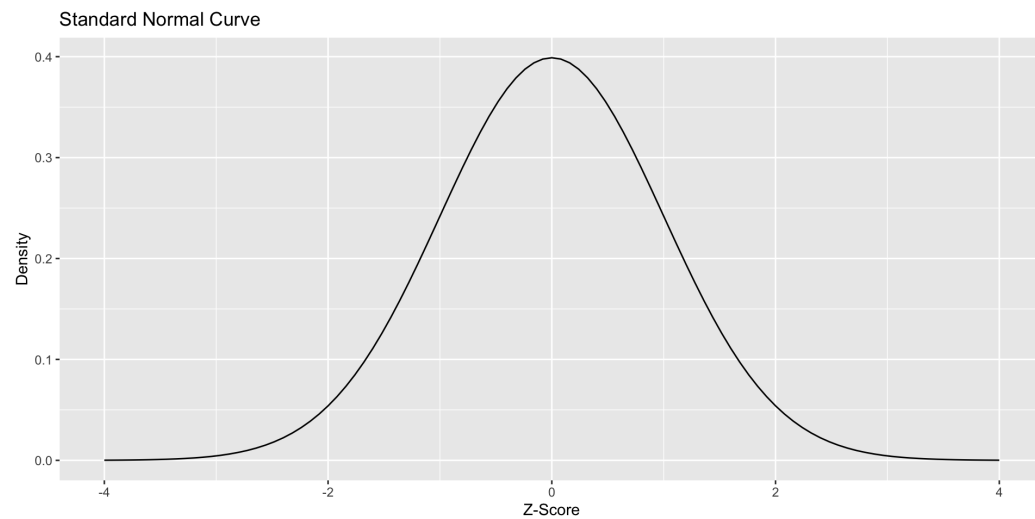
$$X \sim N(\mu, \sigma^2)$$

The Standard Normal Curve

- The standard normal curve is one where the mean = 0, and variance = 1

$$X \sim N(\mu = 0, \sigma^2 = 1)$$

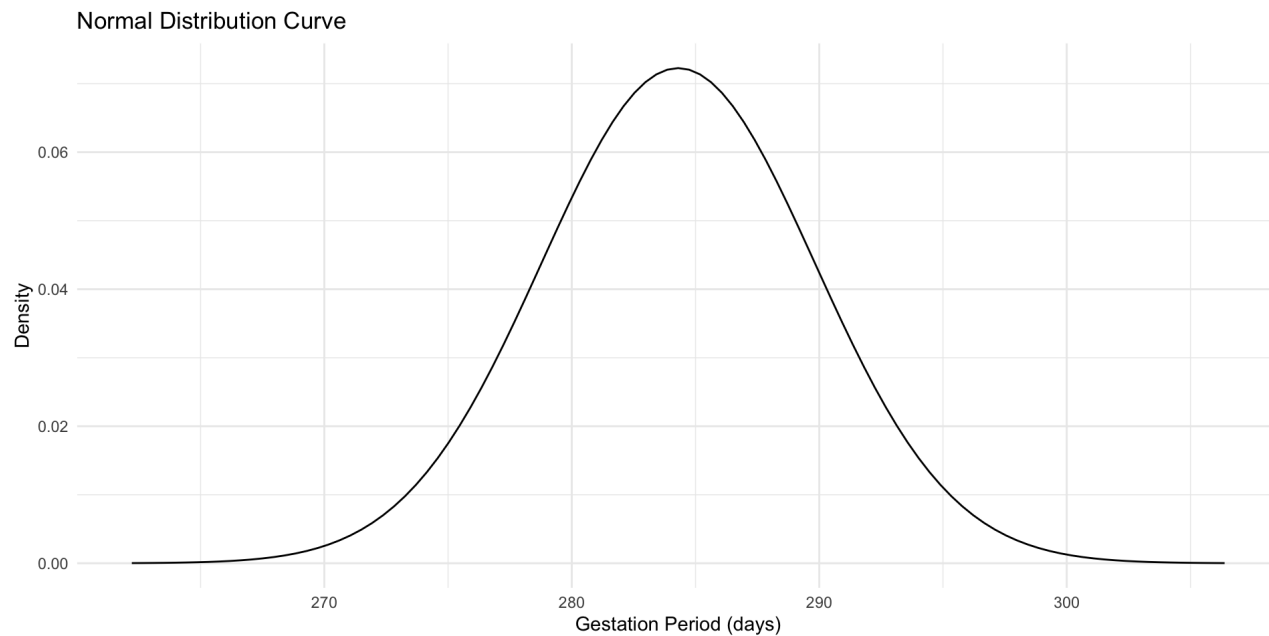
```
1 # Load the necessary libraries
2 library(infer); library(ggplot2)
3 # Create a sequence of values for the x-axis
4 x_values <- seq(-4, 4, by = 0.01)
5 # Create a data frame to hold these values
6 data_frame <- data.frame(x_values)
7 # Use ggplot to plot the standard normal curve
8 ggplot(data_frame, aes(x = x_values)) +
9   stat_function(fun = dnorm, args = list(mean = 0, sd = 1)) +
10  labs(title = "Standard Normal Curve", x = "Z-Score", y = "Density")
```



The General Normal Curve

- Simmental cattle gestation times...

```
1 # Define the mean and standard deviation
2 mean <- 284.3
3 sd <- 5.52
4 # Create a sequence of x values
5 x_values <- seq(mean - 4*sd, mean + 4*sd, length.out = 1000)
6 # Create a ggplot
7 ggplot(data.frame(x_values), aes(x = x_values)) +
8   stat_function(fun = dnorm, args = list(mean = mean, sd = sd)) +
9   labs(title = "Normal Distribution Curve", x = "Gestation Period (days)", y = "Density") +
10  theme_minimal()
```



The General Normal Distribution

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

- PDF

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

- CDF

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(y)dy$$

Types of Normal Probabilities

- There are 3 type of probabilities that we are interested in:
 - Tail probabilities (lower and upper) = Cumulative probabilities
 - Interval probabilities;
 - Inverse probabilities.

