

Exploratory Data Analysis & Visualization

How good is your data?

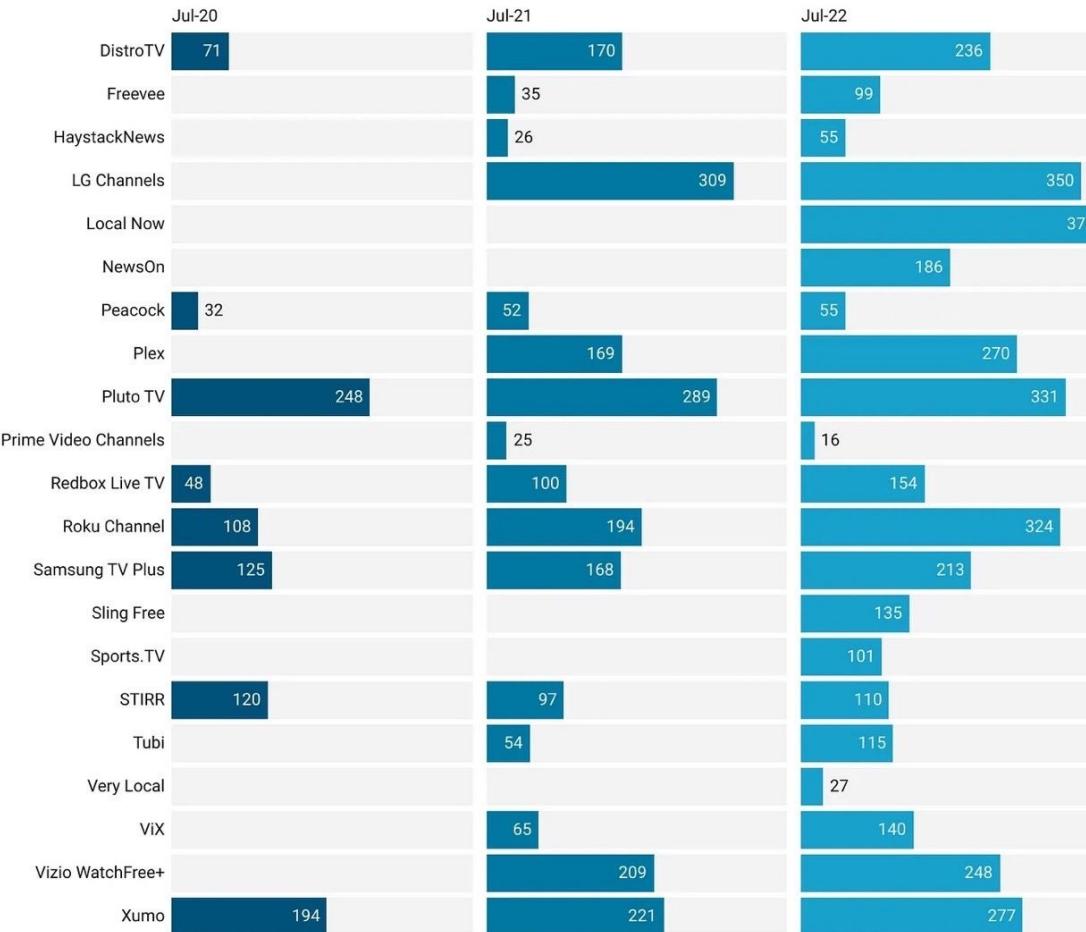
How good is your visualization?

Ben Winjum

Questions about Setups? Materials? Logistics?

Example Visualizations from Discussions

Number of FAST Channels on FAST Services in July 2020, 2021 and 2022

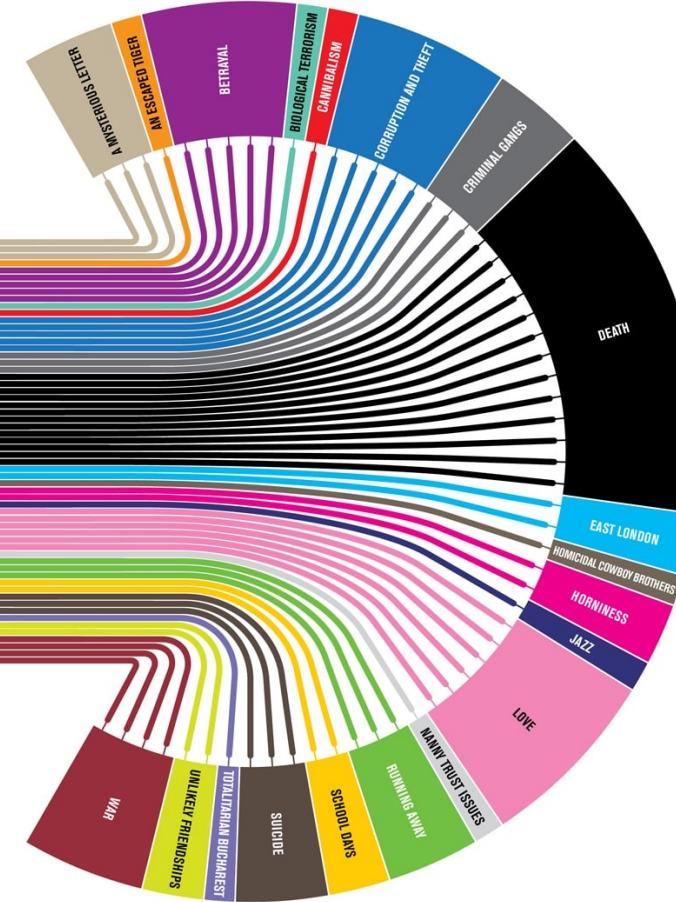
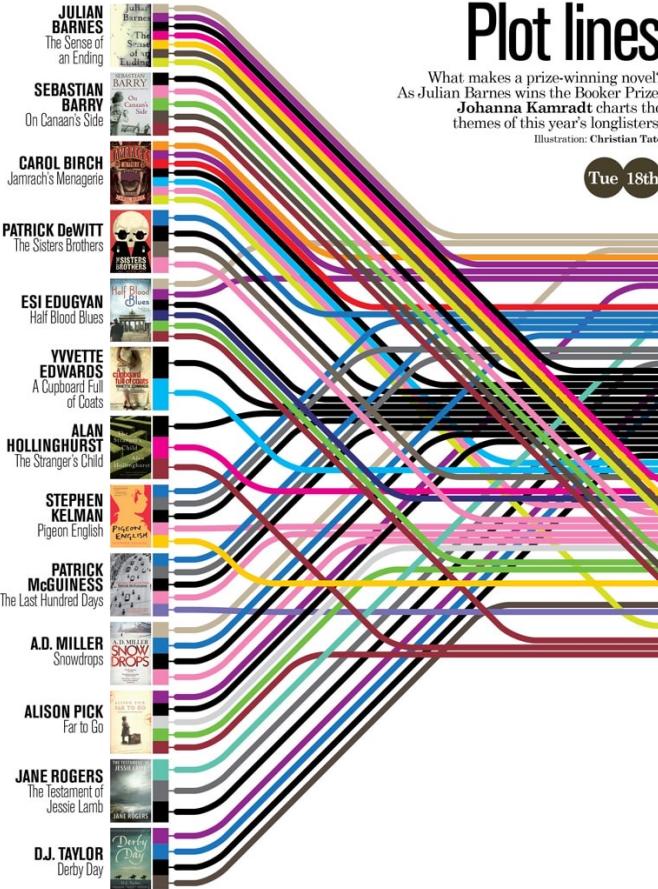


Gerardo

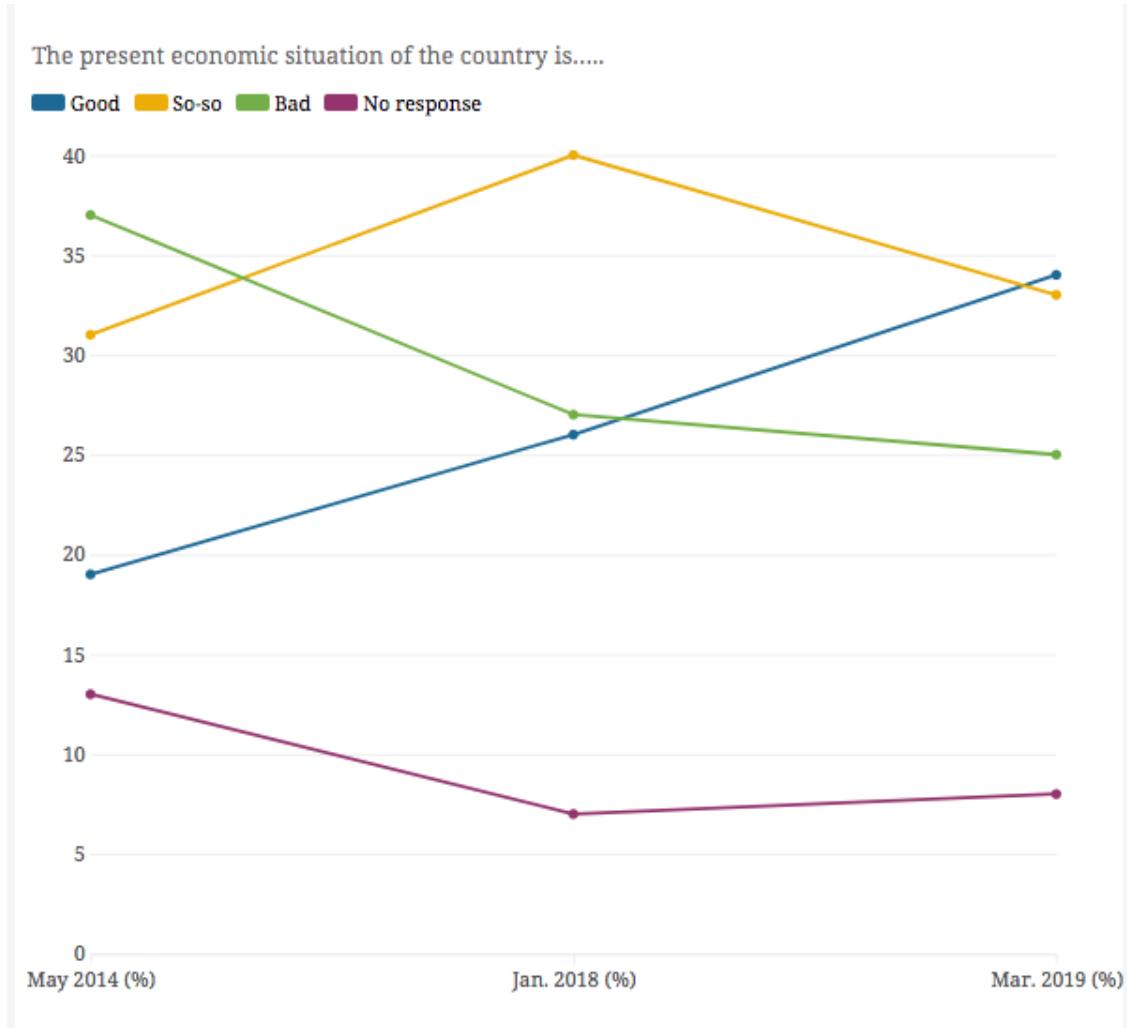
Figures for News by Fire TV are under revision and not included for July

Chart: Gavin Bridge (c/o The FASTMaster) • Source: fastmaster.substack.com • Created with Datawrapper

Example Visualizations from Discussions



Example Visualizations from Discussions

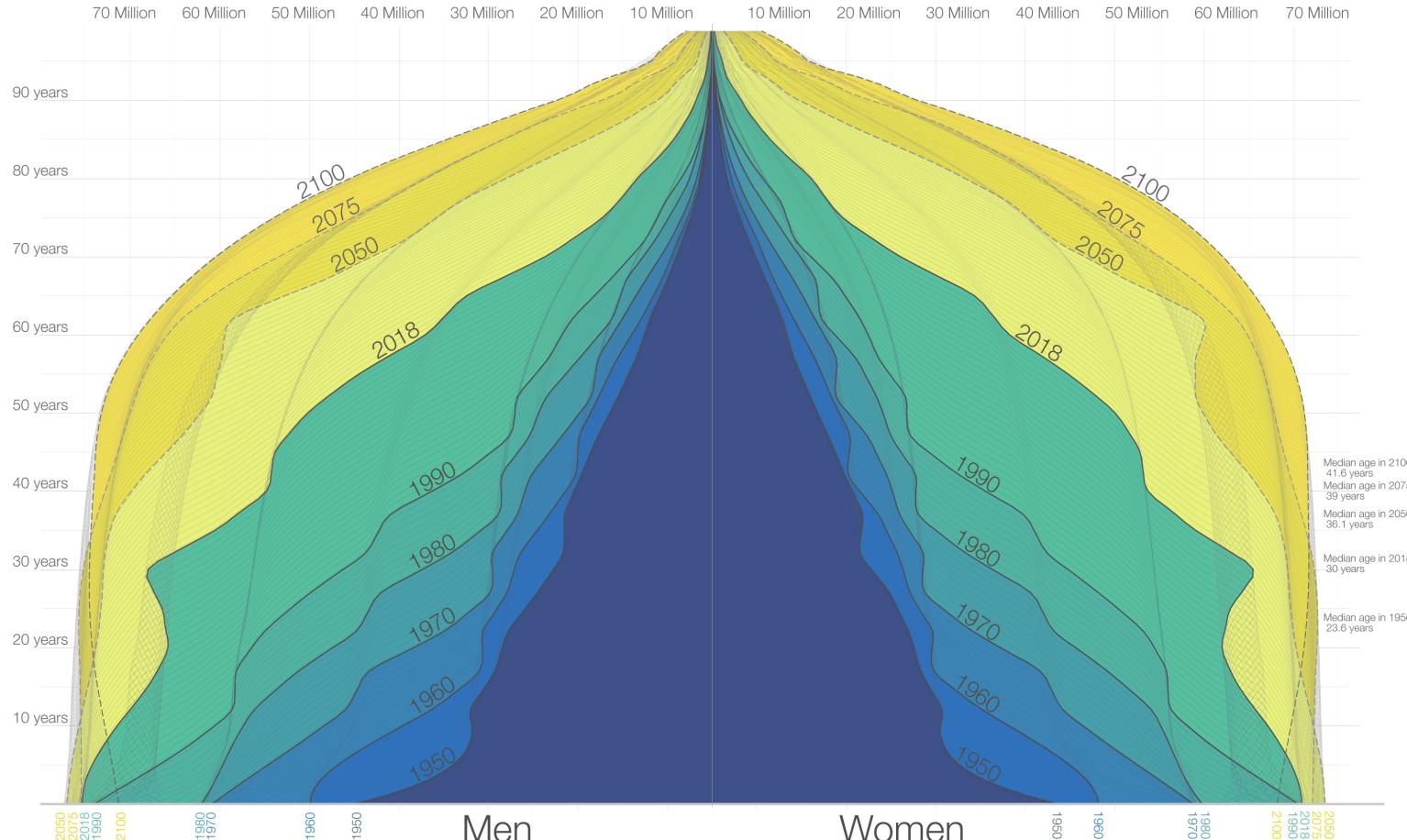


Divya

Example Visualizations from Discussions

The Demography of the World Population from 1950 to 2100

Shown is the age distribution of the world population – by sex – from 1950 to 2018 and the *UN Population Division's* projection until 2100.



