

Chapter 6

Continuous Random Variables

Aims

In this chapter we deal with:

- Continuous random variables
- Probability density functions
- Normal distributions
- Working in standard units (i.e., z-scores)

By the end of this chapter you should:

- know the properties of a probability density function
 - know the properties of a Normal probability density function
 - in particular, the 68–95–99.7 rule
 - be able to use software and computer output to solve
 - Normal probability problems
 - Inverse Normal problems
 - understand and be able to apply z-scores
-

Problem: Diagnosing Spina Bifida

A screening test for spina bifida in a foetus involves measuring the concentration of alpha fetoprotein in the mother's urine.

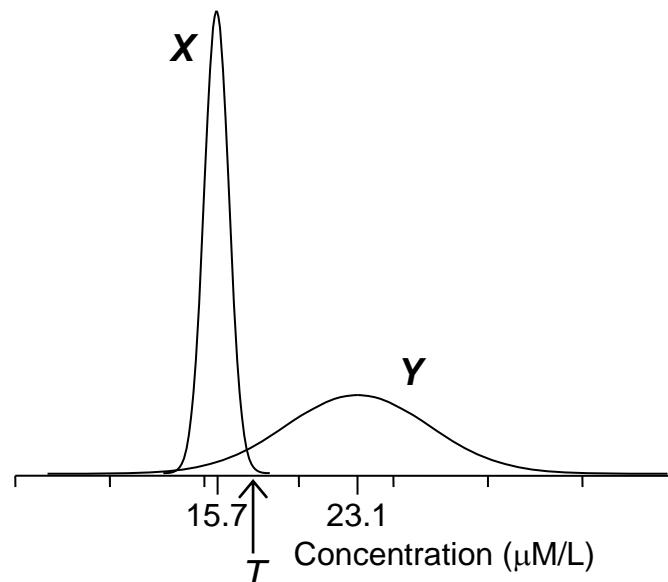
For mothers carrying healthy foetuses, the mean is 15.7 micromoles/litre ($\mu\text{M/L}$) and the standard deviation is 0.7 $\mu\text{M/L}$. For mothers carrying foetuses with spina bifida, the mean is 23.1 $\mu\text{M/L}$ and the standard deviation is 4.1 $\mu\text{M/L}$.

To operate this screening test for spina bifida, medical professionals must set a threshold concentration of alpha fetoprotein, T , say. If the alpha fetoprotein concentration is below T , the foetus is diagnosed as not having spina bifida, whereas if it is above T further testing is required.

Let:

X be the concentration of alpha fetoprotein for a mother carrying a healthy foetus

Y be the concentration of alpha fetoprotein for a mother carrying a foetus with spina bifida.



Suppose T was set at 17.8 $\mu\text{M/L}$:

- What is the probability that a foetus with spina bifida is correctly diagnosed?
- What is the probability that a foetus without spina bifida is correctly diagnosed?

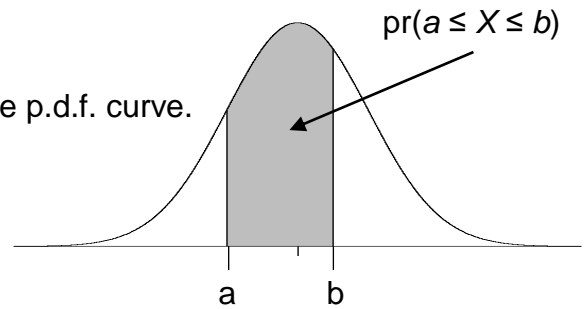
To ensure that 99% of foetuses with spina bifida are correctly diagnosed, at what value should T be set?

Properties of a Probability Density Function (p.d.f.)

1. The p.d.f. curve is always above or on the x-axis.

2. Probabilities are represented by the area under the p.d.f. curve.

$\text{pr}(a \leq X \leq b) =$ the area under the p.d.f. curve between $x = a$ and $x = b$.



3. The total area under a p.d.f. curve = 1.

Endpoints of Intervals

For continuous random variables:

$$\begin{aligned} \text{pr}(a \leq X \leq b) &= \text{pr}(a < X \leq b) \\ &= \text{pr}(a \leq X < b) \\ &= \text{pr}(a < X < b) \end{aligned}$$

In calculations involving continuous random variables, we do not have to worry about whether interval endpoints are included or excluded.

The Normal Distribution

Heights of Male Students

In a recent survey of 544 male Introductory Statistics students the mean height was 177.1cm and the standard deviation (sd) was 7.7cm.

