

# Expository graphs

Roger D. Peng, Associate Professor of Biostatistics

May 18, 2016

# Why do we use graphs in data analysis?

- ▶ To understand data properties
- ▶ To find patterns in data
- ▶ To suggest modeling strategies
- ▶ To “debug” analyses
- ▶ To communicate results

# Expository graphs

- ▶ To understand data properties
- ▶ To find patterns in data
- ▶ To suggest modeling strategies
- ▶ To “debug” analyses
- ▶ To communicate results

# Characteristics of expository graphs

- ▶ The goal is to communicate information
- ▶ Information density is generally good
- ▶ Color/size are used both for aesthetics and communication
- ▶ Expository figures have understandable axes, titles, and legends

# Housing data

The screenshot shows the U.S. Census Bureau website. The header includes the U.S. Department of Commerce, the United States Census Bureau logo, and navigation links for People, Business, Geography, Data, Research, and Newsroom. A search bar is also present. The main content area is titled "American Community Survey" and "Public Use Microdata Sample (PUMS)". A left sidebar contains a table of contents with links to Data Releases, Data Product Descriptions, Documentation, Geography, Downloadable data via FTP, Summary File, and Public Use Microdata Sample (PUMS). The PUMS section is expanded, showing sub-links for About PUMS, PUMS Data, PUMS Documentation, PUMS on DataFerrett, PUMS FAQs, and Custom Tabulations. The main content area includes a section for "Public Use Microdata Sample (PUMS)" with a description of the data, a "Summary products" section, a "Why Use PUMS?" section, a "What's Available and How Can I Access PUMS?" section, a "Need Help with PUMS?" section, and a "Geographic Areas Available" section.

U.S. Department of Commerce  
United States Census Bureau

Home | About Us | Subjects A to Z | FAQs | Help

People | Business | Geography | Data | Research | Newsroom

Search

Census.gov > American Community Survey > Data & Documentation: Public Use Microdata Sample (PUMS)

## American Community Survey

Main | About the Survey | Guidance for Data Users | Data & Documentation | Methodology | Library

- Data Releases
- Data Product Descriptions
- Documentation
- Geography
- Downloadable data via FTP
- Summary File
- Public Use Microdata Sample (PUMS)**
  - About PUMS
  - PUMS Data
  - PUMS Documentation
  - PUMS on DataFerrett
  - PUMS FAQs
- Custom Tabulations

### Public Use Microdata Sample (PUMS)

Print | Share this page | Connect with us

The American Community Survey (ACS) Public Use Microdata Sample (PUMS) files are a set of untabulated records about individual people or housing units. The Census Bureau produces the PUMS files so that data users can create custom tables that are not available through pretabulated (or summary) ACS data products.

**Summary products**, such as the tables and profiles accessible via American FactFinder (AFF), show data that have already been tabulated for specific geographic areas.

**PUMS files**, in contrast, include population and housing unit records with individual response information such as relationship, sex, educational attainment, and employment status.

#### Why Use PUMS?

PUMS files are perfect for people, such as students, who are looking for greater accessibility to inexpensive data for research projects. Social scientists often use the PUMS for regression analysis and modeling applications.

#### What's Available and How Can I Access PUMS?

The Census Bureau produces 1-year, 3-year, and 5-year ACS PUMS files. The 3-year and 5-year PUMS files are multiyear combinations of the 1-year PUMS file with appropriate adjustments to the weights and inflation adjustment factors. The PUMS files are accessible via [American FactFinder](#), the Census Bureau's [FTP site](#), and [DataFerrett](#). Statistical software is needed to use the PUMS files from American FactFinder and the FTP site.

#### Need Help with PUMS?

Learn more about PUMS in the Compass Products [What PUMS Data Users Need to Know](#) handbook and [Introduction to the PUMS](#) training presentation.