Data source: United Nations Population Division – World Population Prospects 2017; Medium Variant.
The data visualization is available at [OurWorldInData.org](https://ourworldindata.org/age-structure), where you find more research on how the world is changing and why.

Licensed under CC-BY by the author Max Roser.

Tristan

Going Through the Assignments: Python

Going Through the Assignments: Python

- First let's review some essential first steps in EDA

Look at your data to see if it passes the common-sense test

- Wrong values? Or metadata values (like number of rows)?
 - Can you detect them? And/or correct them?
 - Does it match with documentation?
 - Do the data have reasonable values?
- Unusable?
 - Does the data allow you to answer your question?
- Missing values?
 - How do you fill them in?
- Messy?
 - Can you use the data?
 - Does it need to be massaged?

Missing data can arise from various places

- A survey was conducted and values were just randomly missed when being entered in the computer.
- A respondent chooses not to respond to a question like 'Have you ever recreationally used opioids?'
- You decide to start collecting a new variable (due to new actions: like a pandemic) partway through the data collection of a study.
- You want to measure the speed of meteors, and some observations are just 'too quick' to be measured properly.

Missing data type 1: Missing Completely at Random

- Examples:
 - a coin is flipped to determine whether an entry is removed.
 - values were just randomly missed when being entered in the computer.
 - a cosmic particle hits your computer and destroys the bit that stores your data element
- Effect if you ignore:
 - there is no effect on inferences (or estimates)
- How to handle:
 - lots of options
 - best to impute (more on that soon)

Missing data type 2: Missing at Random

- Examples:
 - men and women respond to the question "have you ever felt harassed at work?" at different rates (and may be harassed at different rates)
- Effect if you ignore:
 - inferences are biased and predictions are usually worsened
- How to handle:
 - use the information in the other predictors to build a model and impute a value for the missing entry
 - key: we can fix any biases by modeling and imputing the missing values based on what is observed!

Missing data type 3: Missing Not at Random

- Example(s):
 - patients drop out of a study because they experience some really bad side effect that was not measured.
 - cheaters are less likely to respond when asked if you've ever cheated.
- Effect if you ignore:
 - major effects on inferences or predictions.
- How to handle:
 - you can 'improve' things by dealing with it like it is MAR, but you [likely] may never completely fix the bias
 - incorporating a missingness indicator variable may be the best approach (if it is in a predictor)

Correcting for missing data: we will consider basic approaches

- Delete it
 - Drop the feature
 - BUT: Is there so little data that this is reasonable? Or will you be dropping key info?
 - Drop the particular record
 - BUT: Is there a pattern to the data that's missing? And will this cause there to be a resulting bias in your data if you drop the records that have missing data?
 - BUT: If you drop all the records, will what's left be too little to do modeling? And/or make your model brittle?

Missing data

- Delete it
- Fill it with another value (imputation)
 - Mean, median, mode
 - If you have other relevant values, you can use subset for calculations
 - Example, if a male height is missing, substitute with the median value only of other male heights
 - Use regression or classification
 - Fill with values of the “nearest” records
 - Interpolate from “nearest” records

Correcting for missing data

- We are not focused on modeling and will take the basic approach
 - Drop the observations that have any missing values.
 - Use `pd.DataFrame.dropna(axis=0)`
 - Discard the entire feature
 - Use `pd.DataFrame.drop(column, axis=1)`
 - Imputation: substitute a value like the mean or median (if quantitative) or most common class (if categorical) for all missing values
 - Use `pd.DataFrame.fillna(value=x.mean())`

Correcting for missing data

- We are not focused on modeling and will take the basic approach
 - Drop the observations that have any missing values.
 - Use `pd.DataFrame.dropna(axis=0)`
 - Discard the entire feature
 - Use `pd.DataFrame.drop(column, axis=1)`
 - Imputation: substitute a value like the mean or median (if quantitative) or most common class (if categorical) for all missing values
 - Use `pd.DataFrame.fillna(value=x.mean())`
- FYI: there are several alternative approaches to imputation
 - Create a new variable that is an indicator of missingness
 - Hot deck imputation: for each missing entry, randomly select an observed entry in the variable and plug it in.
 - Model the imputation: plug in predicted values from a model based on the other observed predictors.
 - Model the imputation with uncertainty: plug in predicted values plus randomness from a model based on the other observed predictors.

Outliers

- Method of finding
 - Boxplot
 - Scatter plot
 - Histogram – might notice that in the default plot, the scales are way off
- Methods for dealing with them
 - Similar to missing values (except for discarding an entire feature)

Going Through the Assignments: Python

Data

- What is it?
- How do we store it?
- How do we think about it?
- How do we use it?

What is data?

- Webster's Dictionary:
 - factual information (such as measurements or statistics) used as a basis for reasoning, discussion, or calculation
 - information in digital form that can be transmitted or processed
 - information output by a sensing device or organ that includes both useful and irrelevant or redundant information and must be processed to be meaningful

What is data?



KANYE WEST

@kanyewest



Follow

I hate when I'm on a flight and I wake up with
a water bottle next to me like oh great now I
gotta be responsible for this water bottle

RETWEETS

48,234

LIKES

52,776



3:09 PM - 27 Jan 2016

What is data?

Text

KANYE WEST

@kanyewest

Follow

I hate when I'm on a flight and I wake up with a water bottle next to me like oh great now I gotta be responsible for this water bottle

RETWEETS 48,234 LIKES 52,776

3:09 PM - 27 Jan 2016

What is data?

Pictures



KANYE WEST

@kanyewest



Follow

Text

I hate when I'm on a flight and I wake up with
a water bottle next to me like oh great now I
gotta be responsible for this water bottle

RETWEETS

48,234

LIKES

52,776



3:09 PM - 27 Jan 2016

What is data?

Pictures



KANYE WEST

@kanyewest



Follow

Text

I hate when I'm on a flight and I wake up with
a water bottle next to me like oh great now I
gotta be responsible for this water bottle

Numbers

RETWEETS

48,234

LIKES

52,776



3:09 PM - 27 Jan 2016

What is data?

Pictures



KANYE WEST

@kanyewest



Follow

Text

I hate when I'm on a flight and I wake up with
a water bottle next to me like oh great now I
gotta be responsible for this water bottle

Numbers

RETWEETS

48,234

LIKES

52,776



3:09 PM - 27 Jan 2016

Dates

What is data?

Pictures



KANYE WEST

@kanyewest



Follow

Text

I hate when I'm on a flight and I wake up with
a water bottle next to me like oh great now I
gotta be responsible for this water bottle

Numbers

RETWEETS

48,234

LIKES

52,776



3:09 PM - 27 Jan 2016

States

Dates

What is data?

Tweet
(aggregate)

Pictures



KANYE WEST

@kanyewest



Follow

Text

I hate when I'm on a flight and I wake up with
a water bottle next to me like oh great now I
gotta be responsible for this water bottle

Numbers

RETWEETS
48,234

LIKES
52,776

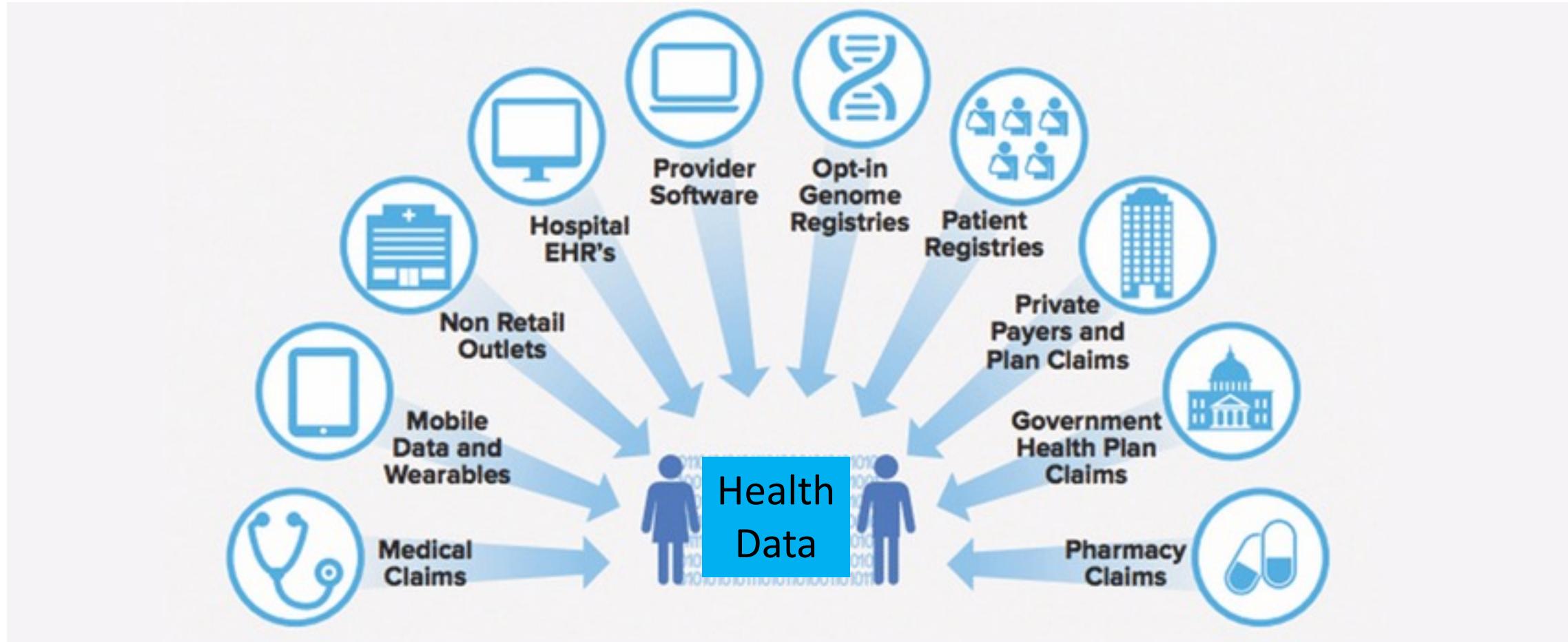


3:09 PM - 27 Jan 2016

States

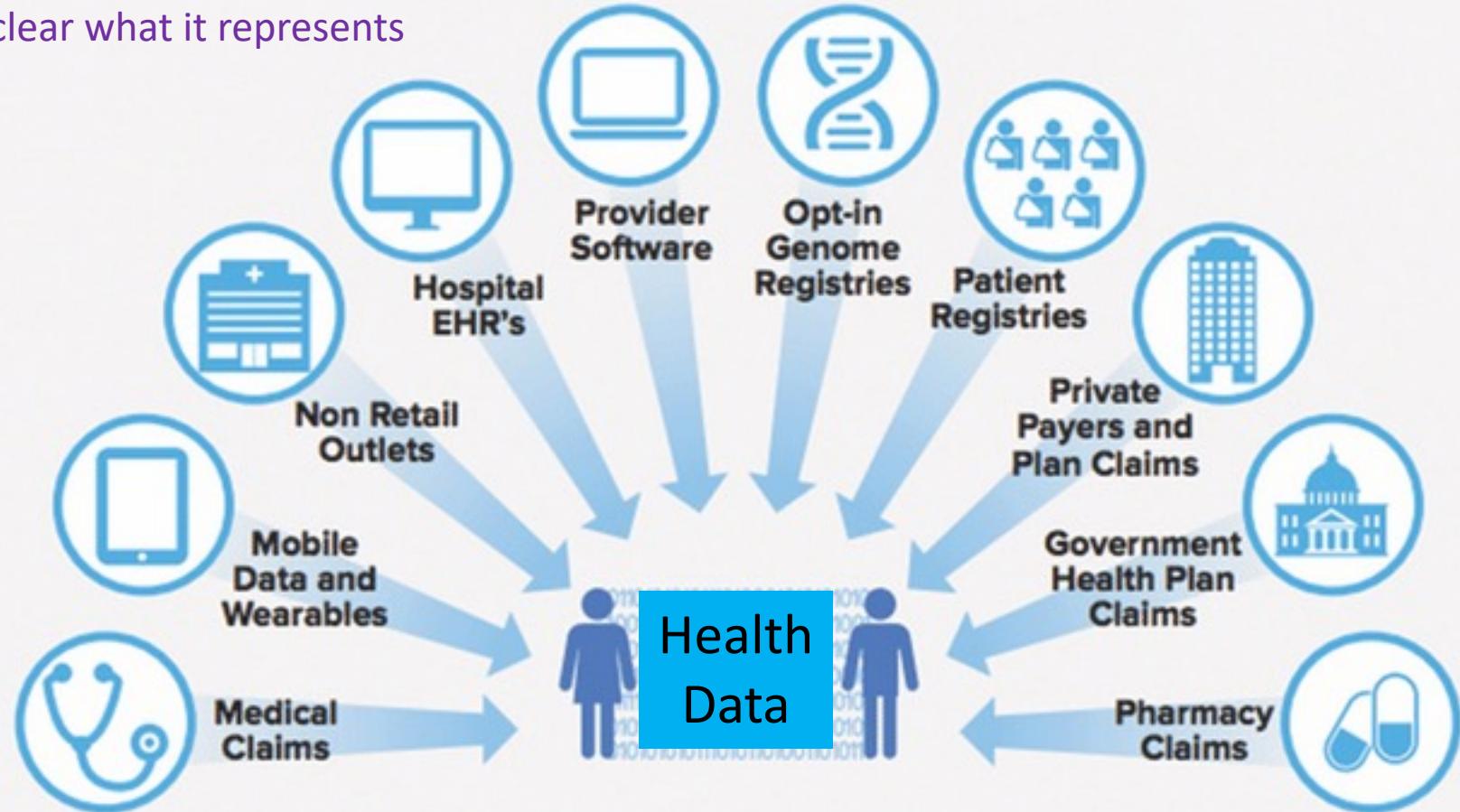
Dates

What is data?