In this survey:

$$\text{mean} - 1 \text{ sd} = 177.1\text{cm} - 7.7\text{cm} = 169.4\text{cm}$$

$$\text{mean} + 1 \text{ sd} = 177.1\text{cm} + 7.7\text{cm} = 184.8\text{cm}$$

95% of students have a height within 1 standard deviation of the mean.

$$\text{mean} - 2 \text{ sd} = 177.1\text{cm} - 2 \times 7.7\text{cm} = 161.7\text{cm}$$

$$\text{mean} + 2 \text{ sd} = 177.1\text{cm} + 2 \times 7.7\text{cm} = 192.5\text{cm}$$

95% of students have a height within 2 standard deviations of the mean.

$$\text{mean} - 3 \text{ sd} = 177.1\text{cm} - 3 \times 7.7\text{cm} = 154.0\text{cm}$$

$$\text{mean} + 3 \text{ sd} = 177.1\text{cm} + 3 \times 7.7\text{cm} = 200.2\text{cm}$$

99.7% of students have a height within 3 standard deviations of the mean.

68–95–99.7 Rule

In a Normal distribution, approximately:

68% of observations are within **1** standard deviation of the mean

95% of observations are within **2** standard deviations of the mean

99.7% of observations are within **3** standard deviations of the mean

Percentage Body Fat for Competitive Cyclists

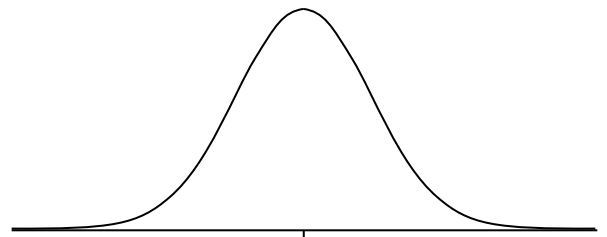
The distribution of percentage body fat for competitive cyclists is modelled by a Normal distribution with a mean of 9 and a standard deviation of 3.

Let X be the percentage body fat of a competitive cyclist.

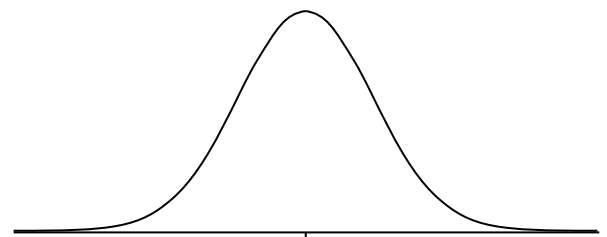
So $X \sim \text{Normal}(\mu = 9, \sigma = 3)$

Based on this model:

- a. approximately 95% of competitive cyclists have percentage body fat somewhere between



- b. for a randomly chosen competitive cyclist, there is a probability of 0.68 that the cyclist has percentage body fat somewhere between



Annual Compound Share Returns for Large Companies

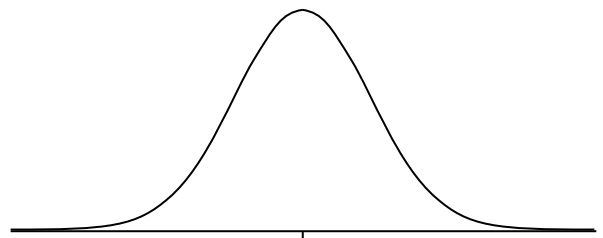
The distribution of annual compound share returns for large companies in the United States (as a percentage) is modelled by a Normal distribution with a mean of 11.0 and a standard deviation of 20.3.

Let X be the annual compound share return for a large company in the United States (as a percentage).

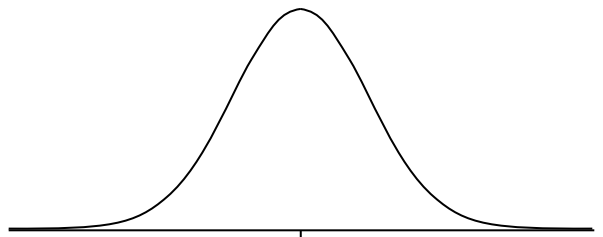
So $X \sim \text{Normal}(\mu = 11.0, \sigma = 20.3)$

Using the computer output given, find the probability that a randomly chosen large US company has an annual compound share return of:

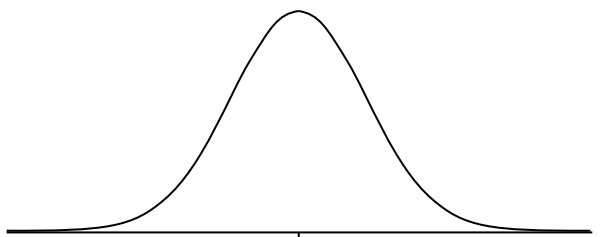
a. less than 0% (i.e., a negative return).



b. more than 50%.



c. between 10% and 30%.



Computer Output

Normal with mean = 11.0 and standard deviation = 20.3

x	P(X ≤ x)
0	0.2940
10	0.4804
30	0.8254
50	0.9726

Spina Bifida

We will model the distributions of alpha fetoprotein with Normal distributions.

