

Objectives

- Understanding that a percentile is a measure of position and that the quartiles and median are also percentiles
- Learning how to describe or compare distributions, with either box plots or percentile plots
- Finding rules for how summary measures are affected when adding or multiplying each data value by a constant

Activity Time: 40–50 minutes

Setting: Paired/Individual Activity or Whole-Class Presentation (use **CensusAtSchool.ftm** for either)

Statistics Prerequisites

- Comparing distributions
- Familiarity with quartiles
- Familiarity with the mean, median, and spread (standard deviation) in a general sense

Statistics Skills

- Comparing distributions using percentile plots and box plots
- Transformations and their effects on summary measures
- Percentiles and using a percentile plot
- Box plots and the five-number summary
- Finding outliers using the $1.5 \cdot IQR$ rule
- Comparing types of plots

AP Course Topic Outline: Part I A, B (2–5), C, E

Fathom Prerequisites: Students should be able to make and split graphs and plot values.

Fathom Skills: Students make box plots and modified box plots, calculate five-number summaries, use sliders to transform attributes, work with percentiles and reading percentile plots, create formulas, split graphs using categorical variables, and use filters.

General Notes: This activity uses the CensusAtSchool data to introduce students to Fathom. It requires little to no knowledge of Fathom or statistics. This activity focuses on two different topics: measures of position, including percentile plots and box plots, and the effects of

transformations on summary measures. Fathom makes it easy for students to work with a large data set and quickly compare distributions. Additionally, students can dynamically change transformation parameters and watch how the distributions change.

Procedure: Try working with the document

CensusAtSchool.ftm before giving it to your students.

This file has 500 cases. If it's too slow on your computers, select about half the cases and delete them. To select many cases, click on the first row number in the case table.

Then scroll halfway down and Shift-click on another row number. All the cases in between will be selected. Choose

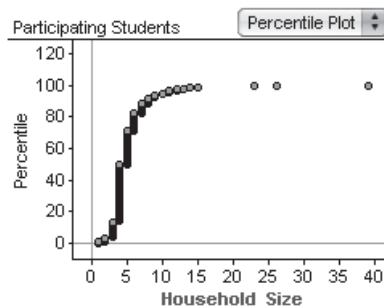
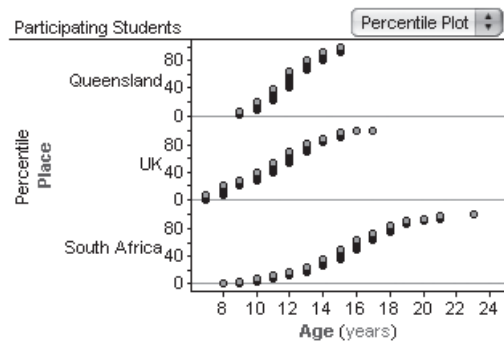
Edit | Delete Cases.

EXAMINE DATA

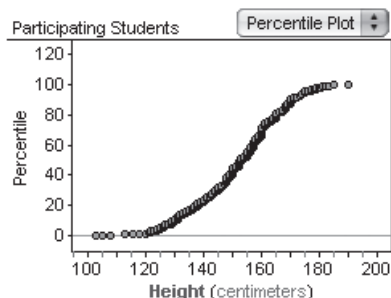
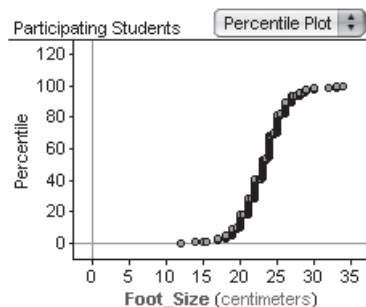
Q1–Q3 The 25th and 75th percentiles are Q_1 and Q_3 respectively. The median is the 50th percentile. The age at the 25th percentile is 10 yr. The age at the 50th percentile is 12 yr. The age at the 75th percentile is 15 yr. The oldest 10% of the population is in the 90th percentile. The age at the 90th percentile is 17 yr, and the age at the 10th percentile is 8 yr. To find these, find points as close as possible to 90 and 10.

Q4 When the percentile plot is steep over an interval, that means there is a concentration of data values in that interval, whereas when it is shallow, that means there are fewer values in that interval. The percentile plot is steep at first and then gradually becomes shallower, so the age distribution is skewed right.

Q5 The split percentile graph shows that the narrowest range of ages is for the Queensland students; the South African students span the greatest range and have all the oldest students; and the youngest students are from the United Kingdom. The percentile plot for South Africa starts out shallow and gets steeper in the middle and then levels out at about 20 years old. The plot for the UK is, from the beginning, fairly steep and looks almost linear, then at the very last two dots it levels out. The plot for Queensland is mostly steeper still but has a small shallow part around 15.