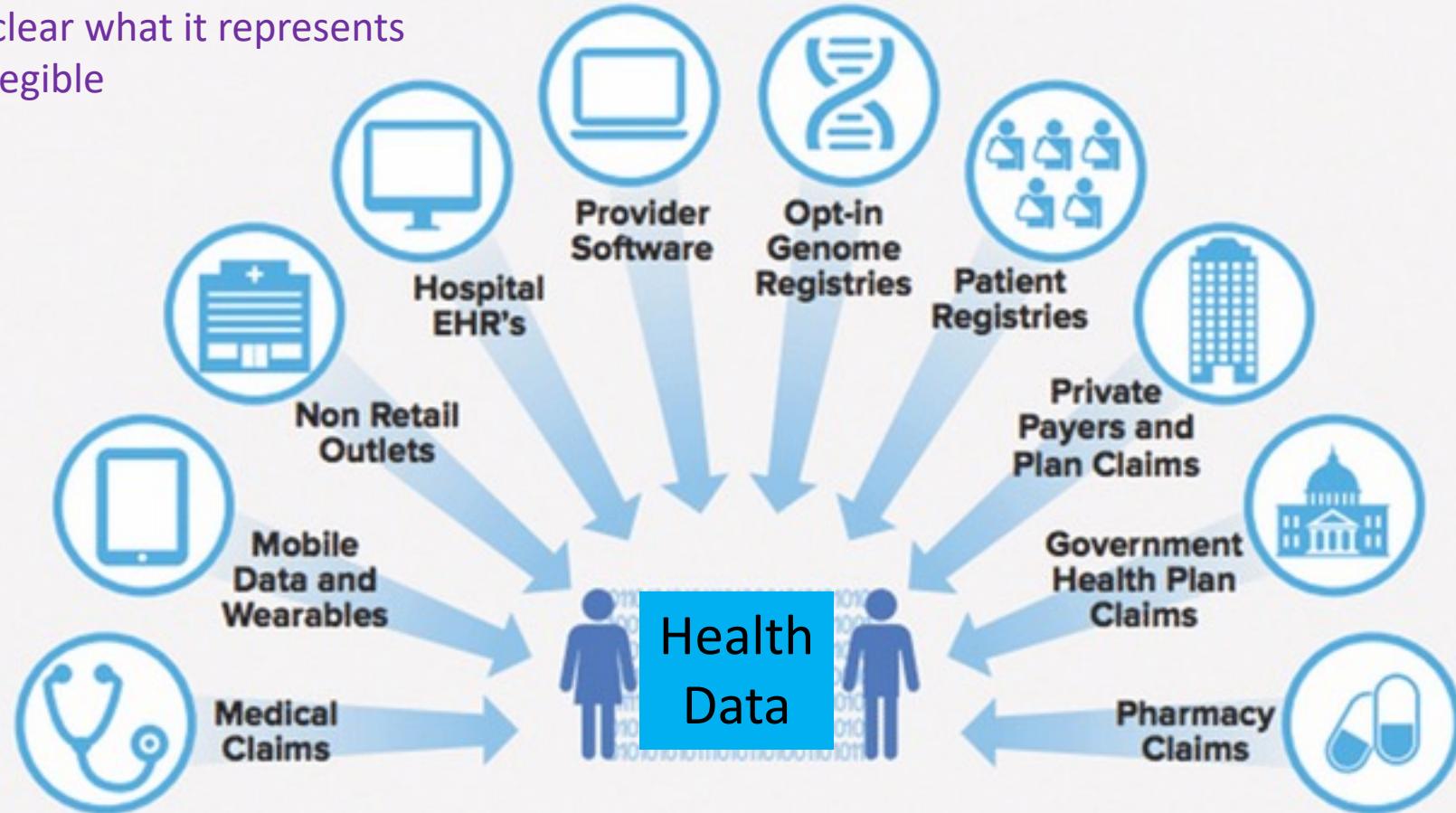
What is data?

It may not be clear what it represents



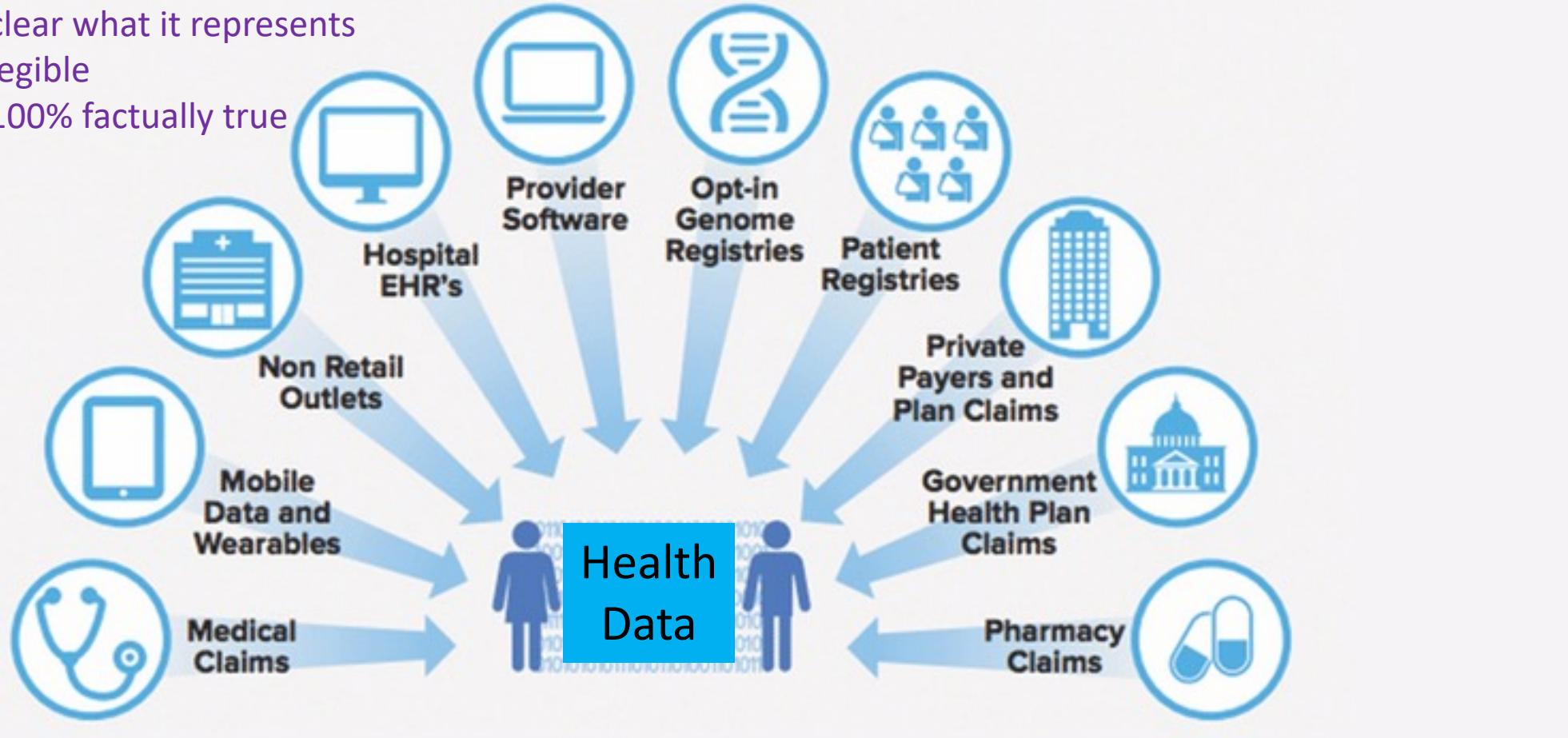
What is data?

It may not be clear what it represents
It may not be legible



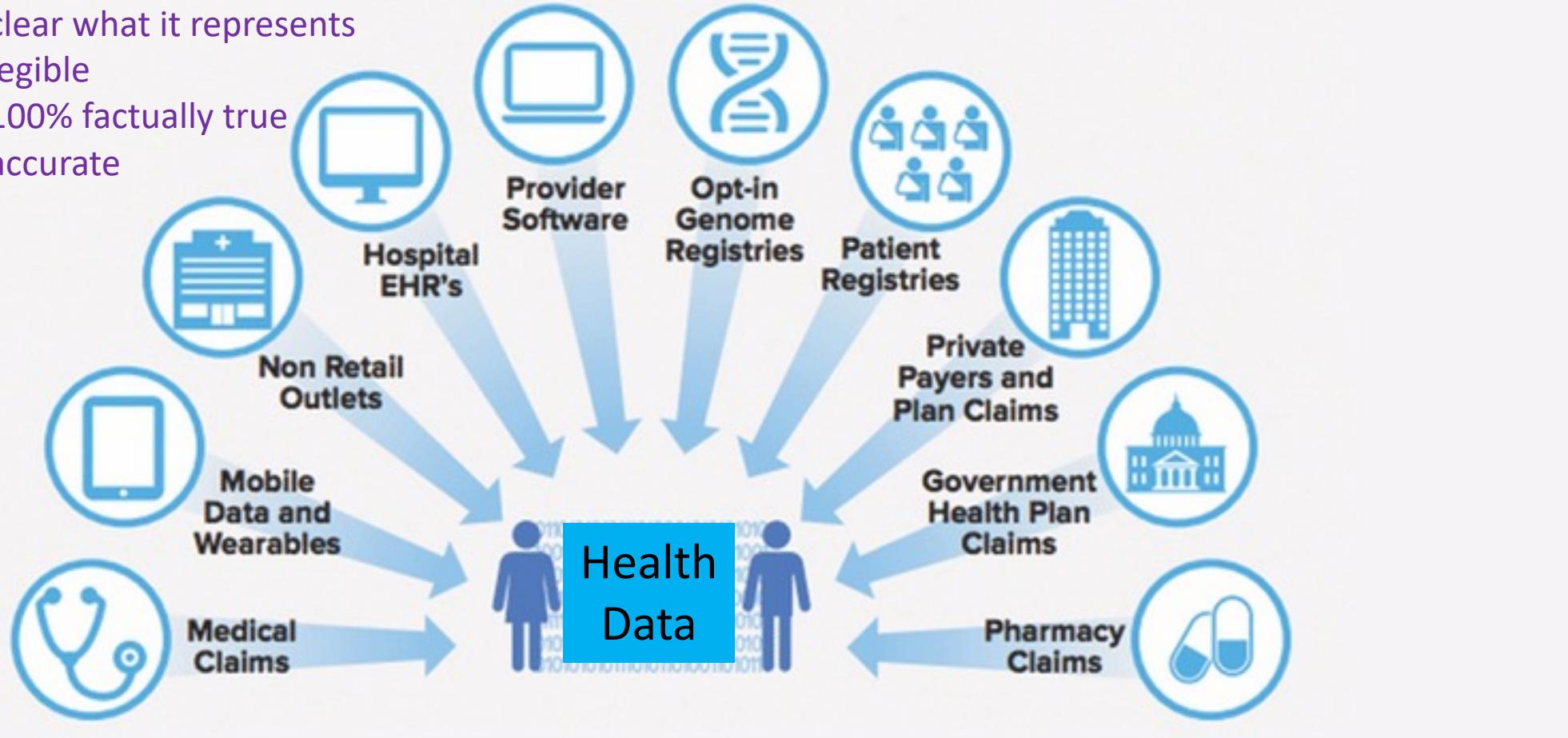
What is data?

It may not be clear what it represents
It may not be legible
It may not be 100% factually true



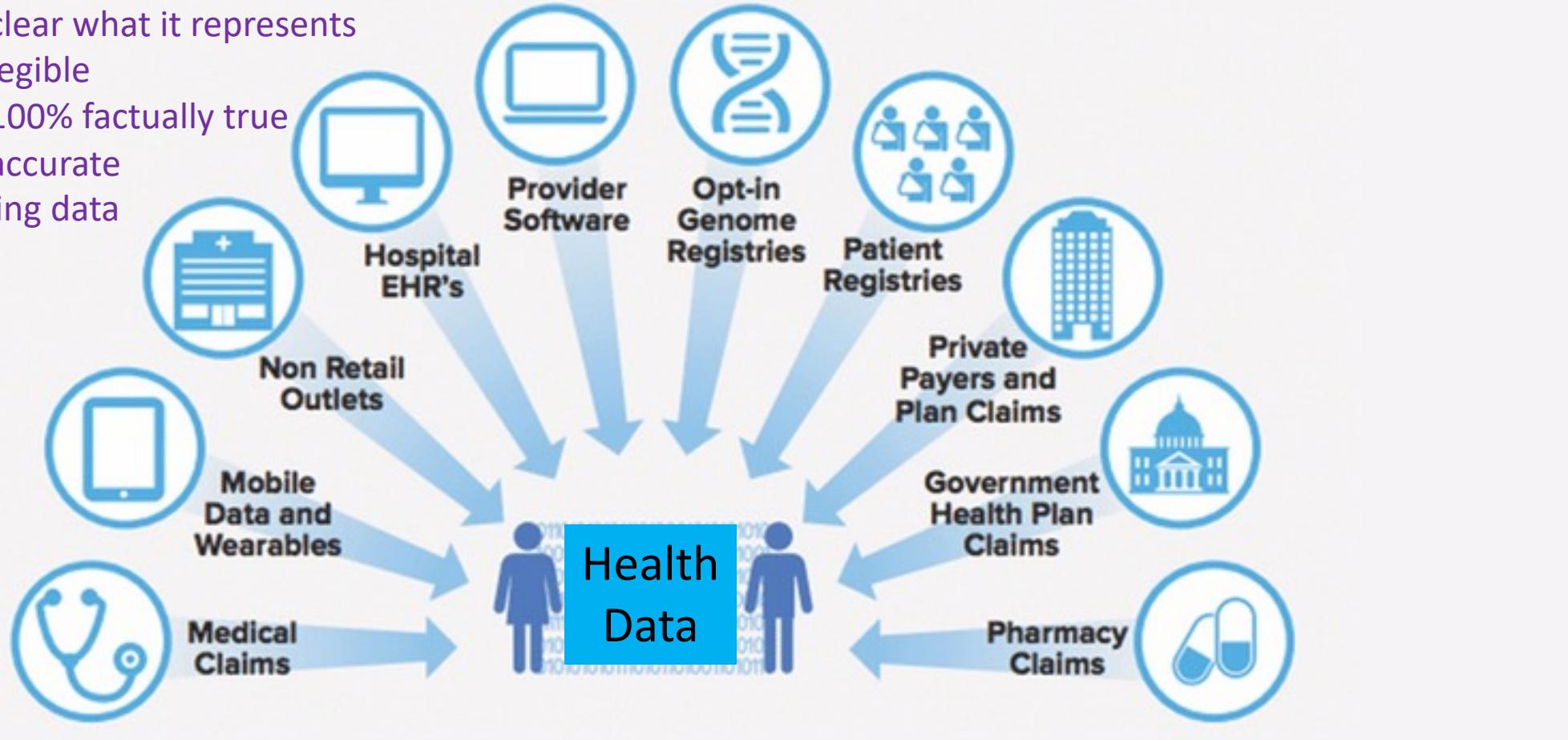
What is data?