Let:

X be the concentration of alpha fetoprotein for a mother carrying a healthy foetus

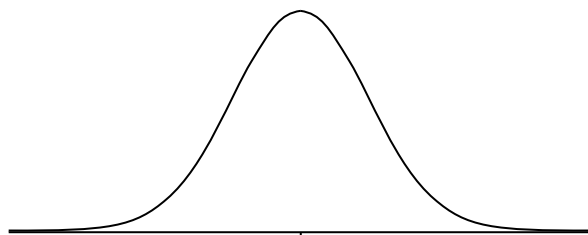
Y be the concentration of alpha fetoprotein for a mother carrying a foetus with spina bifida

So $X \sim \text{Normal}(\mu_X = 15.7, \sigma_X = 0.7)$

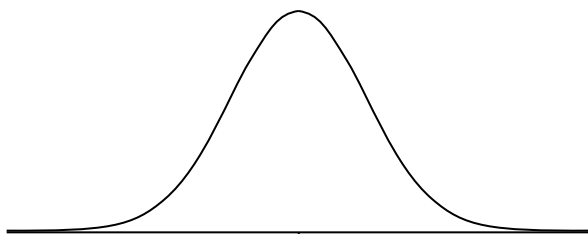
$Y \sim \text{Normal}(\mu_Y = 23.1, \sigma_Y = 4.1)$

Note: The threshold is set at $17.8 \mu\text{M/L}$, i.e., if the concentration of alpha fetoprotein is below $17.8 \mu\text{M/L}$ the foetus would be diagnosed as not having spina bifida.

a. For a foetus without spina bifida, what is the probability that it is correctly diagnosed?



b. What is the probability that a foetus with spina bifida is correctly diagnosed?



Computer Output

Normal with mean = 15.7 and
standard deviation = 0.7

x	P(X <= x)
17.8	0.9987

Normal with mean = 23.1 and
standard deviation = 4.1

x	P(X <= x)
17.8	0.0981

Inverse Normal Problems

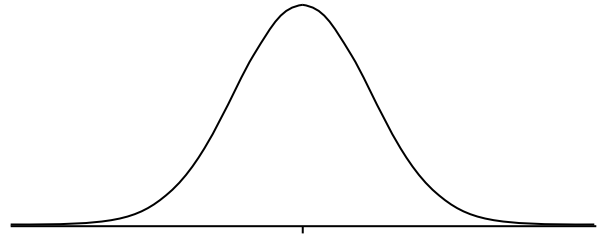
IQ Scores

The distribution of IQ scores is modelled by a Normal distribution with a mean of 100 and a standard deviation of 15.

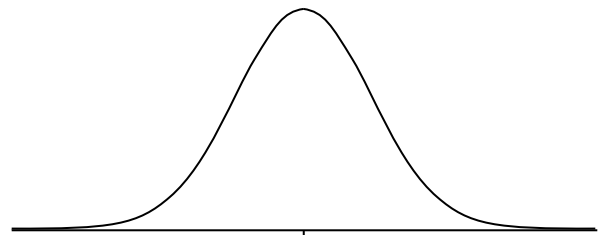
Let X be the IQ of a person.

So $X \sim \text{Normal}(\mu = 100, \sigma = 15)$

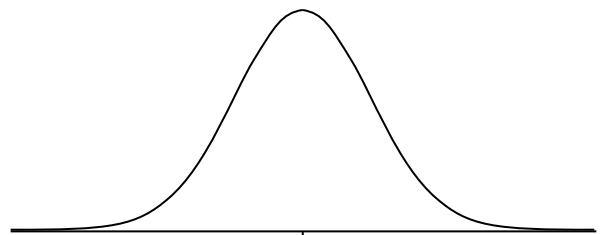
- a. Find the IQ score that the bottom 20% of the population fall below (i.e., the 20th percentile).



- b. What IQ score is exceeded by only the top 10% of the population?



- c. Find the interquartile range for IQ scores.



Computer Output

Normal with mean = 100 and standard deviation = 15

$P(X \leq x)$	x
0.10	80.78
0.20	87.38
0.25	89.88
0.75	110.12
0.80	112.62
0.90	119.22

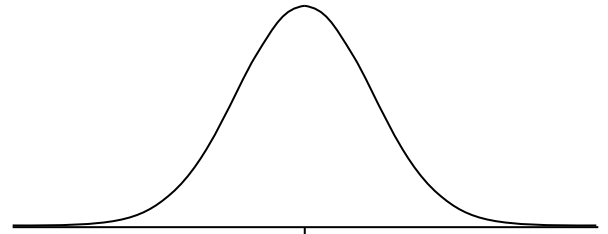
Human Gestation Period

The distribution of the natural human gestation period (in weeks) is modelled by a Normal distribution with a mean of 40 and a standard deviation of 2.3.