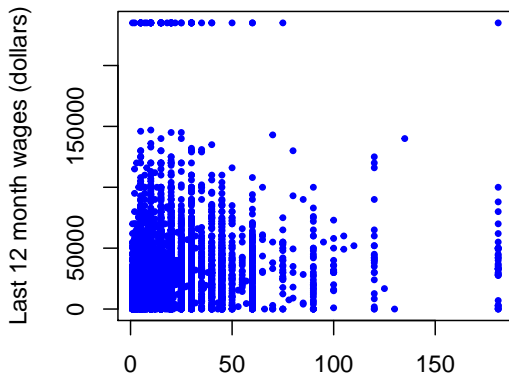
#### Geographic Areas Available

```
pData <- read.csv("./data/ss06pid.csv")
```

## Axes

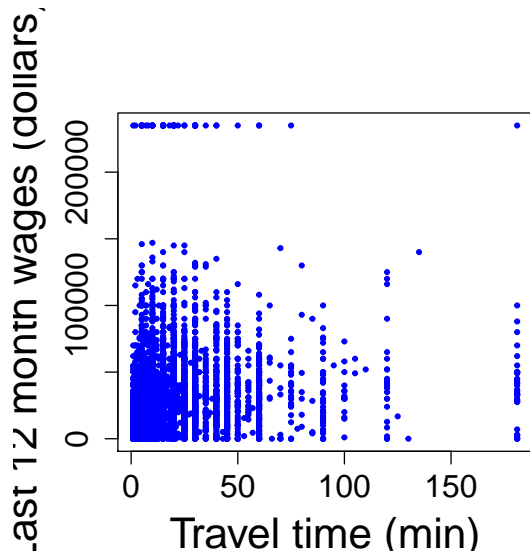
Important parameters: *xlab*, *ylab*, *cex.lab*, *cex.axis*

```
plot(pData$JWMNP, pData$WAGP, pch=19, col="blue", cex=0.5,  
     xlab="Travel time (min)", ylab="Last 12 month wages (dollars)",
```



## Axes

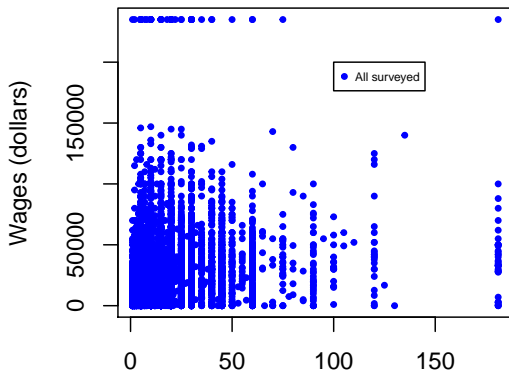
```
plot(pData$JWMNP, pData$WAGP, pch=19, col="blue", cex=0.5,  
     xlab="Travel time (min)", ylab="Last 12 month wages (dollars)",
```



# Legends

- Important parameters: *x,y,legend, other plotting parameters*

```
plot(pData$JWMNP,pData$WAGP,pch=19,col="blue",cex=0.5,xlab=  
legend(100,200000,legend="All surveyed",col="blue",pch=19,c
```

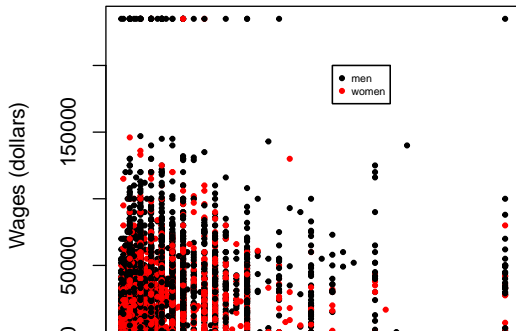




# Legends

```
plot(pData$JWMNP, pData$WAGP, pch=19, cex=0.5, xlab="TT (min)",  
legend(100, 200000, legend=c("men", "women"), col=c("black", "red"))
```

```
## Warning in if (xc < 0) text.width <- -text.width: the co  
## > 1 and only the first element will be used
```