Normal distribution in R

Commands for General Normal: $X \sim N(\mu, \sigma^2)$

	R
Probability density function (pdf) $f(x)$	$dnorm(x, \mu, \sigma)$
Lower Tail probabilities (cumulative, CDF) $F(x) = P(X \leq x)$	$pnorm(x, \mu, \sigma)$
Interval probabilities $P(a \leq X \leq b)$	$pnorm(b, \mu, \sigma) - pnorm(a, \mu, \sigma)$
Inverse probabilities Find q such that $P(X \leq q) = \%$	$qnorm(q, \mu, \sigma)$

Types of Normal Probabilities

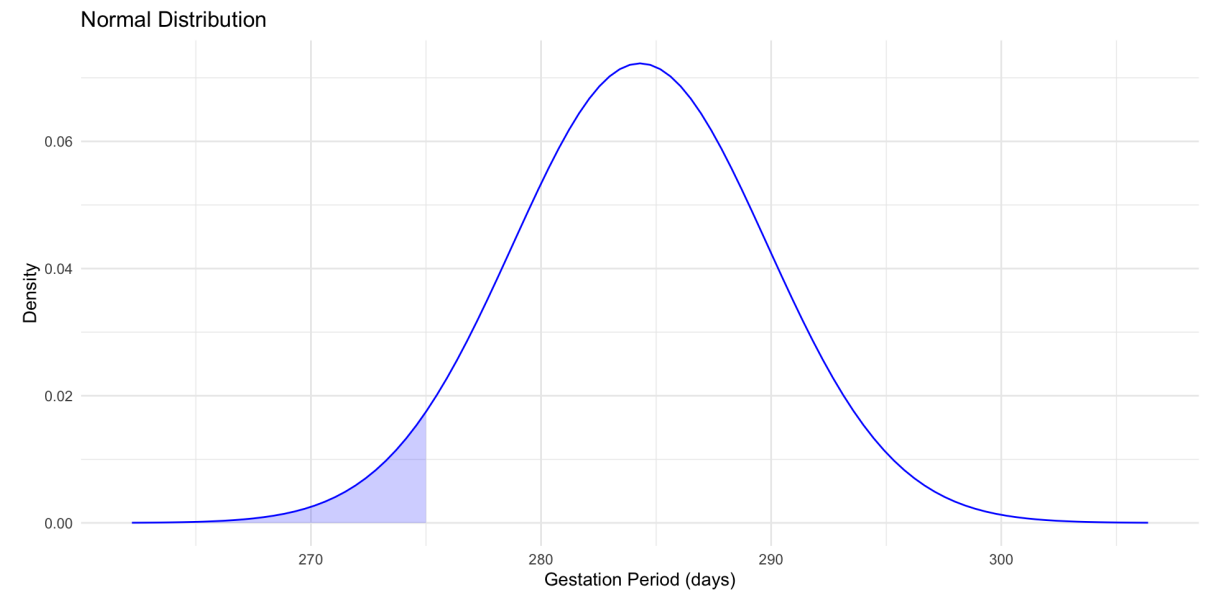
Lower: $P(X \leq 275)$

```
1 pnorm(275, 284.3, 5.52)
```

```
[1] 0.04601526
```

```
1 # Load the ggplot2 library
2 library(ggplot2)
3
4 # Define the mean and standard deviation
5 mean <- 284.3
6 sd <- 5.52
7
8 # Define the x-value up to which the area should be shaded
9 x_shade <- 275
10
11 # Create a data frame for the x-values
12 x_values <- seq(mean - 4*sd, mean + 4*sd, length.out = 100)
13 df <- data.frame(x = x_values)
14
15 # Create the ggplot object
16 p <- ggplot(df, aes(x)) +
17   stat_function(fun = dnorm, args = list(mean = mean, sd = sd)) +
18   geom_area(stat = "function", fun = dnorm, args = list(mean = mean, sd = sd),
19             xlim = c(min(x_values), x_shade), fill = "blue")
20 labs(title = "Normal Distribution", x = "Gestation Period (days)")
21 theme_minimal()
```

```
1 # Print the plot
2 print(p)
```



Types of Normal Probabilities

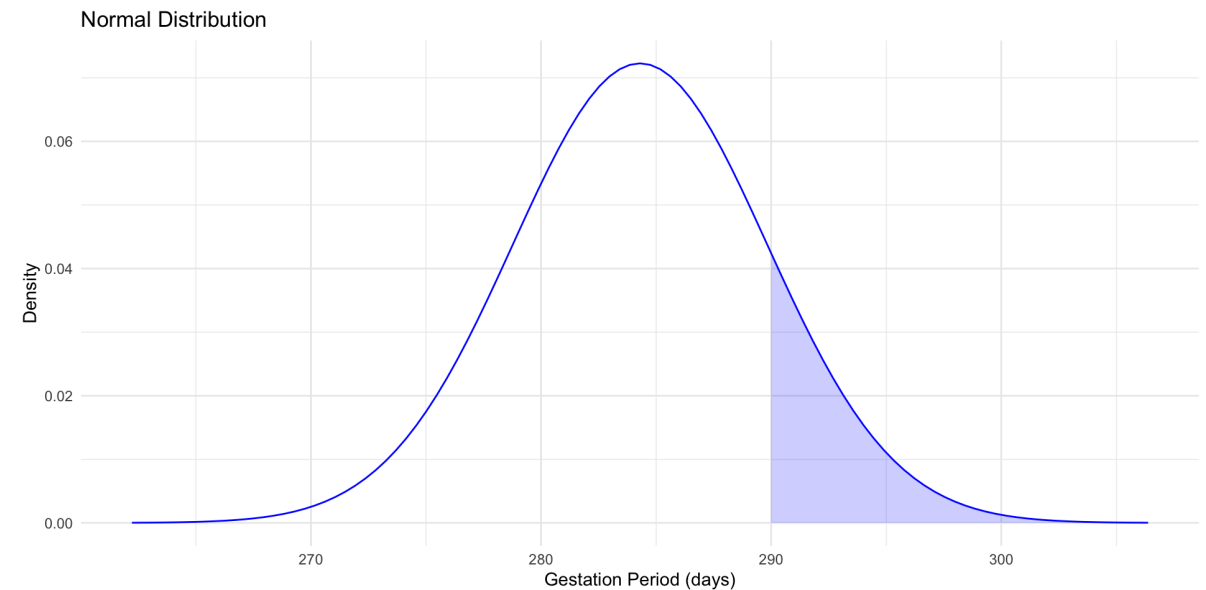
Upper: $P(X \geq 290)$

```
1 1- pnorm(290, 284.3, 5.52)
```

```
[1] 0.1508935
```

```
1 # Load the ggplot2 library
2 library(ggplot2)
3
4 # Define the mean and standard deviation
5 mean <- 284.3
6 sd <- 5.52
7
8 # Define the x-value up to which the area should be shaded
9 x_shade <- 290
10
11 # Create a data frame for the x-values
12 x_values <- seq(mean - 4*sd, mean + 4*sd, length.out = 100)
13 df <- data.frame(x = x_values)
14
15 # Create the ggplot object
16 p <- ggplot(df, aes(x)) +
17   stat_function(fun = dnorm, args = list(mean = mean, sd = sd)) +
18   geom_area(stat = "function", fun = dnorm, args = list(mean = mean, sd = sd),
19             xlim = c(max(x_values), x_shade), fill = "blue")
20 labs(title = "Normal Distribution", x = "Gestation Period (days)")
21 theme_minimal()
```

```
1 # Print the plot
2 print(p)
```