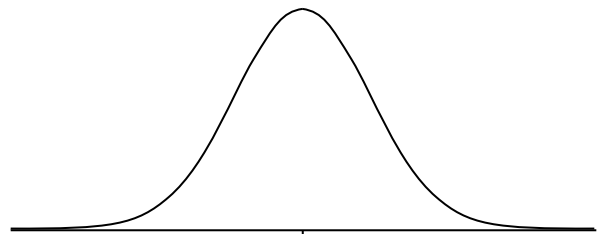
Let X be the natural human gestation period in weeks.

So $X \sim \text{Normal}(\mu = 40, \sigma = 2.3)$

a. Before how many weeks do 10% of births occur?

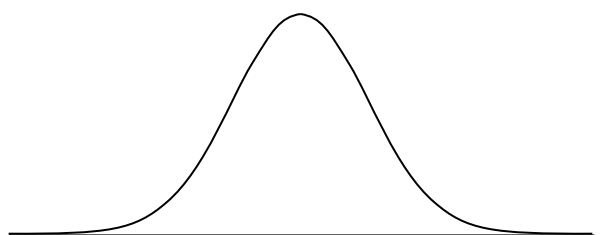


b. What is the range of gestation period for the central 50% of births?



Range is to weeks

c. 95% of births occur after how many weeks?



Computer Output

Normal with mean = 40 and standard deviation = 2.3

$P(X \leq x)$	x
0.05	36.22
0.10	37.05
0.25	38.45
0.75	41.55
0.90	42.95
0.95	43.78

Working in Standard Units

Test and Examination Results

In a recent semester a student obtained the following marks for the Introductory Statistics mid-semester test and exam.

Test: 19/25

Exam: 40/50

In which of these assessments did this student achieve better results in terms of ranking with students in the same course?

The distribution of test marks is modelled by a Normal distribution with mean 13.5 and standard deviation 4.6 and the distribution of exam marks is modelled by a Normal distribution with mean 30.1 and standard deviation 9.4.

Let T be the test mark

E be the exam mark

So $T \sim \text{Normal}(\mu_T = 13.5, \sigma_T = 4.6)$

$E \sim \text{Normal}(\mu_E = 30.1, \sigma_E = 9.4)$

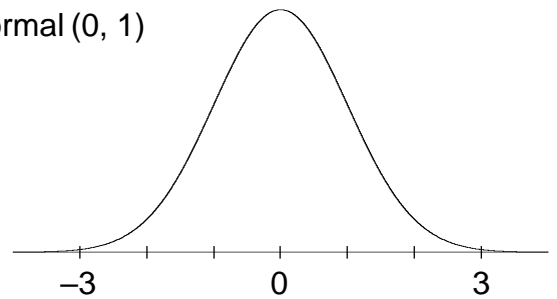
Test: z-score for 19 =

Exam: z-score for 40 =

This student did better in the

in terms of ranking with students in the same course.

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



Percentage Body Fat for Competitive Swimmers

The distribution of percentage body fat for competitive swimmers is modelled by a Normal distribution with a mean of 10 and a standard deviation of 4.

Let X be the percentage body fat of a competitive swimmer.

So $X \sim \text{Normal}(\mu = 10, \sigma = 4)$

Find the z-score for these percentage body fat values for competitive swimmers:

a. $x = 18$ $z =$ (i.e., 18 is sd the mean of 10)

b. $x = 6$ $z =$ (i.e., 6 is sd the mean of 10)

c. $x = 12.6$ $z =$ (i.e., 12.6 is sd the mean of 10)

The z-score for an observation x :

$z =$

Find the z-score for these percentage body fat values for competitive swimmers:

d. $x = 15.7$ $z =$

e. $x = 8.9$ $z =$

Calculate the percentage body fat values (i.e., x -scores) for competitive swimmers with the following z-scores:

f. $z = 2.3$

x is 2.3 the mean so

$x =$

g. $z = -1.8$

x is 1.8 the mean so

$x =$

Continuous Random Variables

If a random variable, X , can take any value in some interval it is called a **continuous random variable**.

There are no gaps between the values a continuous random variable can take.

Continuous random variables measure characteristics of items chosen from a:

- clearly defined finite population, or
- random process producing observations

Examples: time, weight, concentration of alpha fetoprotein, annual share return for a company.