It may not be clear what it represents
It may not be legible
It may not be 100% factually true
It may not be accurate

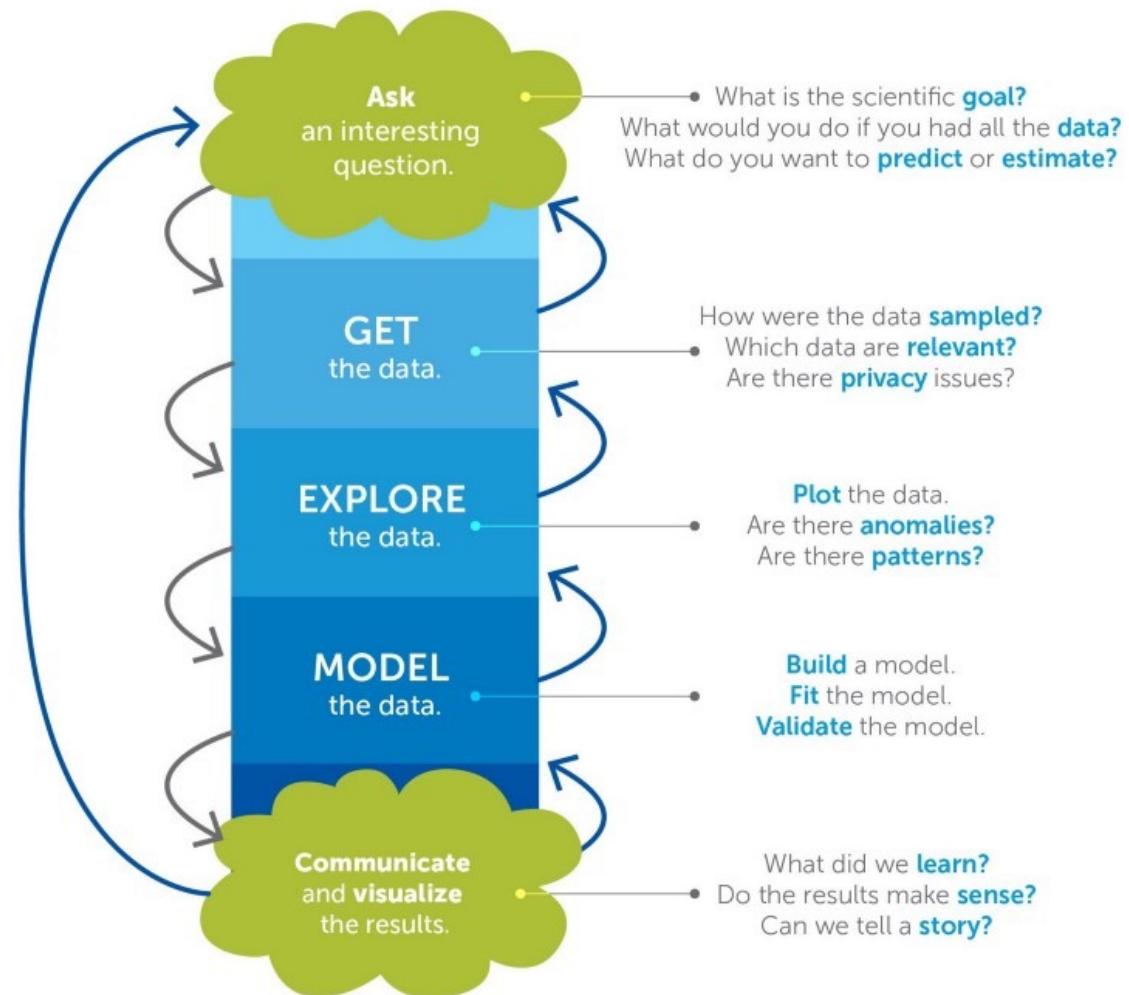


What is data?

It may not be clear what it represents
It may not be legible
It may not be 100% factually true
It may not be accurate
It may be missing data



The Data Science Process



Derived from the work of Joe Blitzstein and Hanspeter Pfister,
originally created for the Harvard data science course <http://cs109.org/>.

Getting data

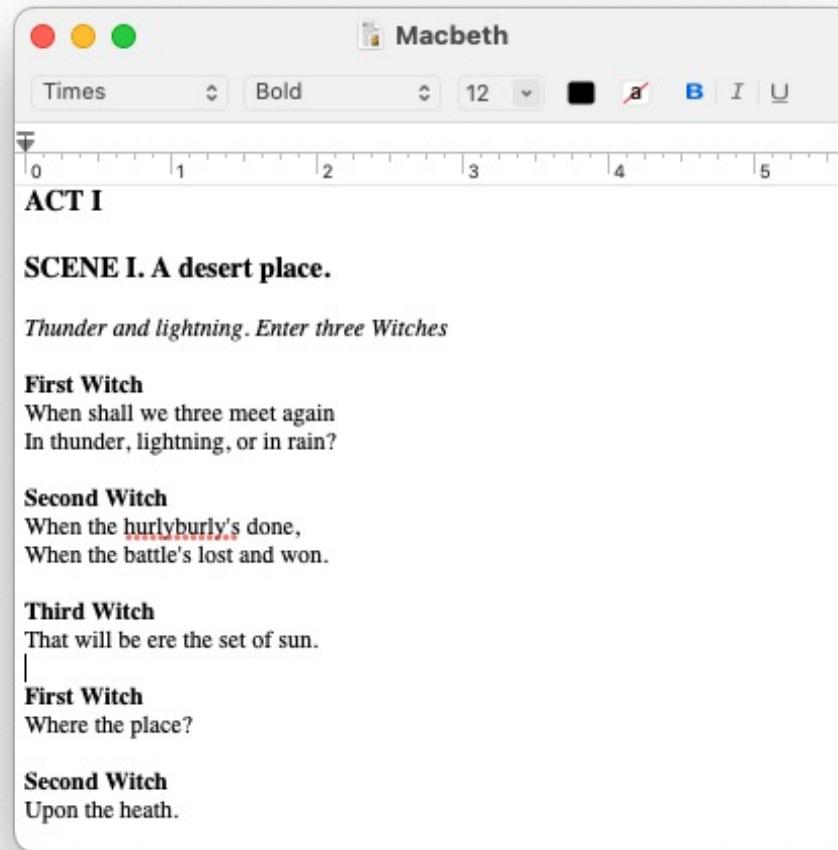
- What data is necessary to answer our question?
- Who collected it?
 - Is it from a reliable source? An authoritative source? (.com, .gov, .org)?
- When and where was it collected?
- How is it stored?
- How much data is necessary?
 - Comprehensive data set or sampled data set?
 - Are there any biases in the collection method?
- Are there restrictions or limitations for using it?
 - Is it licensed?
 - How difficult is it to analyze?
 - What is the allowed usage of data under its license?
- Are there humanistic concerns related to the data?
 - Did its collection need to be approved by an IRB?
 - Does its storage and use fall under HIPAA requirements? FERPA requirements?

Using data: formats

- Text
 - Web standards: html, xml, css, svg, json, ...
 - Source code: c, cpp, h, cs, js, py, java, rb, pl, php, sh, ...
 - Documents: txt, tex, markdown, asciidoc, rtf, ps, ...
 - Configuration: ini, cfg, rc, reg, ...
 - Tabular data: csv, tsv, ...
- Binary
 - Images: jpg, png, gif, bmp, tiff, psd, ...
 - Videos: mp4, mkv, avi, mov, mpg, vob, ...
 - Audio: mp3, aac, wav, flac, ogg, mka, wma, ...
 - Documents: pdf, doc, xls, ppt, docx, odt, ...
 - Archive: zip, rar, 7z, tar, iso, ...
 - Database: mdb, accde, frm, sqlite, ...
 - Executable: exe, dll, so, class, ...

Using data: format examples

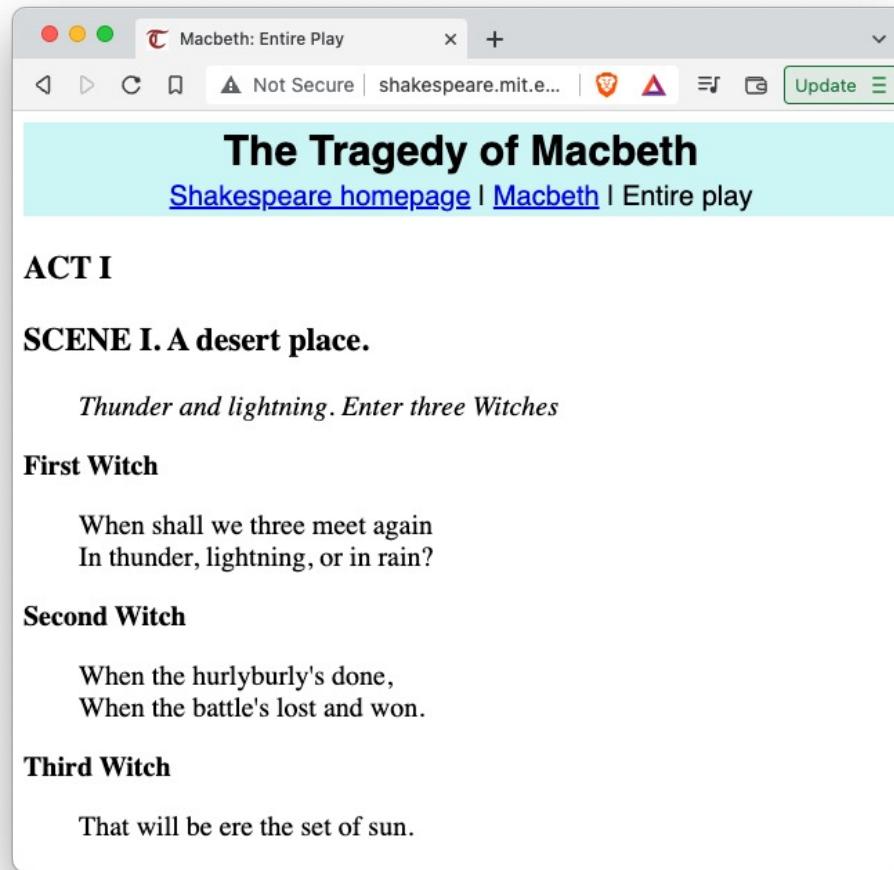
- Text
 - Plain text file (*.txt)



Using data: format examples

- Text

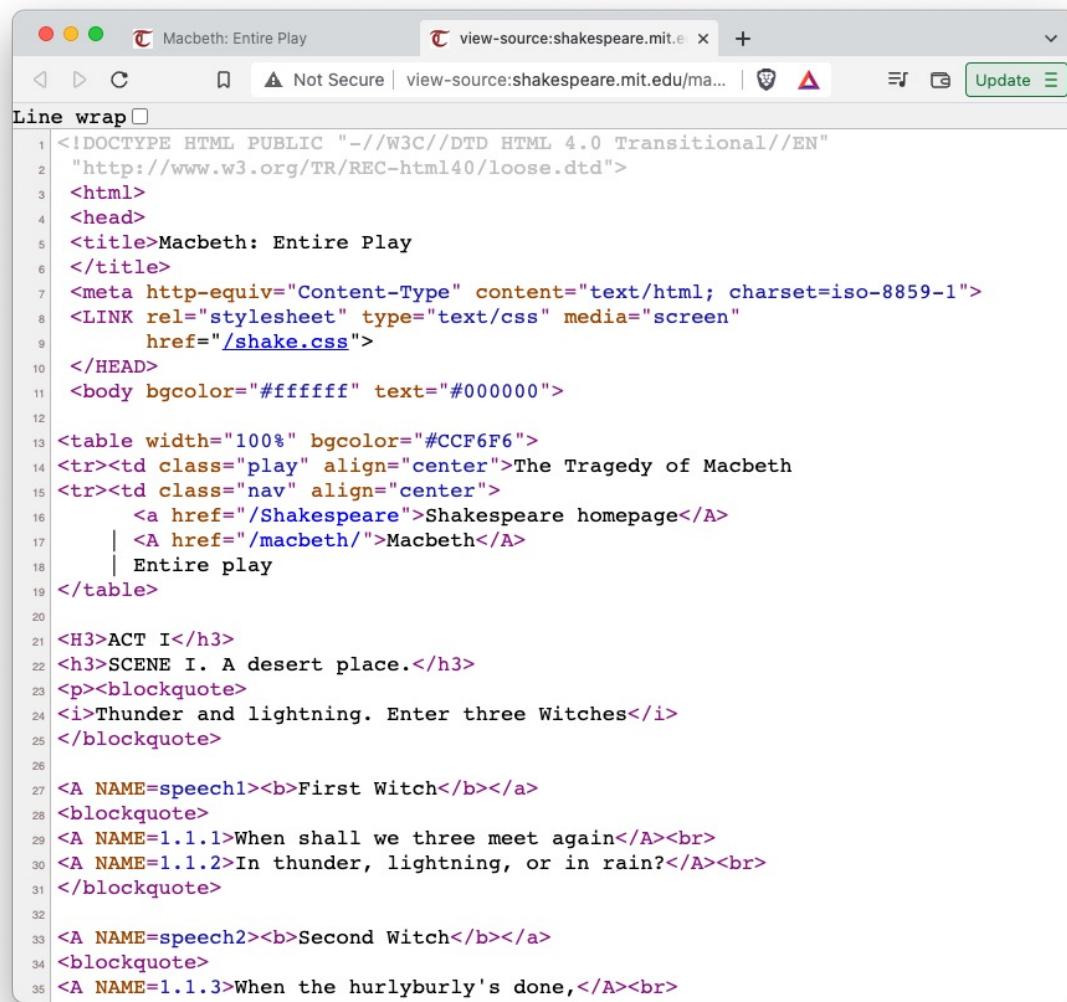
- Plain text file (*.txt)
- HTML file (*.html)



Using data: format examples

- Text

- Plain text file (*.txt)
- HTML file (*.html)



The screenshot shows a web browser window with the title "Macbeth: Entire Play". The address bar indicates the page is "Not Secure" and shows the URL "view-source:shakespeare.mit.edu/macbeth.html". The browser interface includes standard controls like back, forward, and search, along with a "Line wrap" checkbox. The main content area displays the HTML source code for the play. The code is color-coded for syntax highlighting, with tags in blue, attributes in green, and values in purple. The HTML structure includes a head section with a title, meta tags, and a link to a stylesheet. The body section contains a table for navigation, followed by act and scene descriptions, and three blockquote sections for the Witches' speeches.