Types of Normal Probabilities

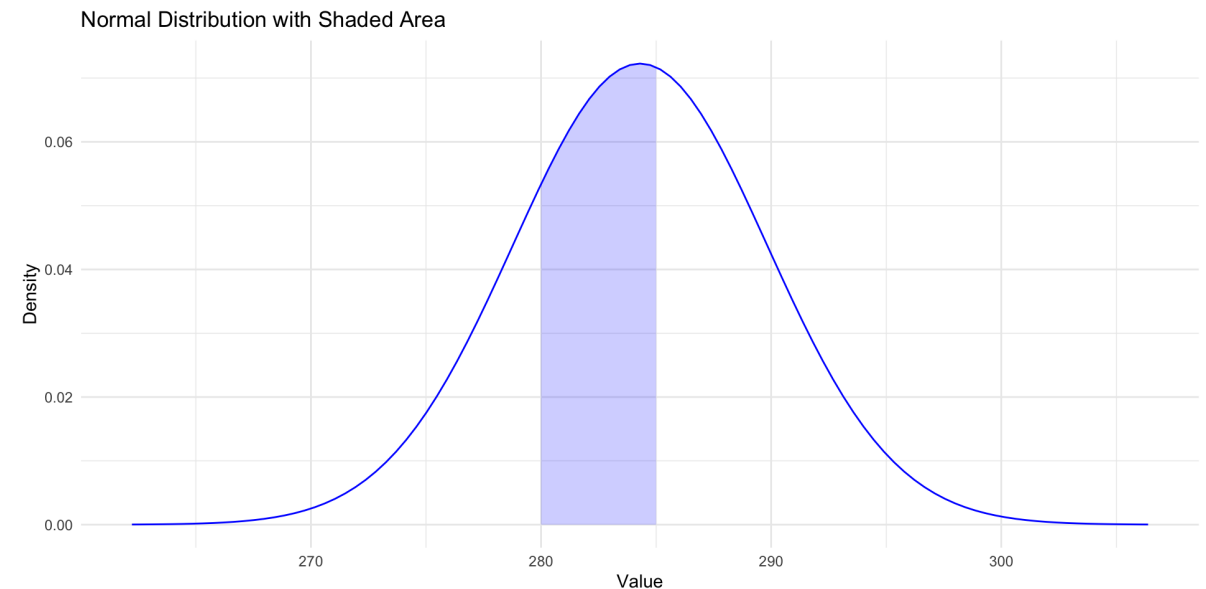
Interval: $P(280 \leq X \leq 285)$

```
1 pnorm(285, 284.3, 5.52)-pnorm(280, 284.3, 5.52)
```

```
[1] 0.3324611
```

```
1 # Load the ggplot2 library
2 library(ggplot2)
3
4 # Define the mean and standard deviation
5 mean <- 284.3
6 sd <- 5.52
7
8 # Define the range for the shaded area
9 lower_bound <- 280
10 upper_bound <- 285
11
12 # Create a data frame for the x values
13 x_values <- seq(mean - 4*sd, mean + 4*sd, length.out = 100)
14 df <- data.frame(x = x_values)
15
16 # Create the ggplot object
17 p <- ggplot(df, aes(x)) +
18   stat_function(fun = dnorm, args = list(mean = mean, sd = sd)) +
19   geom_area(stat = "function", fun = dnorm, args = list(mean = mean, sd = sd),
20           xlim = c(lower_bound, upper_bound), fill = "blue", alpha = 0.2)
21   labs(title = "Normal Distribution with Shaded Area", x = "Value")
22   theme_minimal()
```

```
1 # Print the plot
2 print(p)
```



Types of Normal Probabilities

Inverse: if we know the shaded area = 0.9 (90%), what is x ?

$$P(X \leq x) = 0.9$$

```
1 qnorm(0.9, mean = 284.3, sd = 5.52)
```

```
[1] 291.3742
```

Example

- Let's return to our example for American Simmental cattle where $X \sim N(284.3, 5.522)$,
- What is the probability of a gestation time less than 275 days.

