

Stat645: Week 3

Perceptual properties of continuous variables

Today, we're going to do some explorations into how we perceive different aesthetics, and look at a variation on the histogram that uses some of these ideas.

Perception

Start by loading `pop-selected.csv` in to R from the website. Last week we plotted it with putting time on the x-axis and rank on the y axis. Repeat that plot for reference. Instead of putting rank on the y-axis, we could display the track name and use another aesthetic to display the rank.

List possible aesthetics (other than position) that we could use to display rank, and then create plots using those aesthetics. What do you notice? What aesthetic makes it easiest to see the pattern over time? What makes it hardest? Does combining aesthetics make it better? What makes the original plot so good?

size — both hard!
colour —

shape — would have to be discrete

```
qplot(week, artist.inverted, data = two, size = value, colour = value)
```

combining colour & size doesn't seem to help much
position is easiest by far!

Use `scale_colour_gradient` to vary the colours at the high-end and low-end of the scales. Can you create a scale that makes it easier to see changes over time? Which end of the scale should draw your attention?

```
+ scale_colour_gradient(low = "white", high = "black")
```

Low values most important b/c they reflect more popular songs

You've probably used size as one of the aesthetics - but what is size really mapped to? Create a small experiment to determine whether size is mapped to radius or area? Which do you think it should be mapped to?

```
df <- data.frame(x = 1:5, size = c(1, 2, 4, 8, 16))
```

```
qplot(x, 1, data = df, size = size)
```

size is mapped to radius, not area!

Use `scale_area` to change the mapping so the square root of size is mapped to radius. How does this change the appearance of the plot?

Makes it a little easier to read

Instead of points, we could use bars to display the data. `geom_bar` has a default statistic that counts up observations, so if we want to use bars for this type of data we need to turn that off by specifying `stat = "identity"`:

```
qplot(week, value, data = two, geom = "bar", stat = "identity") +  
  facet_wrap(~ artist.inverted)
```

How does this display compare to the previous displays? What property are we reading for this plot?

Length + position along common scale

What happens if you add on `coord_polar(theta = "y")`? What property do we now perceive?

Area/ Angle - which is harder to perceive

Hanging rootogram

The hanging rootogram is plot developed by Tukey specifically for comparing an empirical distribution to a theoretical distribution, and is designed to answer questions like "does my data come from a normal distribution?" Read about the hanging rootogram on pages 312-315 of <http://www.edwardtufte.com/tufte/tukey>.

What are the important properties of the hanging rootogram and how do they make the desired comparison easier?

- $\sqrt{\text{count}}$ - variance stabilizing transformation
- align top of bars with density - makes comparison much easier - we're really bad at ~~comparing~~ computing distances b/w two lines

Why is overlaying a predicted density on top of a histogram not a good idea? Hint: <http://www.michaelbach.de/ot/size/sinillusion/index.html>

too hard to look at vertical distances
brain tries to compute shortest distance -

We are now going to construct our own hanging rootogram in R, using ggplot2.

Look at the code sample on the website. What do parts 1-4 do?

- 1) Manually bin data & compute ~~counts~~ counts & density
- 2) Draw histogram Manually onto rectangles
- 3) Generate grid of density values - add to plot ~~at~~ as a line (thick & red)
- 4) Compute number predicted by density

Use this data to create your own rootogram and hanging rootogram in R. For an extra challenge, try creating the suspended rootogram and add control limits. (Hint: `geom_hline`)

`ggplot()` +

`geom_rect(aes(data = hist, xmin = x - 0.2, xmax = x + 0.2,`
`ymax = hist nhist, ymin = nhist - n))` +

`geom_line(aes(x, dens), data = norm)`

If you'd like to read more about the hanging rootogram, this article, <http://www.jstor.org/stable/2683341>, is a good start.