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.0 Transitional//EN"
"http://www.w3.org/TR/REC-html40/loose.dtd">
<html>
<head>
<title>Macbeth: Entire Play
</title>
<meta http-equiv="Content-Type" content="text/html; charset=iso-8859-1">
<LINK rel="stylesheet" type="text/css" media="screen"
      href="/shake.css">
</HEAD>
<body bgcolor="#ffffff" text="#000000">

<table width="100%" bgcolor="#CCF6F6">
<tr><td class="play" align="center">The Tragedy of Macbeth
<tr><td class="nav" align="center">
    <a href="/Shakespeare">Shakespeare homepage</A>
    | <a href="/macbeth/">Macbeth</A>
    | Entire play
</table>

<H3>ACT I</h3>
<h3>SCENE I. A desert place.</h3>
<p><blockquote>
<i>Thunder and lightning. Enter three Witches</i>
</blockquote>
<A NAME=speech1><b>First Witch</b></a>
<blockquote>
<A NAME=1.1.1>When shall we three meet again</A><br>
<A NAME=1.1.2>In thunder, lightning, or in rain?</A><br>
</blockquote>
<A NAME=speech2><b>Second Witch</b></a>
<blockquote>
<A NAME=1.1.3>When the hurlyburly's done,</A><br>
```

Using data: format examples

- Text

- Plain text file (*.txt)
- HTML file (*.html)
- Markdown (*.md)

The image shows a code editor with two tabs: 'macbeth.md' and 'macbeth.html'. The 'macbeth.md' tab displays the following Markdown code:

```
1 # ACT I
2 ## SCENE I. A desert place.
3
4 *Thunder and lightning. Enter three Witches*
5
6 **First Witch**
7 When shall we three meet again
8 In thunder, lightning, or in rain?
9
10 **Second Witch**
11 When the hurlyburly's done,
12 When the battle's lost and won.
13
14 **Third Witch**
15 That will be ere the set of sun.
16
17 **First Witch**
18 Where the place?
19
20 **Second Witch**
21 Upon the heath.
```

The 'macbeth.html' tab displays the resulting HTML output:

ACT I

SCENE I. A desert place.

Thunder and lightning. Enter three Witches

First Witch When shall we three meet again In thunder, lightning, or in rain?

Second Witch When the hurlyburly's done, When the battle's lost and won.

Third Witch That will be ere the set of sun.

First Witch Where the place?

Second Witch Upon the heath.

Using data: format examples

- Text

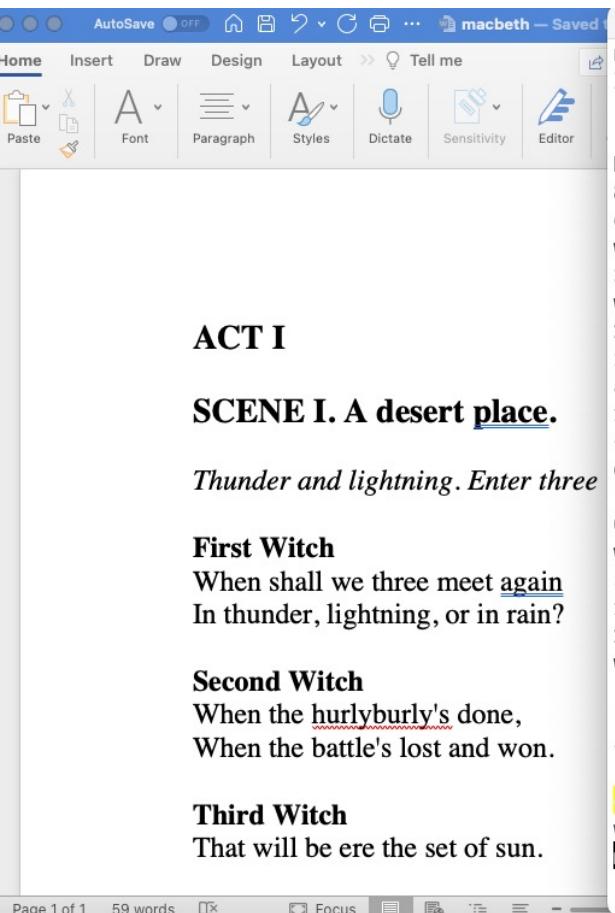
- Plain text file (*.txt)
- HTML file (*.html)
- Markdown (*.md)
- XML (*.xml)

```
<roll_call_vote>
<congress>115</congress>
<session>1</session>
...
<members>
  <member>
    <member_full>Alexander (R-TN)</member_full>
    <last_name>Alexander</last_name>
    <first_name>Lamar</first_name>
    <party>R</party>
    <state>TN</state>
    <vote_cast>Yea</vote_cast>
  ...
</member>
</members>
</roll_call_vote>
```

Using data: format examples

- Text

- Plain text file (*.txt)
 - HTML file (*.html)
 - Markdown (*.md)
 - XML (*.xml)
 - MS Word docs (*.docx)



nts w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> </w:rPr> <w:t>hur
lyburly's</w:t></w:r> <w:proofErr w:type="spellEnd"/> <w:r> <w:rPr> <w:
rFonts w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> </w:rPr> <w:t
xml:space="preserve"> done,</w:t></w:r></w:p> <w:p w14:paraId="0EC
D7433" w14:textId="77777777" w:rsidR="0080320C" w:rsidRDefault="00
80320C" w:rsidP="0080320C"> <w:pPr> <w:autoSpaceDE w:val="0"/> <w:aut
oSpaceDN w:val="0"/> <w:adjustRightInd w:val="0"/> <w:rPr> <w:rFonts
w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> </w:rPr> <w:r>
<w:rPr> <w:rFonts w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> </w:
rPr> <w:t>When the battle's lost and won.</w:t></w:r></w:p> <w:p w
14:paraId="7AEF7063" w14:textId="77777777" w:rsidR="0080320C" w:rs
idRDefault="0080320C" w:rsidP="0080320C"> <w:pPr> <w:autoSpaceDE w:v
al="0"/> <w:adjustRightInd w:val="0"/> <w:
rPr> <w:rFonts w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> <w:b>
<w:bCs/> </w:rPr> </w:pPr> </w:p> <w:p w14:paraId="4EB77814" w14:textI
d="77777777" w:rsidR="0080320C" w:rsidRDefault="0080320C" w:rsidP=
"0080320C"> <w:pPr> <w:autoSpaceDE w:val="0"/> <w:autoSpaceDN w:val="0"
/> <w:adjustRightInd w:val="0"/> <w:rPr> <w:rFonts w:ascii="Times"
w:hAnsi="Times" w:cs="Times"/> </w:rPr> <w:pPr> <w:r> <w:rPr> <w:rFont
s w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> </w:pPr> <w:r> <w:rPr>
<w:t>Third Witch</w:t></w:r></w:p> <w:p w14:paraId="6F8689AE" w
14:textId="77777777" w:rsidR="0080320C" w:rsidRDefault="0080320C"
w:rsidP="0080320C"> <w:pPr> <w:autoSpaceDE w:val="0"/> <w:adjustRightInd
w:val="0"/> <w:rPr> <w:rFonts w:ascii="Times" w:hAnsi="Times" w:cs="Times"/>
</w:rPr> <w:rFonts w:ascii="Times" w:hAnsi="Times" w:cs="Times"/> </w:
rPr> <w:t>That will be ere the set of sun.</w:t></w:r></w:p> <w:p w14:para

<pfile:///Users/bwinjum/Downloads/macbeth.docx:word/document.xml

word/document.xml

macbeth.docx [RO]

Using data: format examples

- Text

- Plain text file (*.txt)
- HTML file (*.html)
- Markdown (*.md)
- XML (*.xml)
 - MS Word docs (*.docx)
- JSON (*.json)

The image shows two side-by-side Jupyter Notebook interfaces. The left interface displays a JSON-formatted cell containing the structure of a Jupyter notebook cell. The right interface shows a Python script cell that imports the 'diabetes' dataset from scikit-learn, prints its description, and then displays the dataset's characteristics.

```
[1]: import numpy as np
import matplotlib.pyplot as plt
import sklearn.datasets

[2]: diabetes = sklearn.datasets.load_diabetes()
print(diabetes.DESCR)

.. _diabetes_dataset:

Diabetes dataset
-----
Ten baseline variables, age, sex, body mass index, average blood pressure, and six blood serum measurements were obtained for each of n = 442 diabetes patients, as well as the response of interest, a quantitative measure of disease progression one year after baseline.

**Data Set Characteristics:** 

:Number of Instances: 442

:Number of Attributes: First 10 columns are numeric predictive values

:Target: Column 11 is a quantitative measure of disease progression one year after baseline

:Attribute Information:
 - age      age in years
 - sex
 - bmi     body mass index
 - bp      average blood pressure
 - s1      tc, total serum cholesterol
```

```
{ "cells": [ { "cell_type": "code", "execution_count": 1, "id": "53c98d7b-05cd-4ca2-ac80-fad0313478cd", "metadata": {}, "outputs": [], "source": [ "import numpy as np\\n", "import matplotlib.pyplot as plt\\n", "import sklearn.datasets" ] }, { "cell_type": "code", "execution_count": 2, "id": "acf6ec46-4832-4a34-a475-62f759979b8e", "metadata": {}, "outputs": [ { "name": "stdout", "output_type": "stream", "text": [ "... _diabetes_dataset:\\n", "\\n", "Diabetes dataset\\n", "-----\\n", "\\n", "Ten baseline variables, age, sex, body mass index, average blood\\n", "pressure, and six blood serum measurements were obtained for each of n =\\n", "442 diabetes patients, as well as the response of interest, a\\n", "quantitative measure of disease progression one year after baseline.\\n", "\\n" ] } ] } ] }
```

Using data: format examples

- **Text**

- Plain text file (*.txt)
- HTML file (*.html)
- Markdown (*.md)
- XML (*.xml)
 - MS Word docs (*.docx)
- JSON (*.json)
- CSV (*.csv)