So we need to calculate the lower Tail probability such that: $P(X \leq 275)$

```
1 pnorm(275, 284.3, 5.52)
```

```
[1] 0.04601526
```

Back to the Standard Normal Curve

- Sometimes it is useful to standardise “data” as it allows us to compare samples that are drawn from populations that may have different means and standard deviations
- Luckily for us we can standardise any general normal distribution $X \sim N(\mu, \sigma^2)$ to a standard normal distribution $X \sim N(0, 1)$
- This was also useful as we could use a set of standard normal tables to calculate probabilities (before computers were readily available).

$$P(X \leq x) = P\left(\frac{X-\mu}{\sigma}, \frac{x-\mu}{\sigma}\right) = P\left(Z \leq \frac{x-\mu}{\sigma}\right)$$

Standard Normal Curve

- Example: if $X \sim N(10, 9)$ find $P(X \leq 14)$

$$P(X \leq x) = P\left(\frac{X-\mu}{\sigma}, \frac{x-\mu}{\sigma}\right) = P\left(Z \leq \frac{x-\mu}{\sigma}\right)$$

$$P(X \leq 14) = P\left(\frac{X-10}{\sqrt{9}}, \frac{14-10}{\sqrt{9}}\right) = P\left(Z \leq \frac{4}{3}\right)$$

```
1 pnorm(14,10,3)
```

```
[1] 0.9087888
```

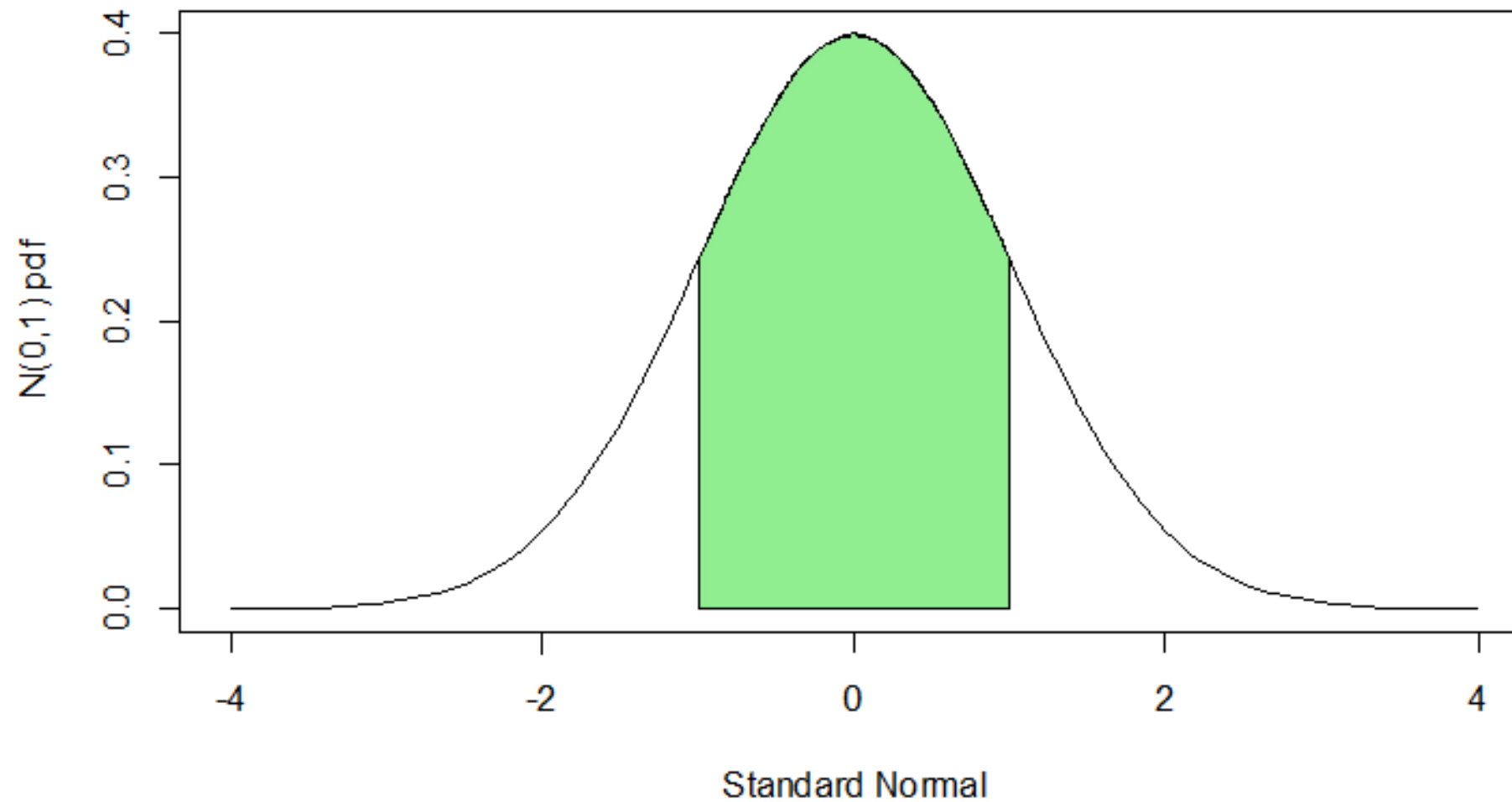
```
1 pnorm(4/3,0,1)
```

```
[1] 0.9087888
```