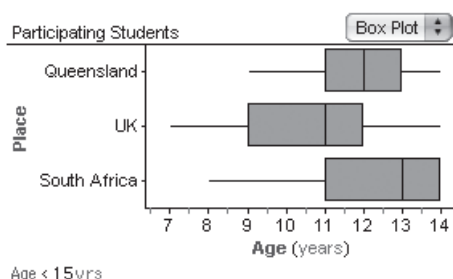
Q6 The percentile plot for *AgeMonths* is naturally similar to the plot for *Age*, but it shows a more detailed distribution. The plot for *Foot_Size* is steep in the center with a shallow part at either end, meaning that most of the data are concentrated between about 19 cm and 25 cm. The plot for *Height* doesn't have such a steep middle section, indicating that the heights are more spread out. The plots for *Household_Size*, *SchoolAgedMales*, and *SchoolAgedFemales* are all quite steep at the beginning, with long shallow ends, meaning that household sizes tend to be small. Similarly, the plot for *Time_to_Travel* has a steep early part and a long shallow end.



INVESTIGATE

- Q7** Median = 12, $Q_1 = 10$, and $Q_3 = 15$. These values should be similar to, if not exactly the same as, the estimates in Q_1 and Q_2 .
- Q8** The outlier is at 23 years. The *IQR* is 5, so $1.5 \cdot IQR = 7.5$. Then $Q_3 + 1.5 \cdot IQR = 22.5$, so 23 is definitely an outlier.
- Q9** The plot again shows that the narrowest range of ages is for the Queensland students; the South African students span the greatest range and have all the oldest students; and the youngest students are from the United Kingdom. The box plots clearly show, however, that all three are fairly mound-shaped and that the students from South Africa are definitely older—75% of the South African students are 14 or older, whereas for both the UK and Queensland, 75% of the students are younger than 14. The spreads (*IQR*) for both South Africa and the UK are the same (4 yr), while Queensland has a spread (*IQR*) of 2 yr. The box plot tells a great deal more about the shape, center, and spread than both the histogram and the dot plot in this case.
- Q10** The outlier (23) is from South Africa. When all of the ages are lumped together into one distribution, 23 years old is quite a bit older than everyone else; however, compared to only the other South African students, the outlier is not that much older. For the whole group, $Q_3 = 15$, but for South Africa, $Q_3 = 18$. So, $Q_3 + 1.5 \cdot IQR = 18 + 1.5(4) = 24$ instead of 22.5.
- Q11** The distributions of *Age* for both South Africa and the UK are skewed toward the younger values. The middle half of the ages of the younger Queensland

students are between 11 and 13. The spreads for the younger UK and South African students are a bit larger: The middle half of the ages are between 9 and 12 for the UK and between 11 and 14 for South Africa. Half of the younger Queensland students have ages above 12 and half 12 or below. Half of the younger South African students have ages above 13 and half 13 or below, and half of the younger UK students have ages above 11 and half 11 or below. Note that for South Africa there is no upper whisker and 50% of the students are between 13 and 14.



Q12 From a box plot, you can see the five-number summary exactly and outliers are clearly marked. These must be estimated from a histogram, which can be difficult. From a histogram, you can estimate the mean by estimating a balance point for the distribution. You cannot do this with a box plot. A histogram will reveal the frequency of the data within an interval. You do not know the exact values but you know how many are within the given boundaries. You know a lower bound and an upper bound, but not necessarily the exact least and greatest values. You know where there are clusters of data and where there are gaps. With a box plot, you get a sense of the basic shape of the distribution but you cannot see clusters or gaps. You cannot see the frequency but you can see the proportions.

10. This group of students was chosen because there are relatively few cases and it is easier to work with the transformed data than with the whole group.

Q13 The majority of dots are captured between the two values—38 out of 52 are captured.

Q14–Q16 When adding the constant a to each height, the whole distribution shifts by that constant. So the range and spread stay the same because the same

constant is added to all values in the range (including the minimum and maximum values). The mean changes to $original\ mean + a$.

Q17–Q19 When multiplying height by the constant a , all three values change. Students will need to resize their plots to see their data. The new mean will be $original\ mean \cdot a$. The new range and the new spread will be the original values multiplied by $|a|$.

EXPLORE MORE

- When adding the constant a to each height, the whole distribution shifts by that constant. So the *IQR* stays the same but the median changes to $original\ median + a$. When multiplying height by the constant a , both values change. Students will need to resize their plots to see their data. The new median will be $original\ median \cdot a$. The new *IQR* will be $original\ IQR \cdot |a|$.
- Taxis have the largest median age associated with them at 16; cars and cycles tie for the lowest median age at 10; walkers have the largest *IQR* at 5 yr. The largest range of older students take taxis, walk, or ride the bus, while the largest range of younger students walk, ride in a car or the bus, or take an unlisted mode of transportation.
- It appears that the South African students have the longest travel time: 50% have a travel time of at least 15 min. The UK and Queensland students both have median travel times of 10 min. This could be due to the large number of students who walk to school in South Africa. The outliers could represent students that live in very rural areas served by regional schools or students on the last stop on a bus route.