```
taken,person,quant,reading
619,dyer,rad,9.82
619,dyer,sal,0.13
622,dyer,rad,7.8
622,dyer,sal,0.09
734,pb,rad,8.41
734,lake,sal,0.05
734,pb,temp,-21.5
735,pb,rad,7.22
735,,sal,0.06
735,,temp,-26.0
751,pb,rad,4.35
751,pb,temp,-18.5
751,lake,sal,0.1
752,lake,rad,2.19
752,lake,sal,0.09
752,lake,temp,-16.0
752,roe,sal,41.6
837,lake,rad,1.46
837,lake,sal,0.21
837,roe,sal,22.5
844,roe,rad,11.25
```

Using data: format examples

- Text

- Plain text file (*.txt)
- HTML file (*.html)
- Markdown (*.md)
- XML (*.xml)
 - MS Word docs (*.docx)
- JSON (*.json)
- CSV (*.csv)
- TSV (*.tsv)

taken	person	quant	reading
619	dyer	rad	9.82
619	dyer	sal	0.13
622	dyer	rad	7.8
622	dyer	sal	0.09
734	pb	rad	8.41
734	lake	sal	0.05
734	pb	temp	-21.5
735	pb	rad	7.22
735		sal	0.06
735		temp	-26.0
751	pb	rad	4.35
751	pb	temp	-18.5
751	lake	sal	0.1
752	lake	rad	2.19
752	lake	sal	0.09
752	lake	temp	-16.0
752	roe	sal	41.6
837	lake	rad	1.46
837	lake	sal	0.21
837	roe	sal	22.5
844	roe	rad	11.25

Essential Python skill: Importing data

– we're going to start with basic files & Pandas

Format Type	Data Description	Reader	Writer
text	CSV	<code>read_csv</code>	<code>to_csv</code>
text	Fixed-Width Text File	<code>read_fwf</code>	
text	JSON	<code>read_json</code>	<code>to_json</code>
text	HTML	<code>read_html</code>	<code>to_html</code>
text	LaTeX		<code>Styler.to_latex</code>
text	XML	<code>read_xml</code>	<code>to_xml</code>