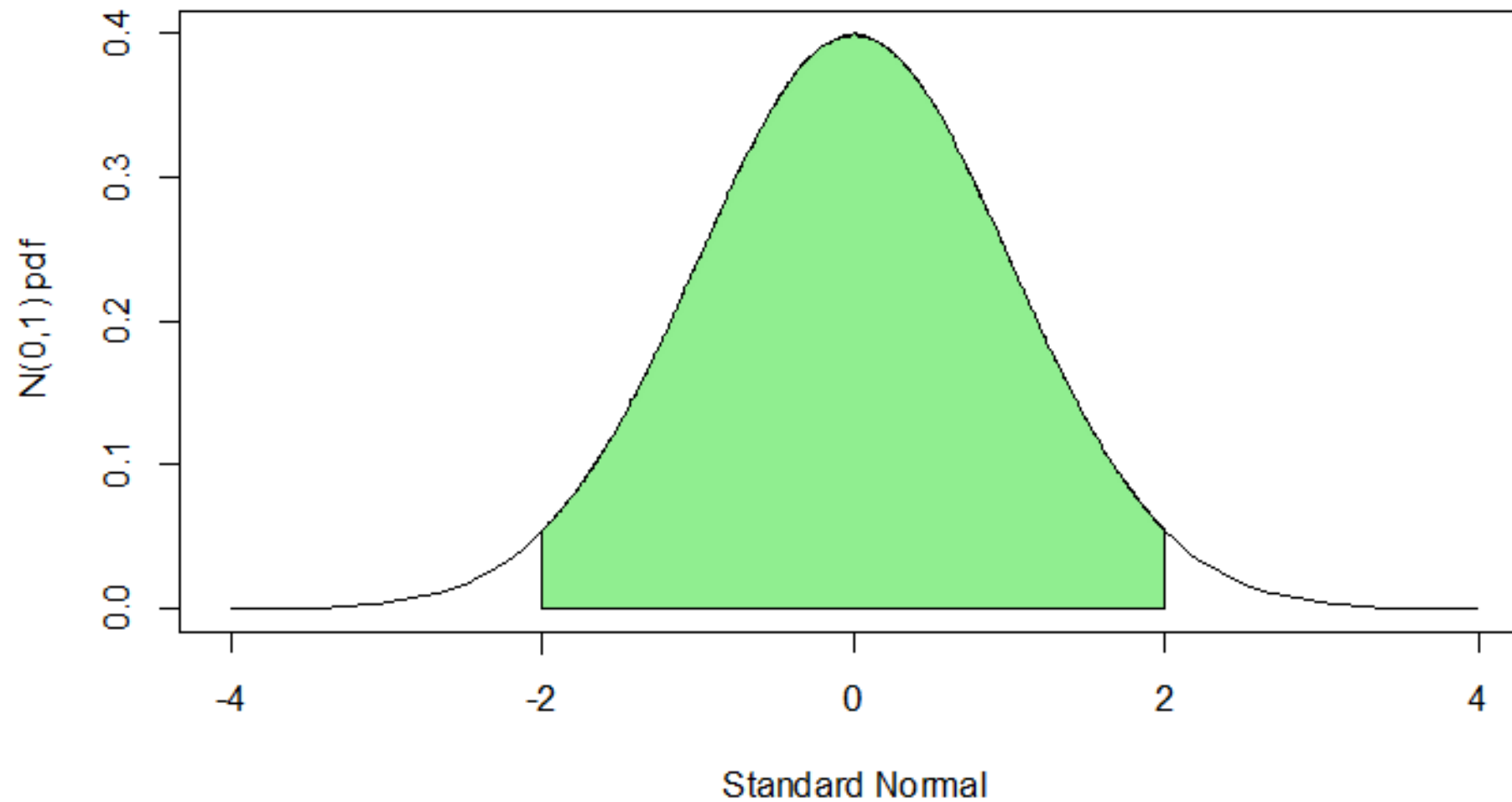
```
1 pnorm(4/3)
```

```
[1] 0.9087888
```

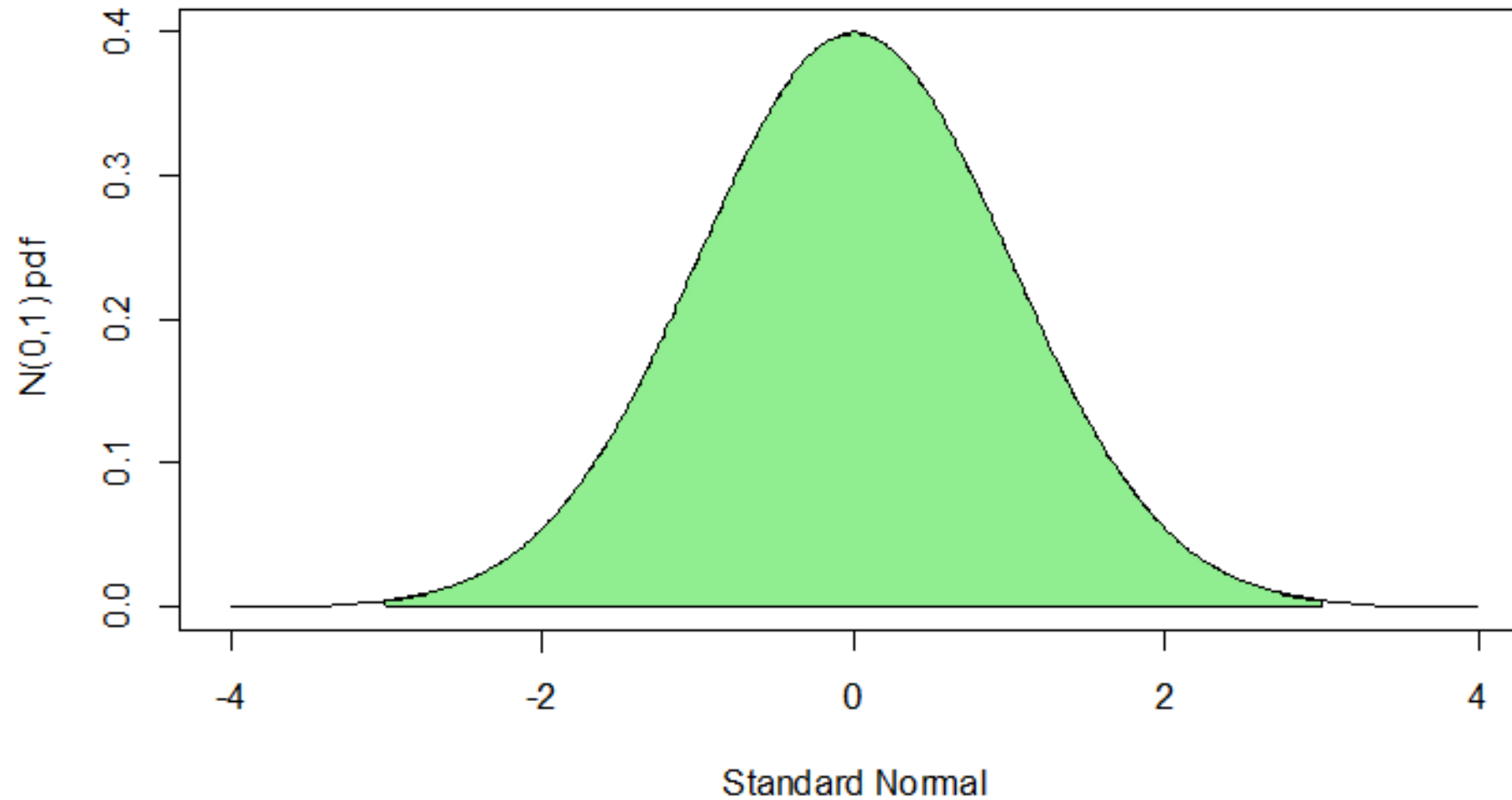
Percentiles of the Standard Normal Curve