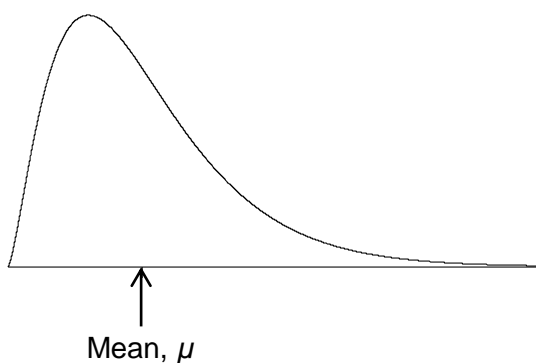
Properties of a Probability Density Function (p.d.f.)

1. The p.d.f. curve is always above or on the x-axis.
2. Probabilities are represented by areas under the p.d.f. curve.
 $\text{pr}(a \leq X \leq b) = \text{area under the p.d.f. curve between } x = a \text{ and } x = b.$
3. The total area under a p.d.f. curve = 1

Mean (Expected Value) and Standard Deviation

The **expected value** of a random variable X , $E(X)$, is:

- the long-run average for X -values
- also called the **population mean**, μ_x (often shortened to μ)
- where the probability density curve balances



The **standard deviation** of a random variable X , $\text{sd}(X)$, is:

- a measure of the variability (spread) of X -values
- roughly, the average distance of the X -values from the population mean
- also called the **population standard deviation**, σ_x (often shortened to σ)

The Normal Distribution

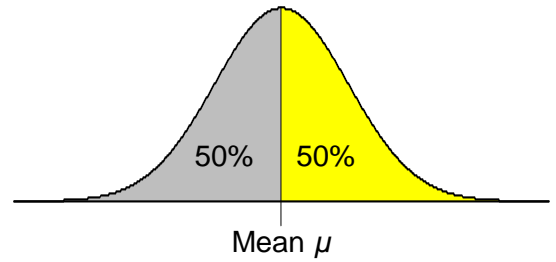
If the distribution of a random variable, X , has a Normal distribution with mean μ and standard deviation σ we write:

$$X \sim \text{Normal}(\mu, \sigma)$$

μ and σ are called the **parameters** of the distribution.

Features of the Normal density curve:

- Symmetric and bell-shaped
- Centred at μ
- σ determines the spread (and hence the height)



68–95–99.7 Rule

In a Normal distribution with mean μ and standard deviation σ , approximately:

68% of observations are within **1** standard deviation of the mean,
i.e., between $\mu - 1\sigma$ and $\mu + 1\sigma$ (or $\mu \pm 1\sigma$)

95% of observations are within **2** standard deviations of the mean,
i.e., between $\mu \pm 2\sigma$

99.7% of observations are within **3** standard deviations of the mean,
i.e., between $\mu \pm 3\sigma$

Obtaining Normal Probabilities

Use statistical software, e.g. Excel, SPSS.

When obtaining Normal probabilities we give the software an **x-value** and it returns a **probability**.

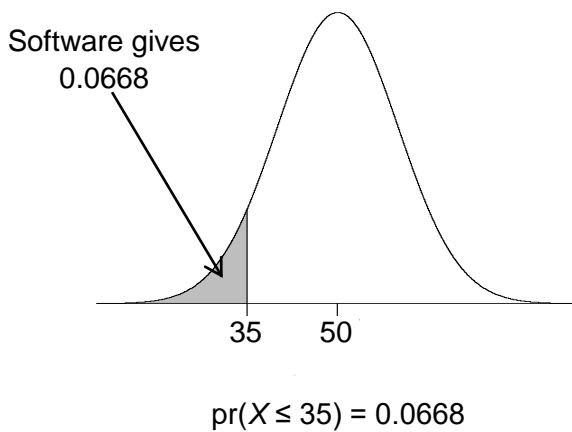
Most software calculates the value of $\text{pr}(X \leq x)$, the cumulative or lower-tail probability.

Method for obtaining Normal probabilities:

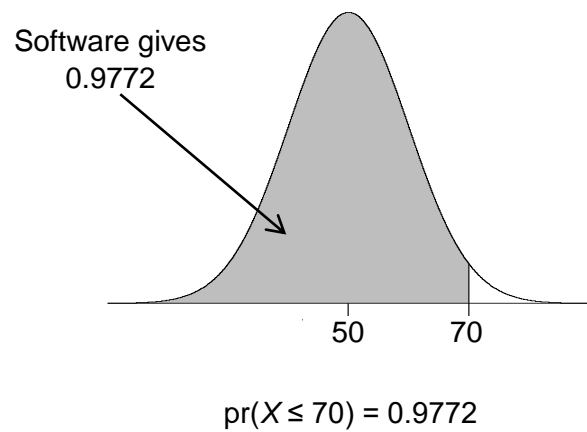
1. Sketch a Normal curve, marking on the mean and value(s) of interest.
2. Shade the area under the curve corresponding to the required probabilities.
3. Obtain the desired probabilities from the lower-tail probabilities provided by the software.

$$X \sim \text{Normal}(\mu = 50, \sigma = 10)$$

Find $\text{pr}(X \leq 35)$



Find $\text{pr}(X \leq 70)$



Inverse Normal Problems

Also use statistical software, e.g. Excel, SPSS.