text	Local clipboard	<code>read_clipboard</code>	<code>to_clipboard</code>
binary	MS Excel	<code>read_excel</code>	<code>to_excel</code>
binary	OpenDocument	<code>read_excel</code>	
binary	HDF5 Format	<code>read_hdf</code>	<code>to_hdf</code>
binary	Feather Format	<code>read_feather</code>	<code>to_feather</code>
binary	Parquet Format	<code>read_parquet</code>	<code>to_parquet</code>
binary	ORC Format	<code>read_orc</code>	
binary	Stata	<code>read_stata</code>	<code>to_stata</code>
binary	SAS	<code>read_sas</code>	
binary	SPSS	<code>read_spss</code>	
binary	Python Pickle Format	<code>read_pickle</code>	<code>to_pickle</code>
SQL	SQL	<code>read_sql</code>	<code>to_sql</code>
SQL	Google BigQuery	<code>read_gbq</code>	<code>to_gbq</code>

Code time

Messy data?

	Lord of the Rings	Chronicles of Narnia	Dark Tower
Book #1	700	150	2300
Book #2	823	235	1600
Book #3	1432	176	666

This makes Python happy

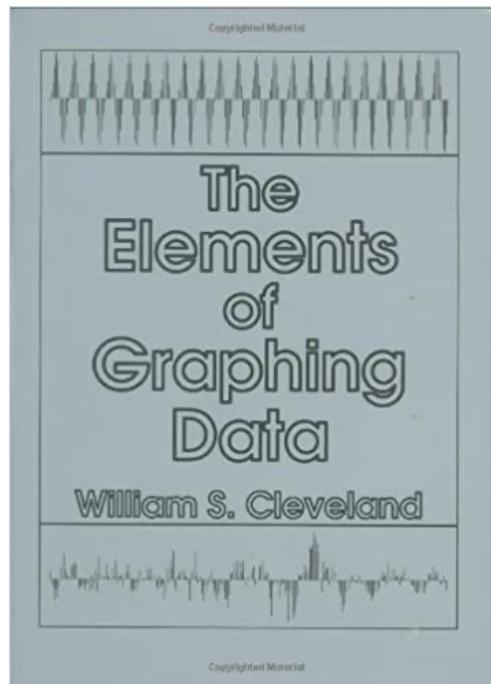
ID	Book Title	Book Number	Page Count
1	Lord of the Rings	Book #1	700
2	Lord of the Rings	Book #2	823
3	Lord of the Rings	Book #3	1432
4	Chronicles of Narnia	Book #1	150
5	Chronicles of Narnia	Book #2	235
6	Chronicles of Narnia	Book #3	176
7	Dark Tower	Book #1	2300
8	Dark Tower	Book #2	1600
9	Dark Tower	Book #3	666

Going Through the Assignments: Tableau

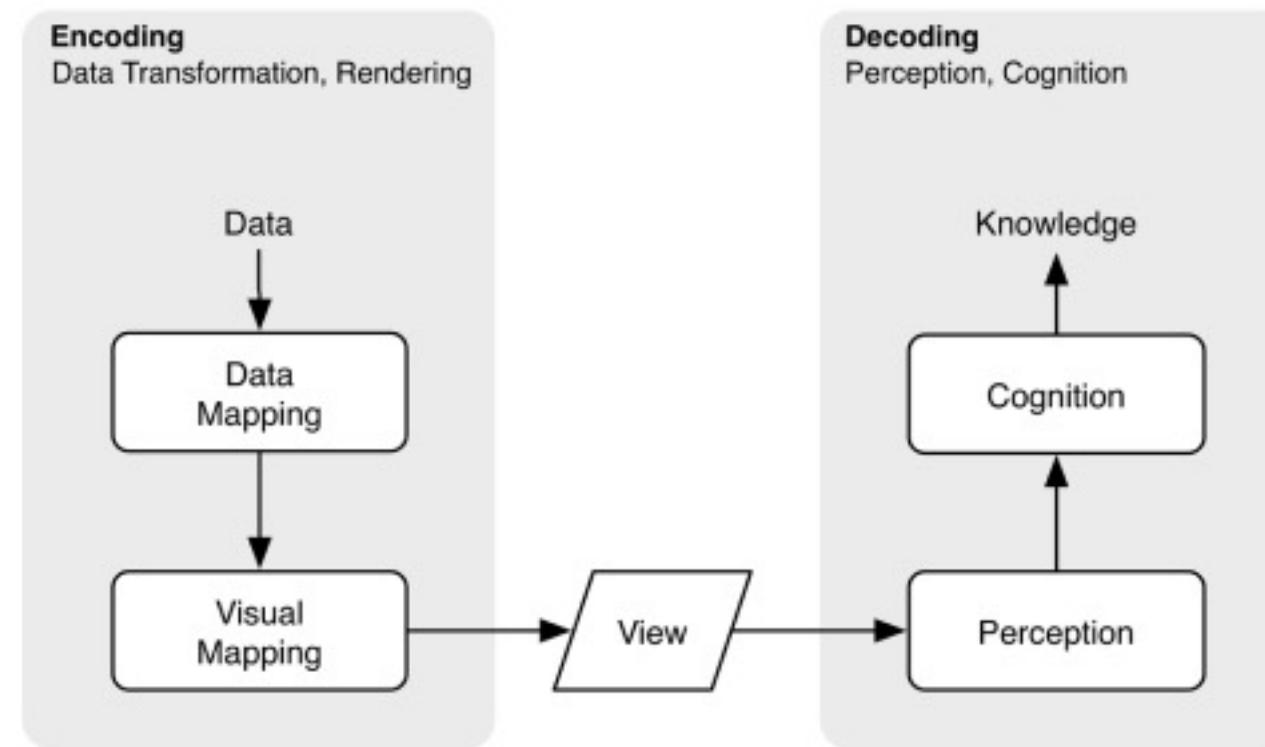
Elements of Graphing Data

Elements of Graphing Data