Percentiles of the Standard Normal Curve

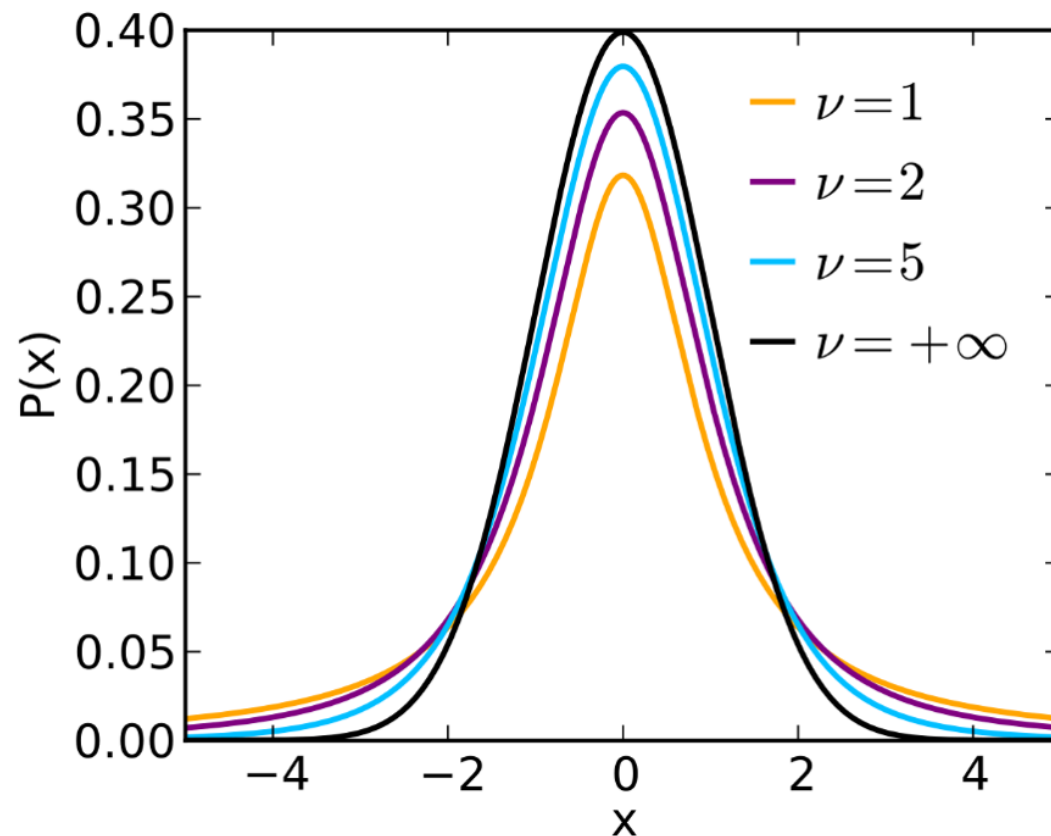


Percentiles of the Standard Normal Curve



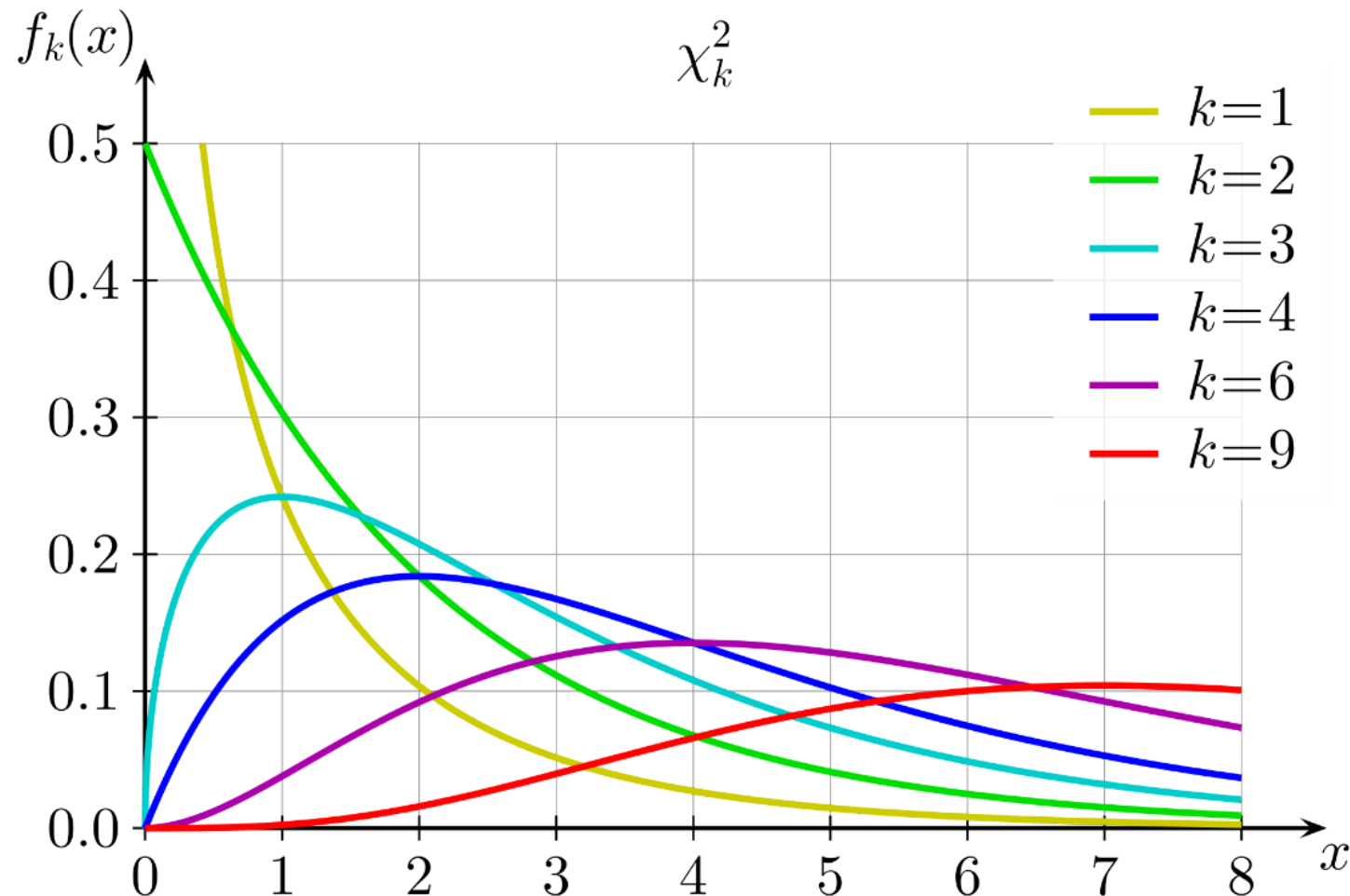
Not so Normal Distributions

- Student T: The Student T distribution models a symmetric bell-shaped variable with thicker tails than a Normal.
- We say the variable $X \sim t_n$, with n degrees of freedom.
- It has an extra parameter = n which is related to population size
- The T distribution is used for the 1 and 2 sample T-tests which are really important in the next few weeks.



Not so Normal Distributions

- The Chi-Squared distribution models a variable which can only take positive values and is skewed in distribution.
- We say the variable $X \sim \chi_n^2$, with n degrees of freedom.
- The Chi-Squared distribution is used for the Chi-Squared Test which you will cover in the next few weeks.



Sampling distributions

- Rye grass root growth (in mg dry weight) follows the distribution $X \sim N(300, 502)$.
- 1. One measurement is taken: how likely is it that the dry weight exceeds 320 mg?
- 2. 10 measurements are taken: how likely is it that the sample mean exceeds 320 mg?



Sampling distributions

- Here, we are dealing with 2 distributions:
 - 1. Measurement: $X \sim N(300, 50^2)$
 - 2. Sample Mean of 10 measurements: $\bar{X} = \frac{1}{10} \sum_{i=1}^{10} X_i \sim \dots$

How does the sampling distribution occur?

- http://onlinestatbook.com/stat_sim/sampling_dist/
- We have a population X
 - We take a sample of size n and we calculate the mean \bar{x}_1
 - We take another sample of size n and we calculate the mean \bar{x}_2
 - We take another sample of size n and we calculate the mean \bar{x}_3 ... If we sample all possibilities, then the sampling distribution of $\bar{X} = \frac{1}{10} \sum_{i=1}^{10} X_i$ is the distribution of $\{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots\}$

Distribution for a sample mean

- if $X \sim N(\mu, \sigma^2)$
- then $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$