When solving inverse Normal problems we give the software a **probability** and it returns an **x-value**.

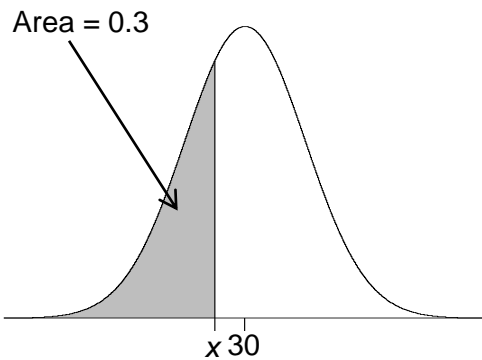
Most software requires the value of $\text{pr}(X \leq x)$, the cumulative or lower-tail probability, to be given.

Method for solving inverse Normal problems:

1. Sketch a Normal curve, marking on the mean.
2. Shade the area under the curve corresponding to the given probability.
3. Obtain the desired x -value from the lower-tail probability given to the software.

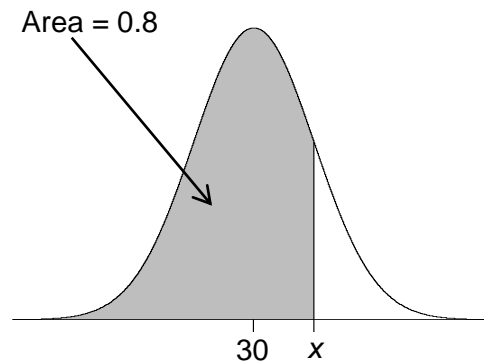
$$X \sim \text{Normal}(\mu = 30, \sigma = 5)$$

Find x , where $\text{pr}(X \leq x) = 0.3$



Software gives
 $x = 27.38$

Find x , where $\text{pr}(X \leq x) = 0.8$



Software gives
 $x = 34.21$

Standard Units (z-scores)

The **z-score** for an observation, x , is the number of standard deviations x is from the mean, μ .

If x is an observation from a Normal distribution with mean μ and standard deviation σ then the **z-score** is

$$z = \frac{x - \mu}{\sigma}$$

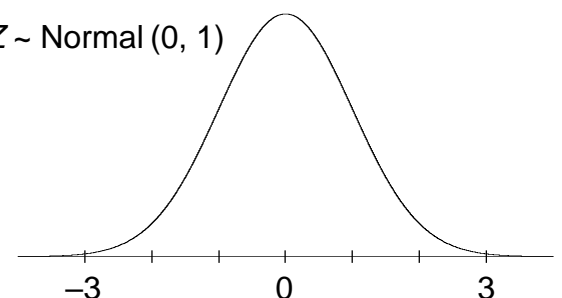
- If x is above the mean then the z-score is positive.
- If x is below the mean then the z-score is negative.

Standard Normal Distribution

If $X \sim \text{Normal}(\mu, \sigma)$ and $Z = \frac{X - \mu}{\sigma}$ then $Z \sim \text{Normal}(0, 1)$.

$Z \sim \text{Normal}(0, 1)$ is called the **standard Normal distribution**.

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



Sample Exam / Cecil Test Questions

Questions 1 to 8 refer to the following information.

For mothers carrying foetuses with spina bifida, the distribution of the concentration of alpha fetoprotein ($\mu\text{M/L}$) in a mother's urine is modelled by a Normal distribution with mean 23.1 and standard deviation 4.1.

Use the following computer output to answer Questions 1 to 6.

Normal with mean = 23.1 and standard deviation = 4.1

x	P(X ≤ x)	P(X ≤ x)	x
17	0.0684	0.0500	16.356
18	0.1068	0.1000	17.846
19	0.1587	0.2000	19.649
20	0.2248	0.2248	20.000
21	0.3043	0.2400	20.204
22	0.3942	0.2500	20.335
23	0.4903	0.2600	20.462
31	0.9730	0.6785	25.000
32	0.9850	0.7500	25.865
33	0.9921	0.8000	26.551
		0.9000	28.354
		0.9500	29.844

- For a randomly selected mother carrying a foetus with spina bifida, the probability that the concentration of alpha fetoprotein in her urine is less than $20.0 \mu\text{M/L}$ is approximately:
 - 0.1587
 - 0.6957
 - 0.7752
 - 0.2248
 - 0.3043
- For mothers carrying foetuses with spina bifida, the proportion who have an alpha fetoprotein concentration greater than $18.0 \mu\text{M/L}$ in their urine is approximately:
 - 0.8413
 - 0.1068
 - 0.8932
 - 0.1587
 - 0.9316