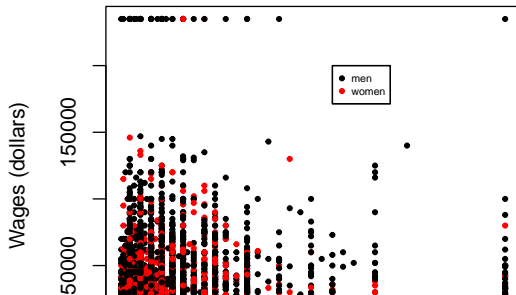


# Titles

```
plot(pData$JWMNP,pData$WAGP,pch=19,cex=0.5,xlab="CT (min)",  
     ylab="Wages (dollars)",col=pData$SEX,main="Wages earned  
legend(100,200000,legend=c("men","women"),col=c("black","red"))
```

```
## Warning in if (xc < 0) text.width <- -text.width: the co  
## > 1 and only the first element will be used
```

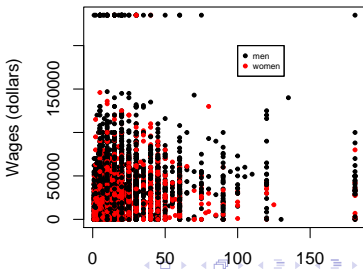
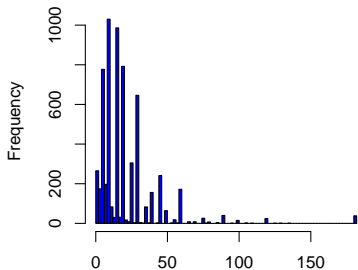
## Wages earned versus commute time



# Multiple panels

```
par(mfrow=c(1,2))  
hist(pData$JWMNP,xlab="CT (min)",col="blue",breaks=100,main="Histogram of CT (min)")  
plot(pData$JWMNP,pData$WAGP,pch=19,cex=0.5,xlab="CT (min)",ylab="Wages (dollars)",  
legend(100,200000,legend=c("men", "women"),col=c("black", "red")))
```

```
## Warning in if (xc < 0) text.width <- -text.width: the coefficient  
## > 1 and only the first element will be used
```

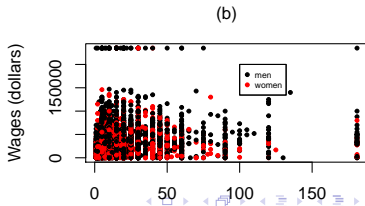
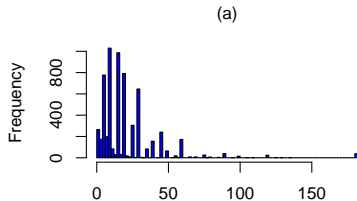


# Adding text

```
par(mfrow=c(1,2))  
hist(pData$JWMNP,xlab="CT (min)",col="blue",breaks=100,main="a")  
mtext(text="(a)",side=3,line=1)  
plot(pData$JWMNP,pData$WAGP,pch=19,cex=0.5,xlab="CT (min)",  
legend(100,200000,legend=c("men","women"),col=c("black","red"))
```

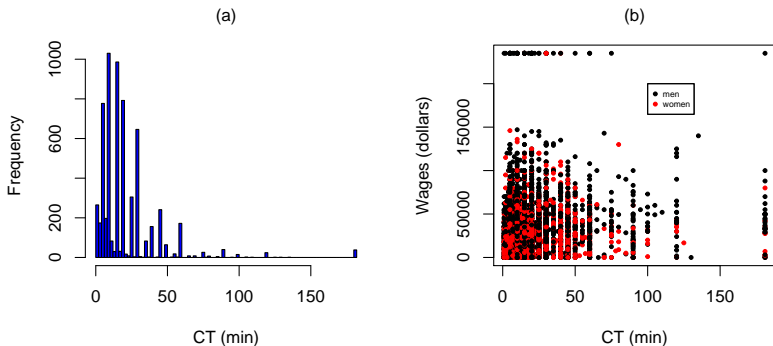
```
## Warning in if (xc < 0) text.width <- -text.width: the coefficient  
## > 1 and only the first element will be used
```

```
mtext(text="(b)",side=3,line=1)
```



## Figure captions

```
## Warning in if (xc < 0) text.width <- -text.width: the co  
## > 1 and only the first element will be used
```