- Note that we call
 - σ the standard deviation such that $sd(X) = \sigma$, and
 - σ/\sqrt{n} the standard error such that $sd(\bar{X}) = \sigma/\sqrt{n}$
- The standard error is important for making inference on a sample populations i.e. how close your sample mean \bar{x} is to the population mean μ

Example

- Rye grass root growth (in mg dry weight) follows the distribution $X \sim N(300, 50^2)$.
 - 1. One measurement is taken: how likely is it that the dry weight exceeds 320 mg?
 - 2. 10 measurements are taken: how likely is it that the sample mean exceeds 320 mg?

Example

- 1. $X = \text{Rye grass root growth} \sim N(300, 50^2)$

$$P(X > 320) = P\left(\frac{X-\mu}{\sigma}, \frac{x-\mu}{\sigma}\right) = P\left(\frac{X-300}{50} > \frac{320-300}{50}\right) = P(Z > 0.4) = 1 - P(Z < 0.4) \approx 1 - 0.66 = 0.34$$

```
1 1-pnorm(0.4)
```

```
[1] 0.3445783
```

Example

- 1. \bar{X} = Rye grass root growth $\sim N(300, \frac{50^2}{10})$

$$P(\bar{X} > 320) = P\left(\frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}}, \frac{x-\mu}{\frac{\sigma}{\sqrt{n}}}\right) = P\left(\frac{\bar{X}-300}{\frac{50}{\sqrt{10}}} > \frac{320-300}{\frac{50}{\sqrt{10}}}\right) = P(Z > 1.26) = 1 - P(Z < 1.26) \approx 1 - 0.90 = 0.10$$

```
1 1-pnorm(1.26)
```

```
[1] 0.1038347
```


Central limit theorem

- what about if we sample from any old population such that $X \sim ??(\mu, \sigma^2)$?
- The it follows that $X \approx N(\mu, \sigma^2)$ and this is as a result of the Central Limit Theorem
- The CLT is the most important result in this course, and in much of statistical theory.
- The CLT requires few assumptions:
 - We must have a 'big enough' sample size n ;
 - We must have finite variance $\sigma^2 < \infty$ What is a 'big enough' sample size? Some textbooks give a rule of thumb (eg $n > 25$ or $n > 30$), but it all depends on the type of distribution. If X is reasonably symmetric, then n could be small; if X is highly asymmetric, then n could be larger.