3. For mothers carrying foetuses with spina bifida, the proportion who have an alpha fetoprotein concentration between $22.0 \mu\text{M/L}$ and $32.0 \mu\text{M/L}$ in their urine is approximately:
- a. 0.4827
 - b. 0.4947
 - c. 0.6878
 - d. 0.5788
 - e. 0.5908
4. For mothers carrying foetuses with spina bifida, the concentration of alpha fetoprotein in urine for which 25% lie below is approximately:
- a. $20.335 \mu\text{M/L}$
 - b. $25.865 \mu\text{M/L}$
 - c. $20.204 \mu\text{M/L}$
 - d. $0.679 \mu\text{M/L}$
 - e. $20.462 \mu\text{M/L}$
5. For mothers carrying foetuses with spina bifida, the concentration of alpha fetoprotein in urine for which 20% lie above is approximately:
- a. $19.649 \mu\text{M/L}$
 - b. $0.775 \mu\text{M/L}$
 - c. $28.354 \mu\text{M/L}$
 - d. $26.551 \mu\text{M/L}$
 - e. $0.225 \mu\text{M/L}$
6. For mothers carrying foetuses with spina bifida, the range of the central 90% of concentrations of alpha fetoprotein in urine is approximately between:
- a. $19.6 \mu\text{M/L}$ and $26.6 \mu\text{M/L}$
 - b. $16.4 \mu\text{M/L}$ and $29.8 \mu\text{M/L}$
 - c. $17.8 \mu\text{M/L}$ and $29.8 \mu\text{M/L}$
 - d. $16.4 \mu\text{M/L}$ and $28.4 \mu\text{M/L}$
 - e. $17.8 \mu\text{M/L}$ and $28.4 \mu\text{M/L}$

7. Which **one** of the following statements about mothers carrying foetuses with spina bifida is **false**?
- a. Approximately 68% have alpha fetoprotein concentrations between 19.0 $\mu\text{M/L}$ and 27.2 $\mu\text{M/L}$.
 - b. Approximately 95% have alpha fetoprotein concentrations between 14.9 $\mu\text{M/L}$ and 31.3 $\mu\text{M/L}$.
 - c. Approximately 16% have alpha fetoprotein concentrations greater than 27.2 $\mu\text{M/L}$.
 - d. More than 5% have alpha fetoprotein concentrations less than 14.9 $\mu\text{M/L}$.
 - e. Almost all have alpha fetoprotein concentrations between 10.8 $\mu\text{M/L}$ and 35.4 $\mu\text{M/L}$.
8. A mother who is carrying a foetus with spina bifida has an alpha fetoprotein concentration of 16.8 $\mu\text{M/L}$. What is the z-score for this observed concentration?
- a. 1.54
 - b. 11.17
 - c. -0.55
 - d. -1.54
 - e. 0.55

Questions 9 and 10 refer to the following information.

Hypholoma Capnoides is a pleasant tasting mushroom which looks very much like the generally taller, poisonous fungus *Sulphur Tuft*. Let H be the height (in centimetres) of a *Hypholoma Capnoides* mushroom and let S be the height (in centimetres) of a *Sulphur Tuft* fungus. The distribution of H is modelled by a Normal distribution with mean 6.5 and standard deviation 1.76 and the distribution of S is modelled by a Normal distribution with mean 9.5 and standard deviation 1.25.

In summary, $H \sim \text{Normal}(\mu_H = 6.5, \sigma_H = 1.76)$ and $S \sim \text{Normal}(\mu_S = 9.5, \sigma_S = 1.25)$.

9. Which **one** of the following statements is **false**?

The proportion of *Hypholoma Capnoides* mushrooms that are:

- a. taller than 11.0 cm is less than the proportion of *Sulphur Tuft* fungi that are taller than 12.5 cm.
- b. taller than 5.0 cm is less than the proportion of *Sulphur Tuft* fungi that are taller than 8.0 cm.
- c. taller than 9.5 cm is greater than the proportion of *Sulphur Tuft* fungi that are shorter than 6.5 cm.
- d. shorter than 8.0 cm is less than the proportion of *Sulphur Tuft* fungi that are taller than 8.5 cm.
- e. taller than 4.0 cm is greater than the proportion of *Sulphur Tuft* fungi that are shorter than 11.0 cm.

10. Which **one** of the following statements is **false**?

- a. About 2.5% of *Sulphur Tuft* fungi are taller than 12.0 cm.
- b. The proportion of *Hypholoma Capnoides* mushrooms that are taller than 10.5 cm is greater than the proportion of *Sulphur Tuft* fungi that are taller than 12.5 cm.
- c. About 16% of *Hypholoma Capnoides* mushrooms are shorter than 4.74 cm.
- d. *Hypholoma Capnoides* mushrooms are generally shorter than and more variable in height than *Sulphur Tuft* fungi.
- e. The proportion of *Hypholoma Capnoides* mushrooms that are shorter than 2.5 cm is less than the proportion of *Sulphur Tuft* fungi that are shorter than 6.0 cm.