**Figure 1. Distribution of commute time and relationship to wage earned by sex** (a) Commute times in the American Community Survey (ACS) are right skewed. (b) Commute times do

# Colorblindness

Vischeck: VischeckImage x

Go to home page. www.vischeck.com/vischeck/vischeckImage.php

## Vischeck

Home  
Vischeck  
•Run Images  
•Run Webpages  
Daltonize  
Examples  
Downloads  
Info & Links  
FAQ  
About Us

**Try Vischeck on Your Image Files**

Select the type of color vision to simulate:

☒ Deuteranope (a form of red/green color deficit)  
☐ Protanope (another form of red/green color deficit)  
☐ Tritanope (a blue/yellow deficit- very rare)

Image file:  unnamed-chunk-6.png


Notes:

- Vischeck accepts most common image formats. However, we recommend that you use PNG or JPEG format for uploading large images as these tend to transfer faster.
- For PowerPoint slides, you can save all your slides as PNG images with "Save As..." and run Vischeck on each slide.
- If you have many images to process, consider [downloading](#) Vischeck to run on your own computer.)
- Uploading a large file may take a while - please be patient!

Please read our [terms of use](#) before using Vischeck.

User quotes:  
I'd just like to say thank you very much for your time and effort helping people like me with colorblindness. Whenever I meet new co-workers or fellow students I have to repeatedly describe what I can and cannot distinguish- it gets so annoying. I'd just like to say thank you for your effort and understanding.  
-Brad C.

☐ Web ☒ Vischeck

 SUPPORT WIKIPEDIA

www.vischeck.com

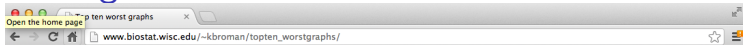
<http://www.vischeck.com/>

# Graphical workflow

- ▶ Start with a rough plot
- ▶ Tweak it to make it expository
- ▶ Save the file
- ▶ Include it in presentations

Saving files in R is done with graphics *devices*. Use the command `?Devices` to see a list. Here we will go over the most popular devices.

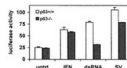
# Something to avoid



## The top ten worst graphs

With apologies to the authors, we provide the following list of the top ten worst graphs in the scientific literature. As these examples indicate, good scientists can make mistakes.

1. Roeder K (1994) DNA fingerprinting: A review of the controversy (with discussion). *Statistical Science* 9:222-278, Figure 4  
[\[The article\]](#) [\[The figure\]](#) [\[Discussion\]](#)
2. Wittke-Thompson JK, Pluzhnikov A, Cox NJ (2005) Rational inferences about departures from Hardy-Weinberg equilibrium. *American Journal of Human Genetics* 76:967-986, Figure 1  
[\[The article\]](#) [\[Fig 1AB\]](#) [\[Fig 1CD\]](#) [\[Discussion\]](#)
3. Epstein MP, Satten GA (2003) Inference on haplotype effects in case-control studies using unphased genotype data. *American Journal of Human Genetics* 73:1316-1329, Figure 1  
[\[The article\]](#) [\[The figure\]](#) [\[Discussion\]](#)
4. Mykland P, Tierney L, Yu B (1995) Regeneration in Markov chain samplers. *Journal of the American Statistical Association* 90:233-241, Figure 1  
[\[The article\]](#) [\[The figure\]](#) [\[Discussion\]](#)
5. Hummer BT, Li XL, Hassel BA (2001) Role for p53 in gene induction by double-stranded RNA. *J Virol* 75:7774-7777, Figure 4  
[\[The article\]](#) [\[The figure\]](#) [\[Discussion\]](#)



http:

//www.biostat.wisc.edu/~kbroman/topten\_worstgraphs/



# Something to aspire to



<http://www.facebook.com/notes/facebook-engineering/visualizing-friendships/469716398919>

## Further resources

- ▶ How to display data badly
- ▶ The visual display of quantitative information
- ▶ Creating more effective graphs
- ▶ R Graphics Cookbook
- ▶ ggplot2: Elegant Graphics for Data Analysis
- ▶ Flowing Data