Central limit theorem

- We can demonstrate this in R (you will do it in the practical)

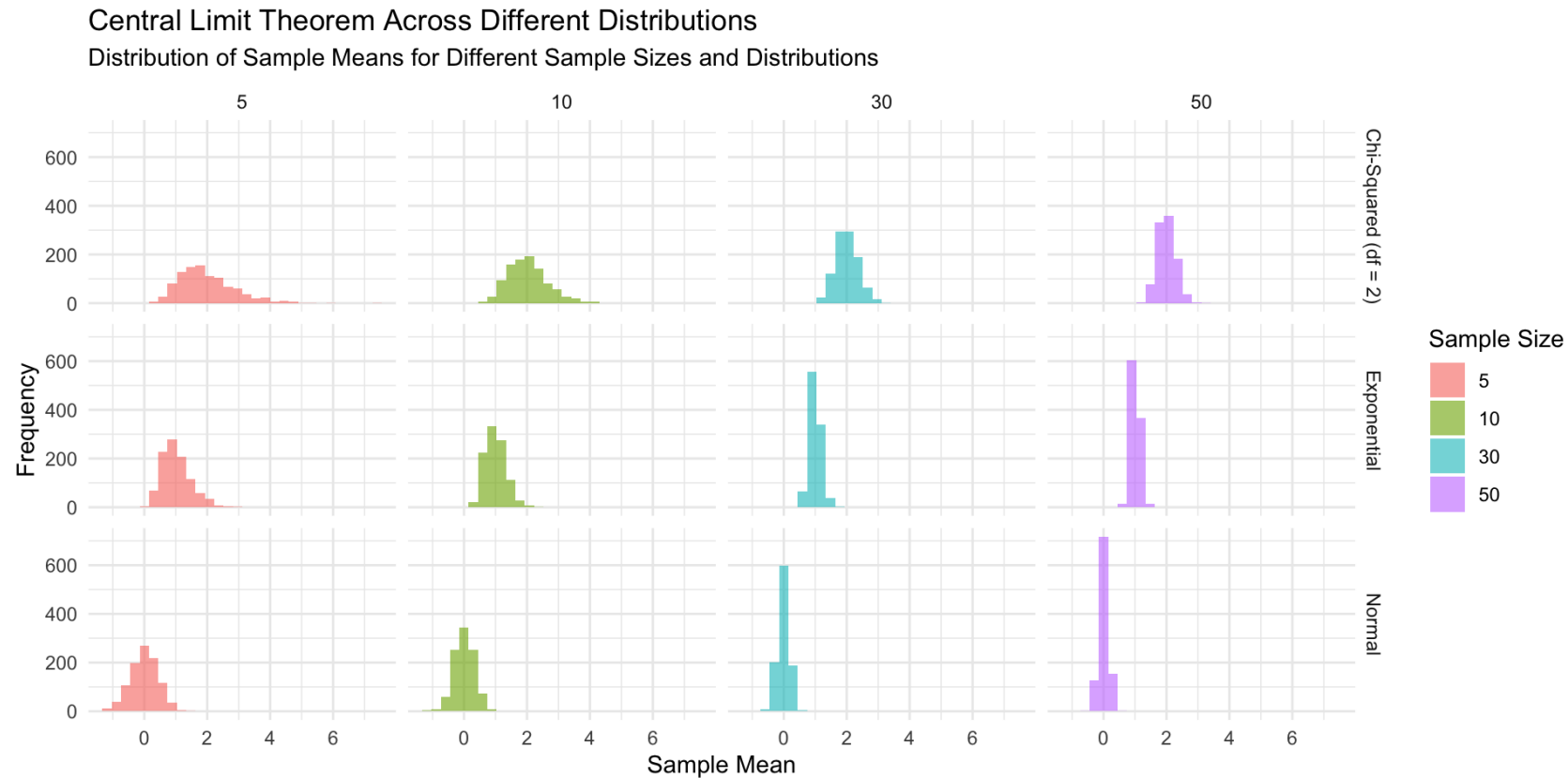
```
1 # Set parameters
2 set.seed(123)
3 number_of_samples <- 1000
4 sample_sizes <- c(5, 10, 30, 50)
5 distributions <- list(
6   "Normal" = rnorm,
7   "Exponential" = rexp,
8   "Chi-Squared (df = 2)" = function(n) rchisq(n, df = 2)
9 )
10
11 # Function to generate sample means from different distributions
12 generate_sample_means <- function(sample_size, number_of_samples, dist_function) {
13   apply(1:number_of_samples, function(x) mean(dist_function(sample_size)))
14 }
15
16 # Generate sample means for different sample sizes and distributions
17 sample_means_list <- lapply(distributions, function(dist_function) {
18   lapply(sample_sizes, generate_sample_means, number_of_samples = number_of_samples, dist_function = dist_function)
19 })
20
21 # Convert to data frame for plotting
22 sample_means_df <- do.call(rbind, lapply(names(distributions), function(dist_name) {
23   do.call(rbind, lapply(1:length(sample_sizes), function(i) {
24     data.frame(
```

Central limit theorem

```

1 # Plot
2 ggplot(sample_means_df, aes(x = `Sample Mean`, fill = as.factor(`Sample Size`))) +
3   geom_histogram(position = "identity", alpha = 0.6, bins = 30) +
4   facet_grid(Distribution ~ `Sample Size`) +
5   theme_minimal() +
6   labs(title = "Central Limit Theorem Across Different Distributions",
7        subtitle = "Distribution of Sample Means for Different Sample Sizes and Distributions",
8        x = "Sample Mean", y = "Frequency", fill = "Sample Size")

```



Thanks!

This presentation is based on the [SOLES Quarto reveal.js template](#) and is licensed under a [Creative Commons Attribution 4.0 International License](#).

title: “ENVX1002-2024-Lecture-Topic04” format: revealjs editor: visual