Answers: (See Section D: Lecture and Tutorial Answers or the fill-ins on CD-ROM or Cecil)

Tutorial

1. The probability distribution function of a continuous random variable is represented by a *density curve*. The following quiz is about the density curve.
 - (a) How are probabilities represented?
 - (b) What is the total area under the density curve?
 - (c) When we calculate probabilities for a continuous random variable, does it matter whether interval endpoints are included or excluded?
 - (d) What are the parameters of the Normal distribution?
2. The natural gestation period for human births, X , has a mean of about 266 days and a standard deviation of about 16 days. Assume that X is Normally distributed with a mean of

Cumulative Distribution Function

Normal with mean = 266.000 and standard deviation = 16.0000

x	P(X ≤ x)	x	P(X ≤ x)
244.0000	0.0846	279.0000	0.7917
245.0000	0.0947	280.0000	0.8092
246.0000	0.1056	281.0000	0.8257
254.0000	0.2266	286.0000	0.8944
255.0000	0.2459	287.0000	0.9053
256.0000	0.2660	288.0000	0.9154

266 days and a standard deviation of 16 days.

Use the computer output above to answer the following questions.

Calculate the proportion of women who carry their babies for:

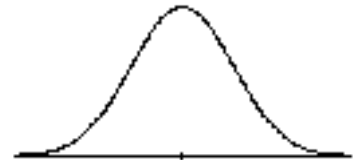
- (a) less than 245 days (i.e., deliver at least 3 weeks early).



(b) between 255 and 280 days.



(c) longer than 287 days (i.e., the baby is more than 3 weeks overdue).



3. A medical trial was conducted to investigate whether a new drug extended the life of a patient who had lung cancer. Assume that the survival time (in months) for patients on this drug is Normally distributed with a mean of 31.1 months and a standard deviation of 16.0 months.

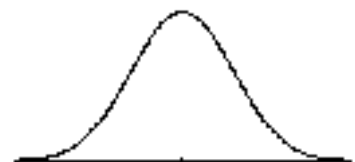
Use the following computer output to answer the questions below.

Inverse Cumulative Distribution Function

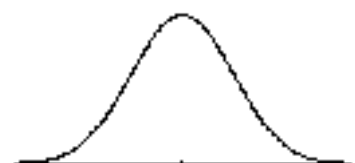
Normal with mean = 31.1000 and standard deviation = 16.0000

P(X <= x)	x
0.1000	10.5952
0.2000	17.6341
0.4000	27.0464
0.6000	35.1536
0.8000	44.5659
0.9000	51.6048

(a) Calculate the number of months beyond which 80% of the patients survive.



(b) Calculate the range of the central 80% of survival times.



4. Complete the following statements about Normal distributions.

(a) z-scores measure the number of _____ an observation is away from the _____.

(b) z-scores are Normally distributed with mean _____ and standard deviation _____.

(c) For Normal distributions:

(i) Approximately _____% of values lie within 1 standard deviation of the mean.

(ii) Approximately _____% of values lie within 2 standard deviations of the mean.

(iii) Approximately _____% of values lie within 3 standard deviations of the mean.

5. *Automotive News* reported the results of a study investigating the time taken to assemble a car. It was found that the time taken to assemble a Japanese car at a plant situated in Japan, J , was Normally distributed with a mean of 20.3 hours and a standard deviation of 1.5 hours and that the time taken to assemble a Japanese car at a plant situated in America, A , was Normally distributed with a mean of 19.6 hours and a standard deviation of 1.9 hours.

I.e. $J \sim \text{Normal}(\mu = 20.3, \sigma = 1.5)$ and $A \sim \text{Normal}(\mu = 19.6, \sigma = 1.9)$

State whether each of the following statements are true or false?

Statement 1:

The time taken to assemble cars in Japan is generally longer and less variable than for cars assembled in America.

Statement 2:

About 68% of cars assembled in America take between 15.8 and 23.4 hours to assemble.

Statement 3:

About 2.5% of cars assembled in Japan take longer than 23.3 hours to assemble.

Statement 4:

The proportion of assembly times longer than 21.8 hours for cars assembled in Japan is approximately the same as the proportion of assembly times shorter than 17.7 hours for cars assembled in America.

Statement 5:

The proportion of assembly times shorter than 22.4 hours for cars assembled in Japan is greater than the proportion of assembly times longer than 16.4 hours for cars assembled in America.

Answers: (See Section D: Lecture and Tutorial Answers or the CD-ROM or Cecil)