- The following slides all make use of points and graphics taken from Cleveland's "The Elements of Graphing Data"
- I highly recommend reading and thinking about this book



Good graphics must be visually clear and understandable



Graphical Perception

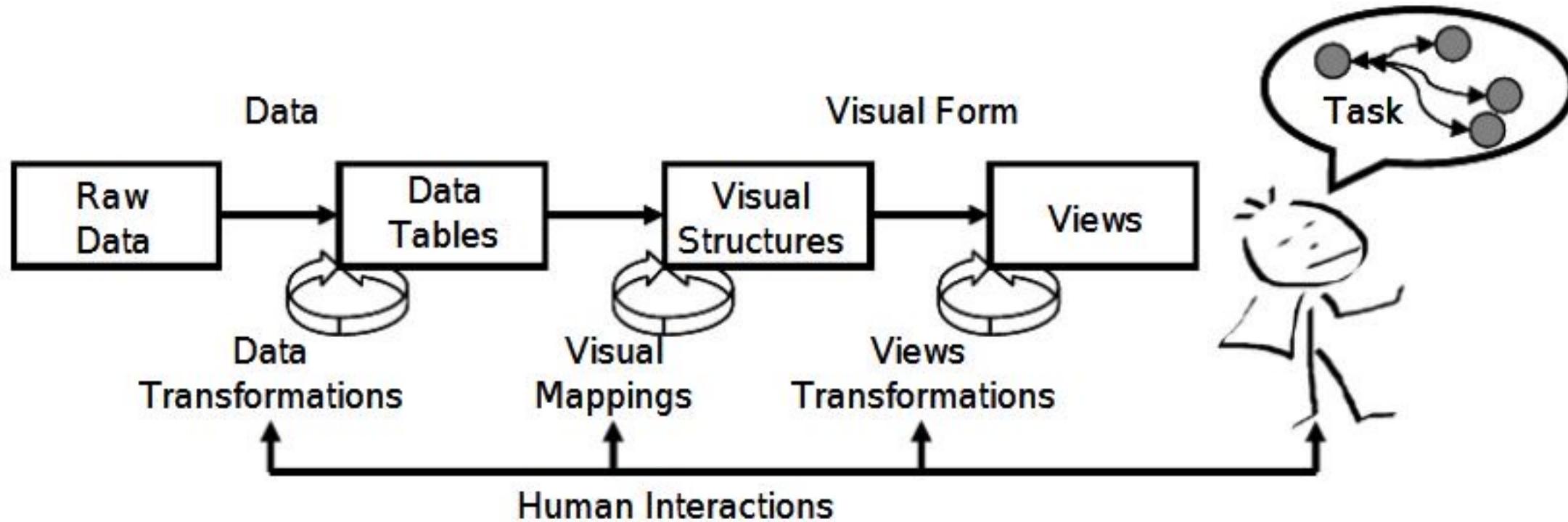
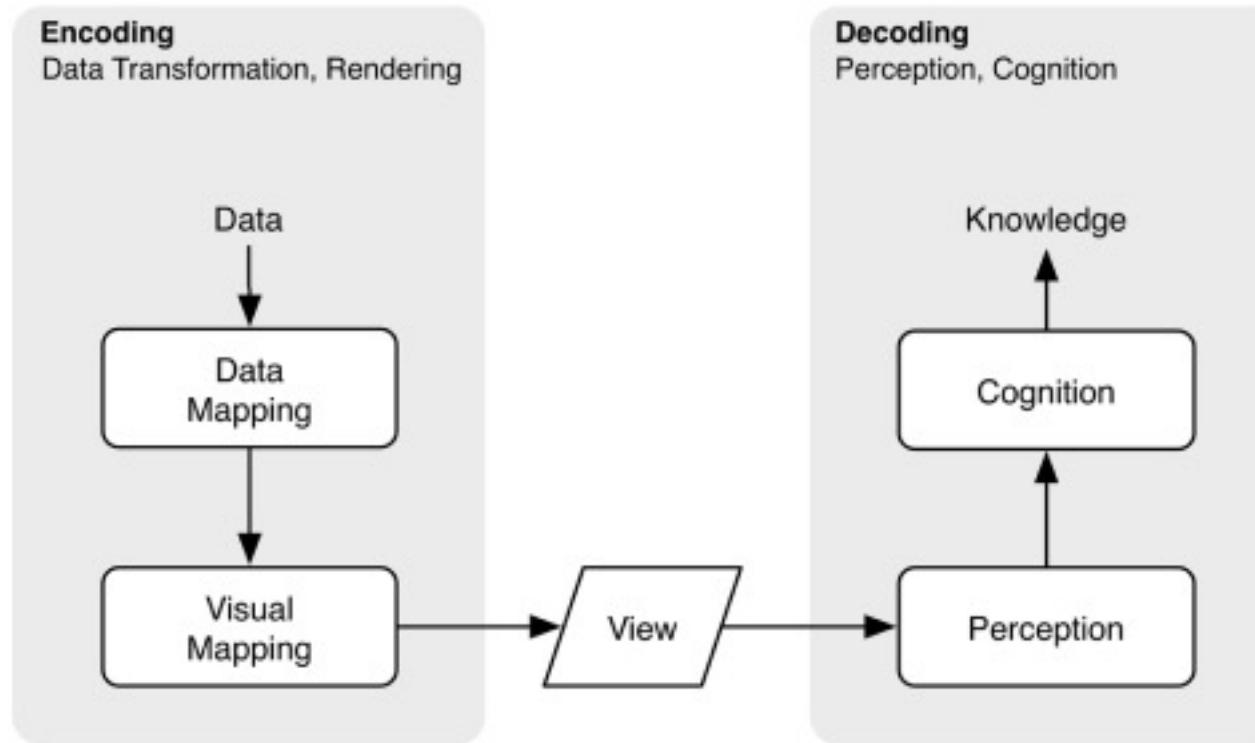


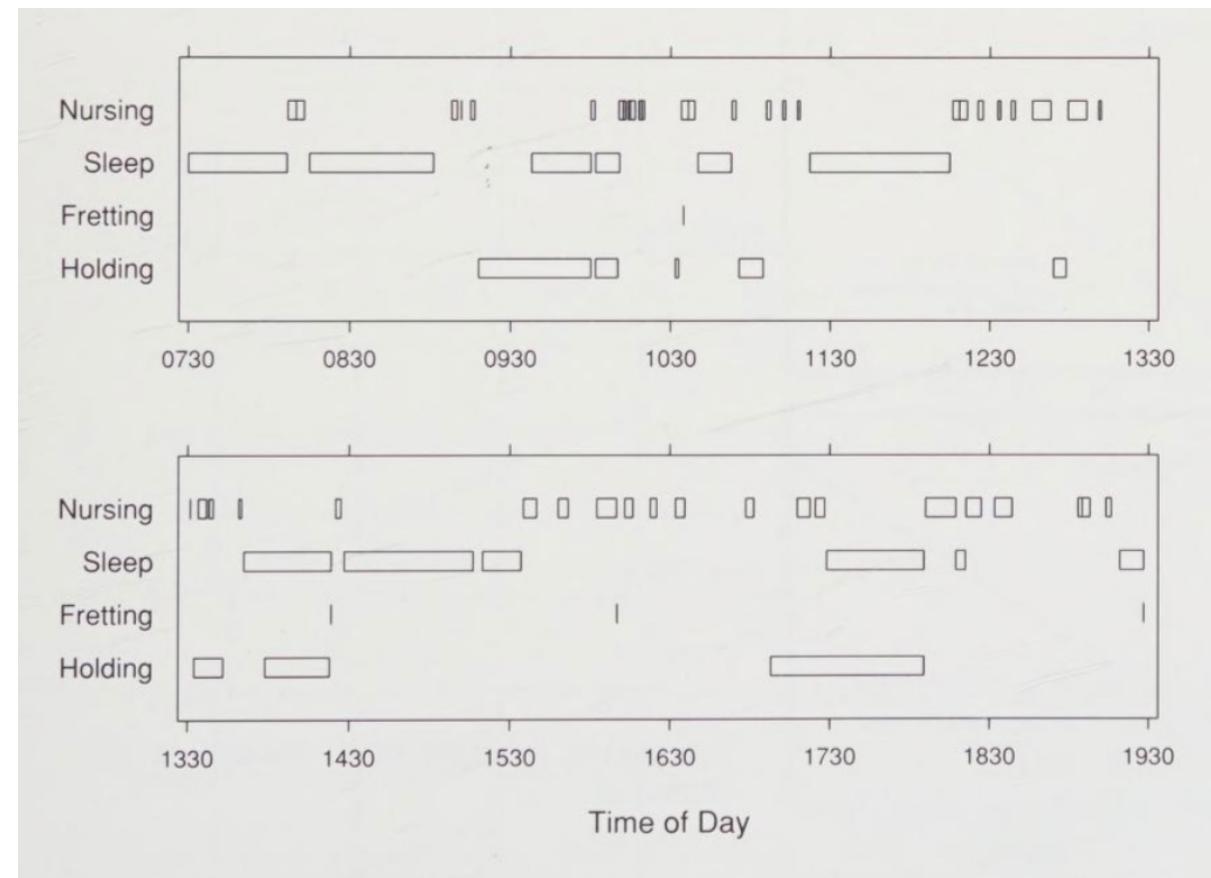
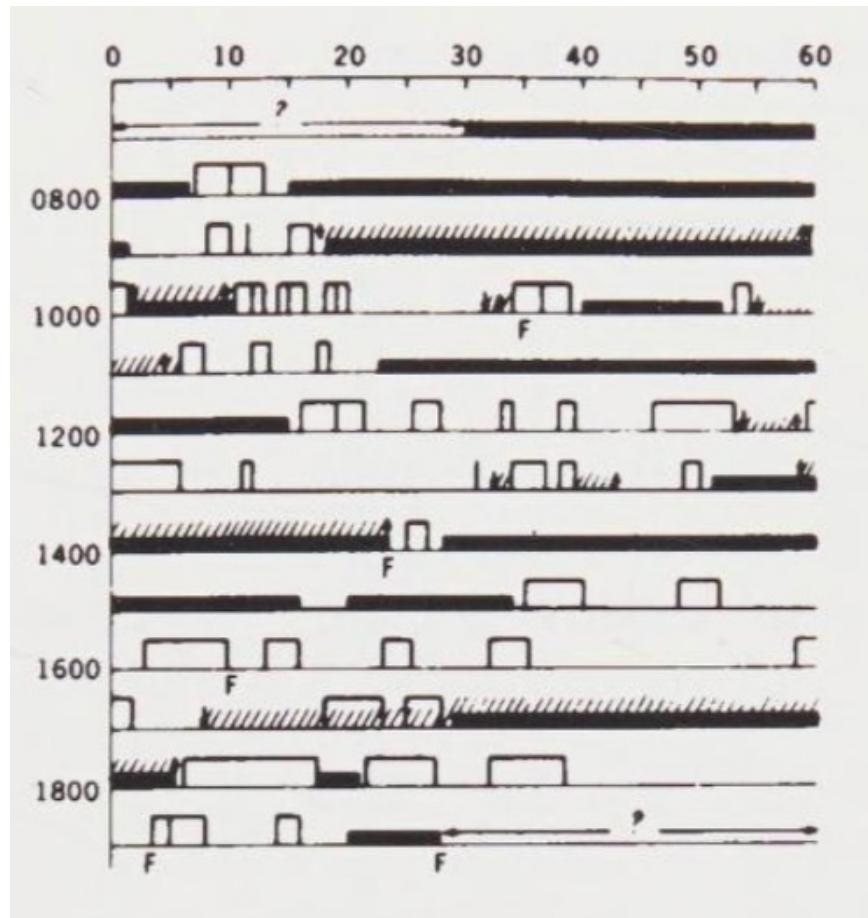
Figure 1 - Reference model for visualization (Card 1999).

Graphical Perception

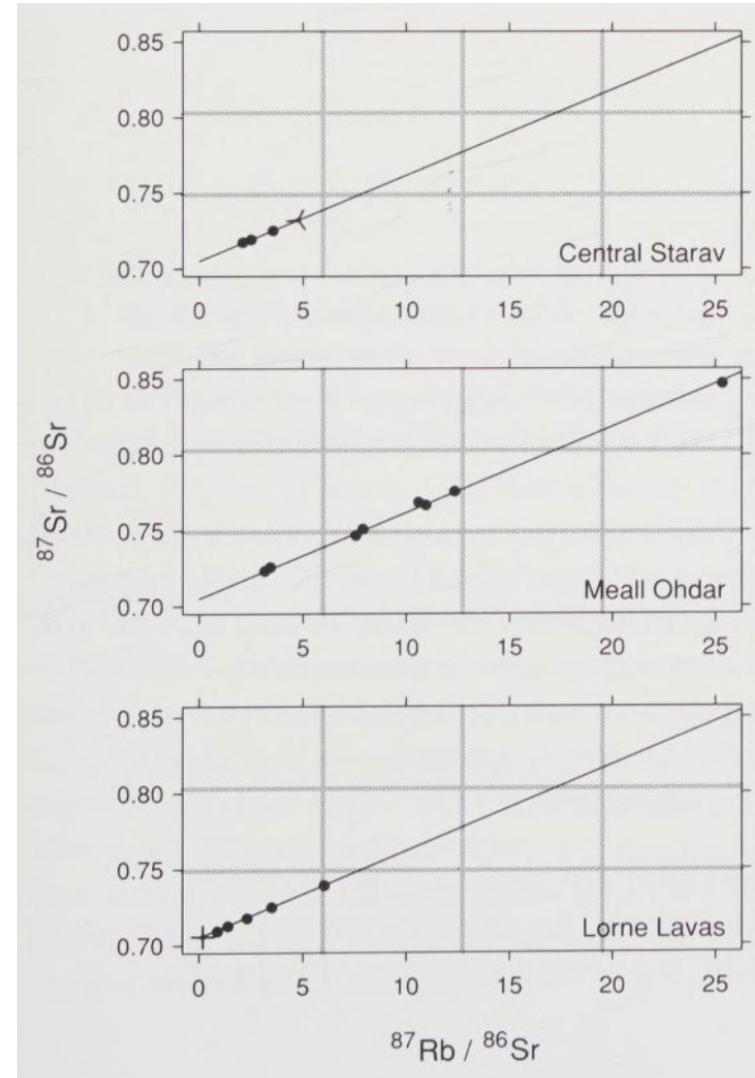
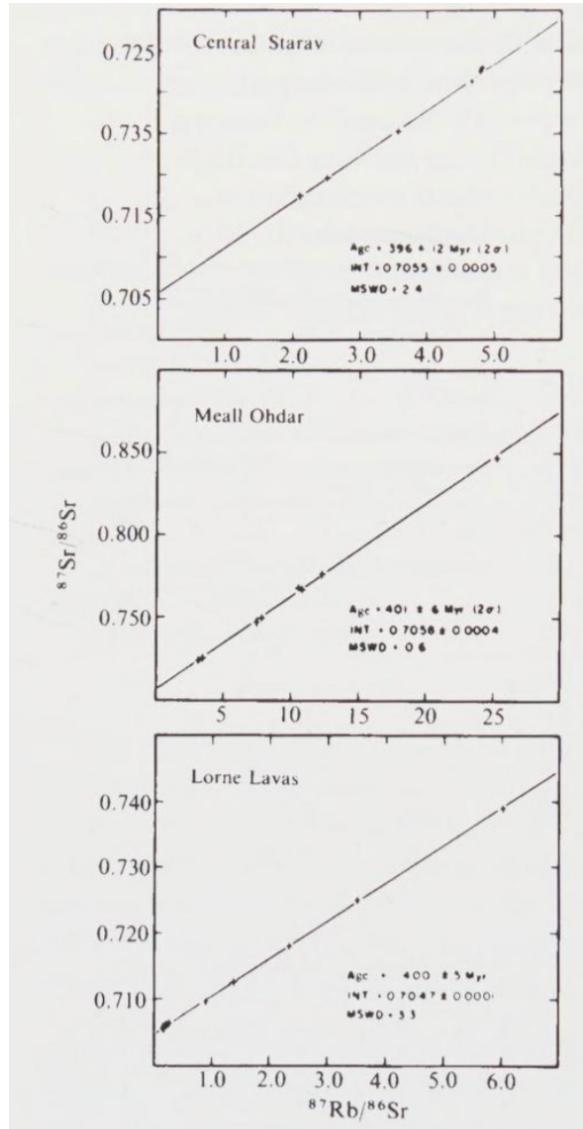


Clear Visual

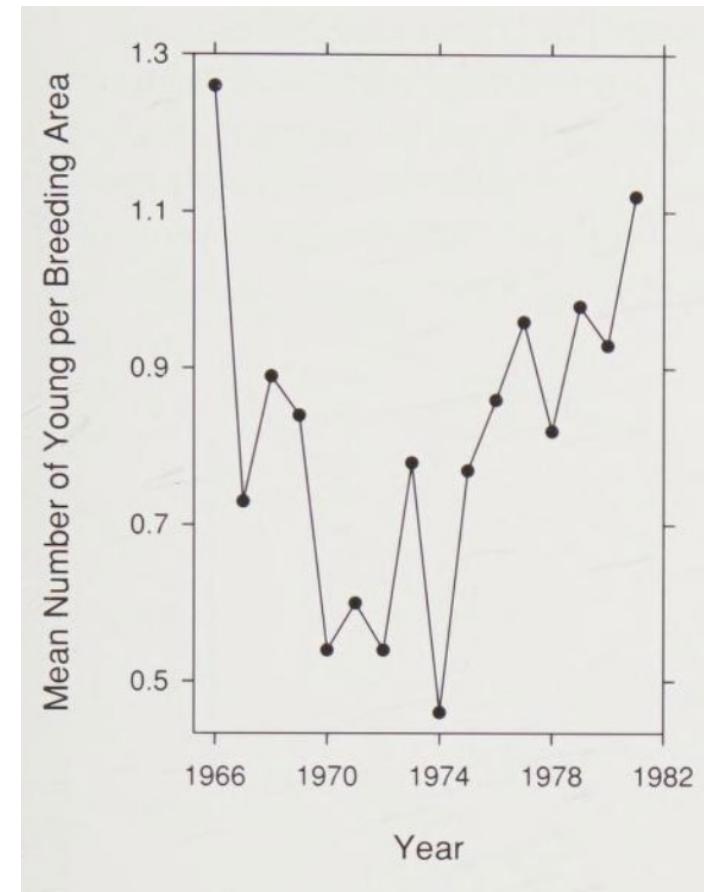
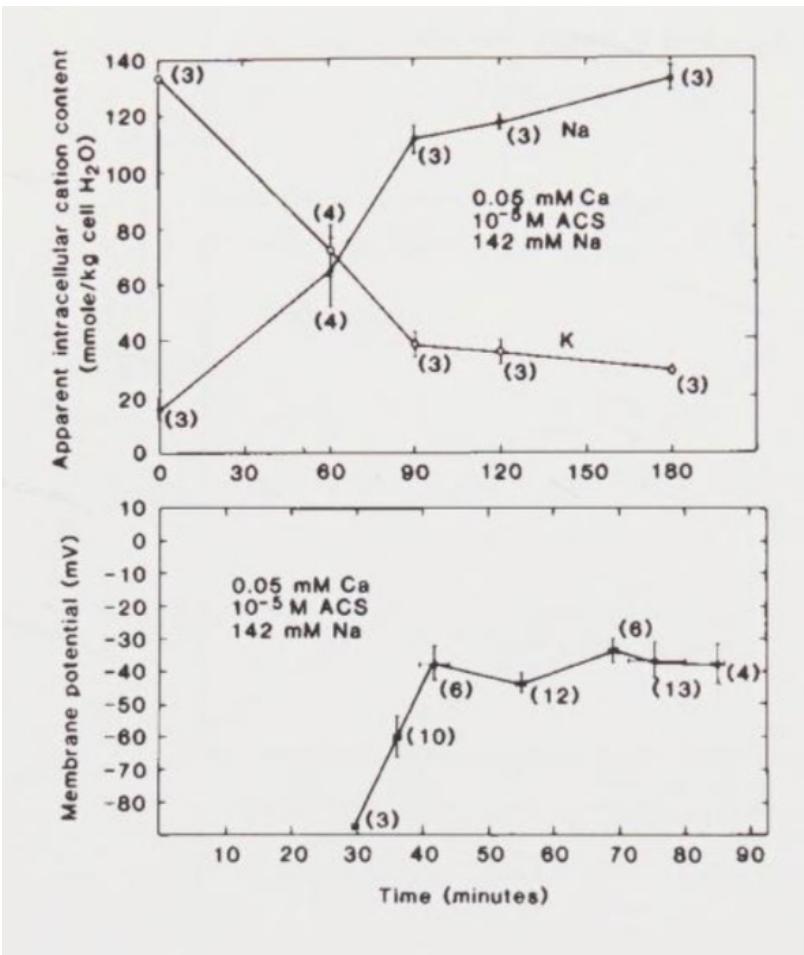
Make the data stand out; avoid superfluity



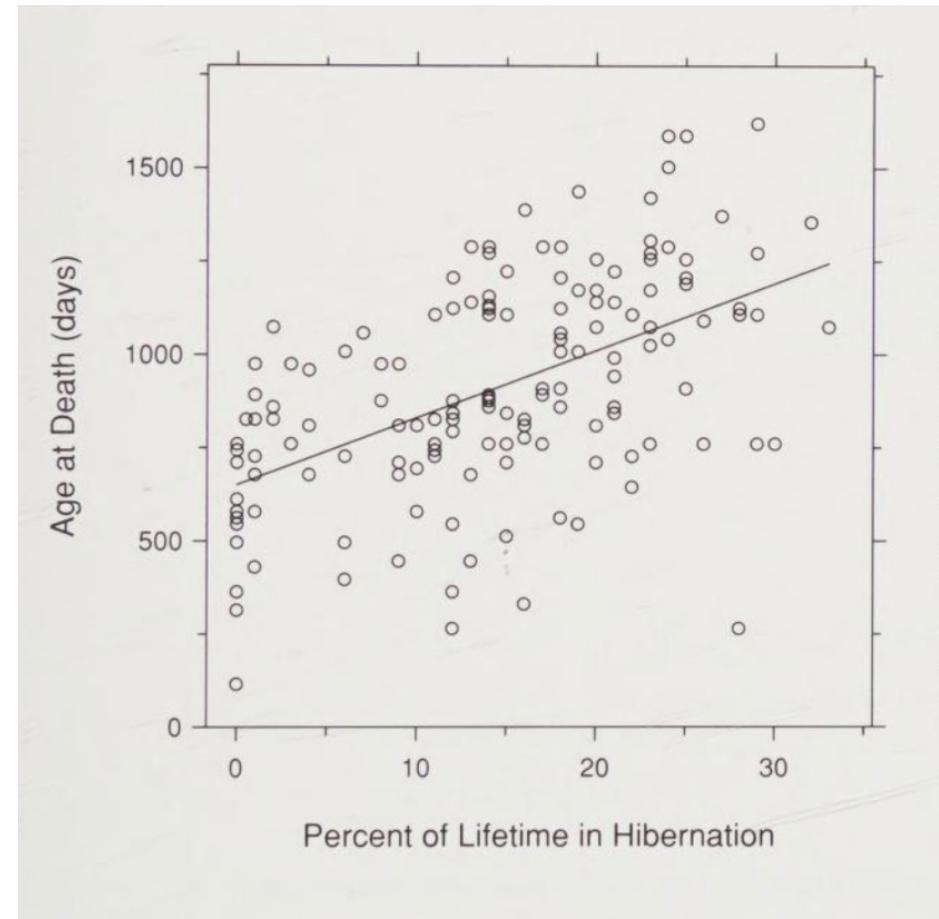
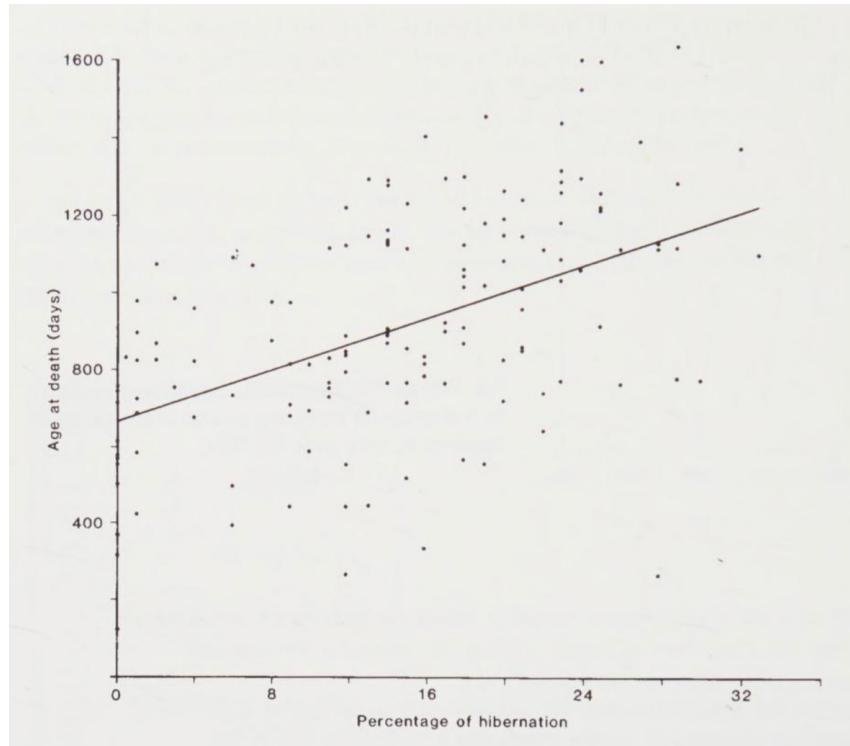
Use visually prominent graphical elements to show the data



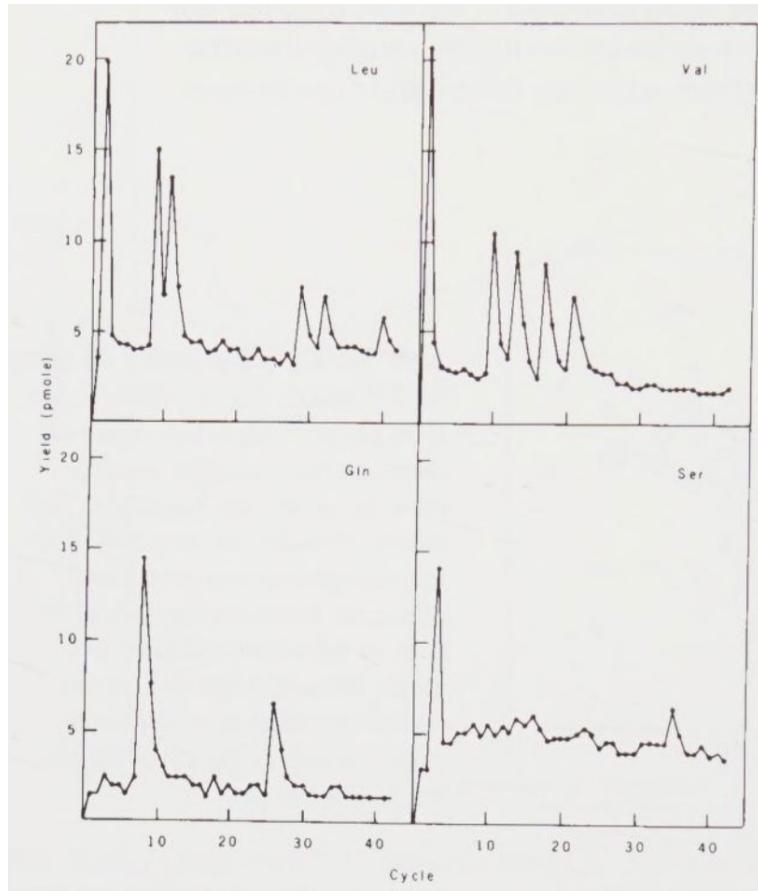
Use visually prominent graphical elements to show the data



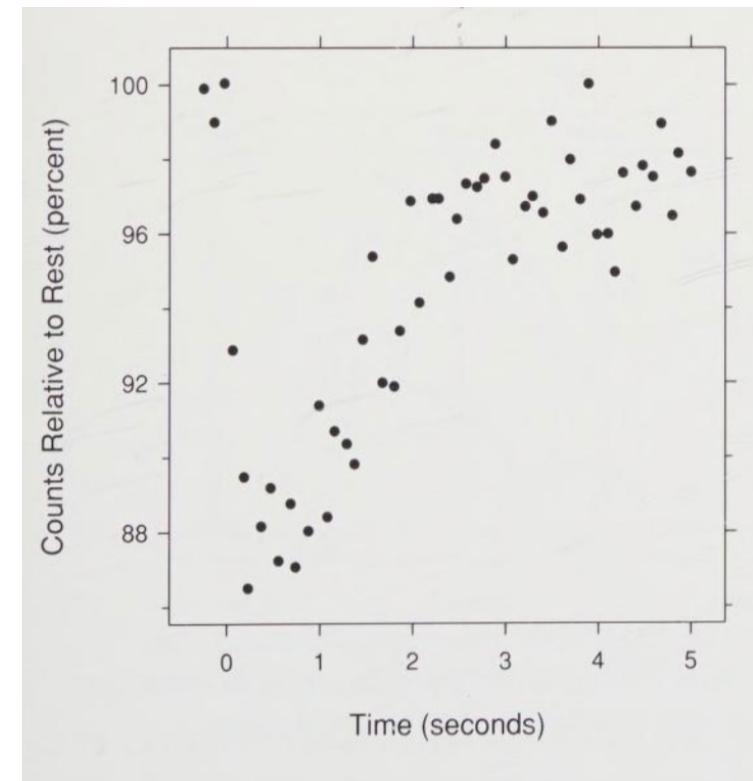
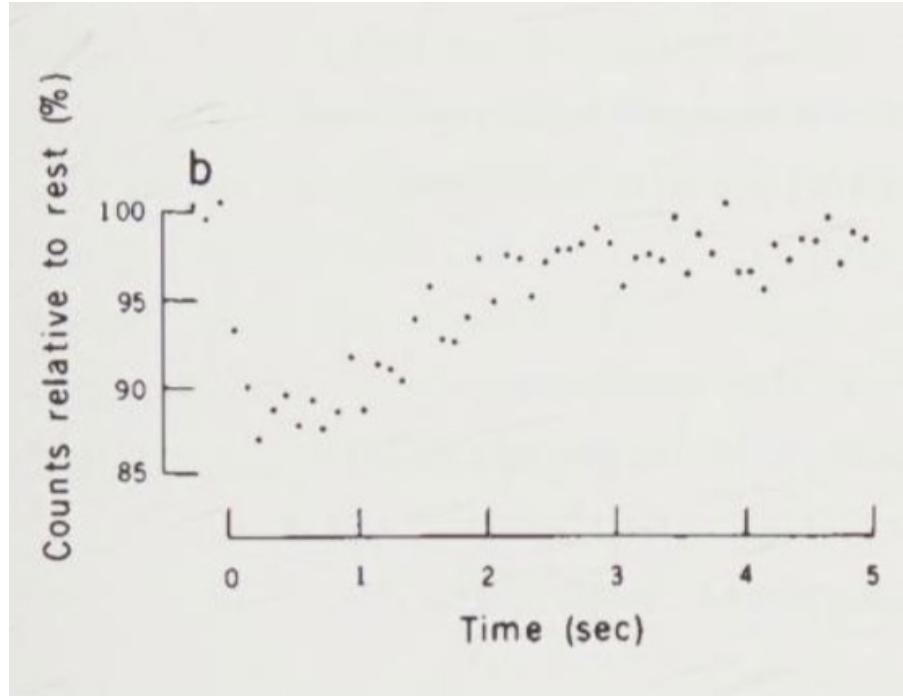
Use a pair of scale lines for each variable; make the data rectangle slightly smaller than the scale-line rectangle; tick marks should point outward



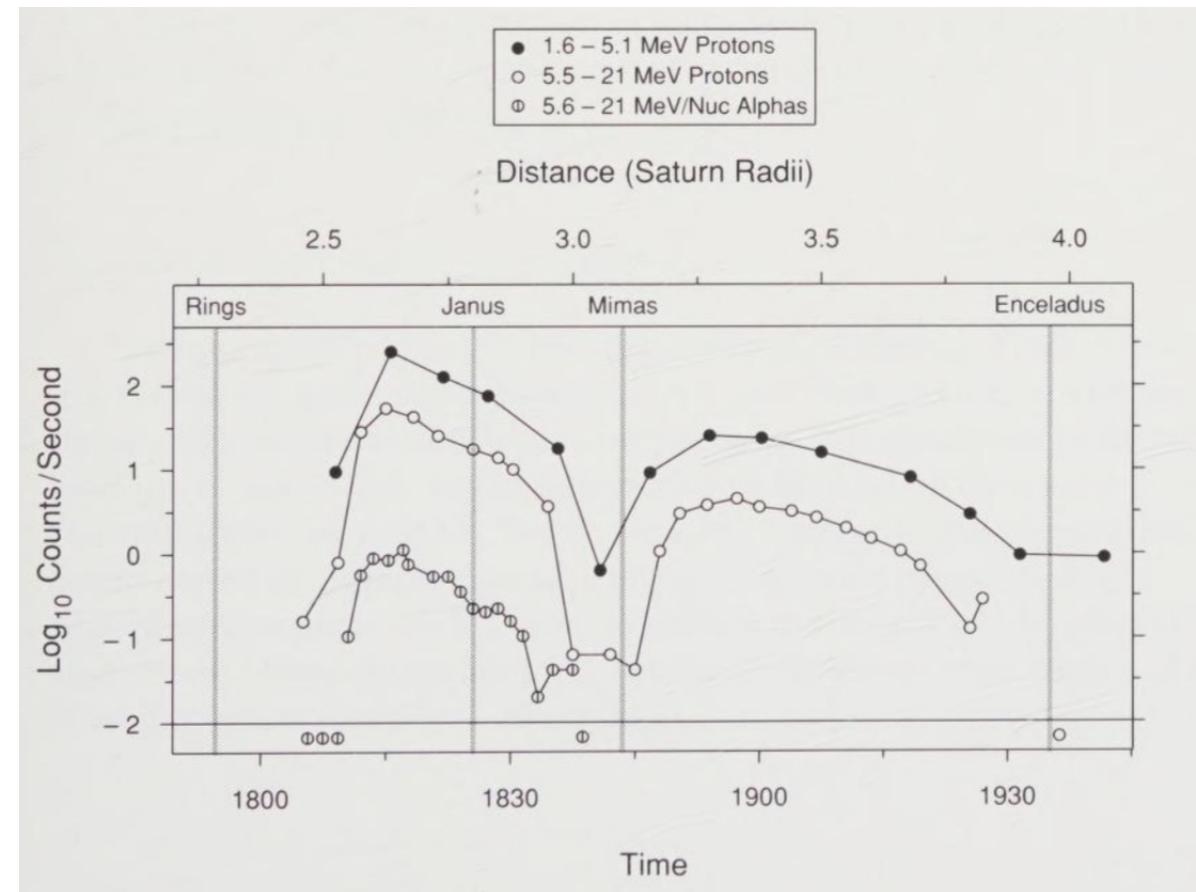
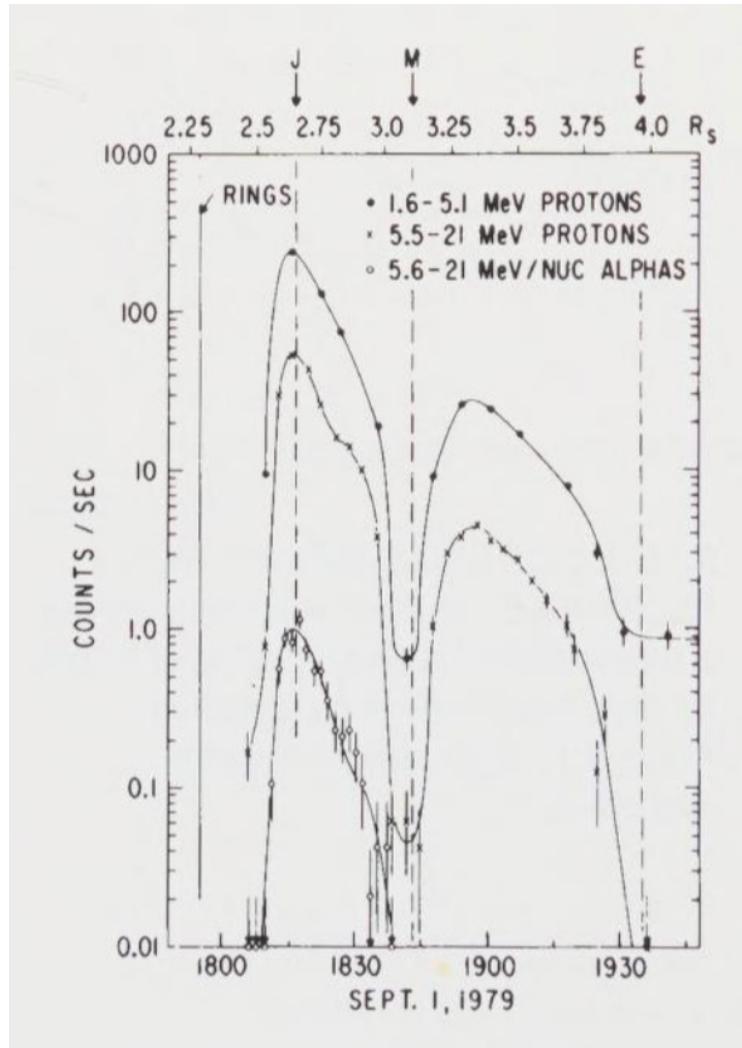
Use a pair of scale lines for each variable; make the data rectangle slightly smaller than the scale-line rectangle; tick marks should point outward



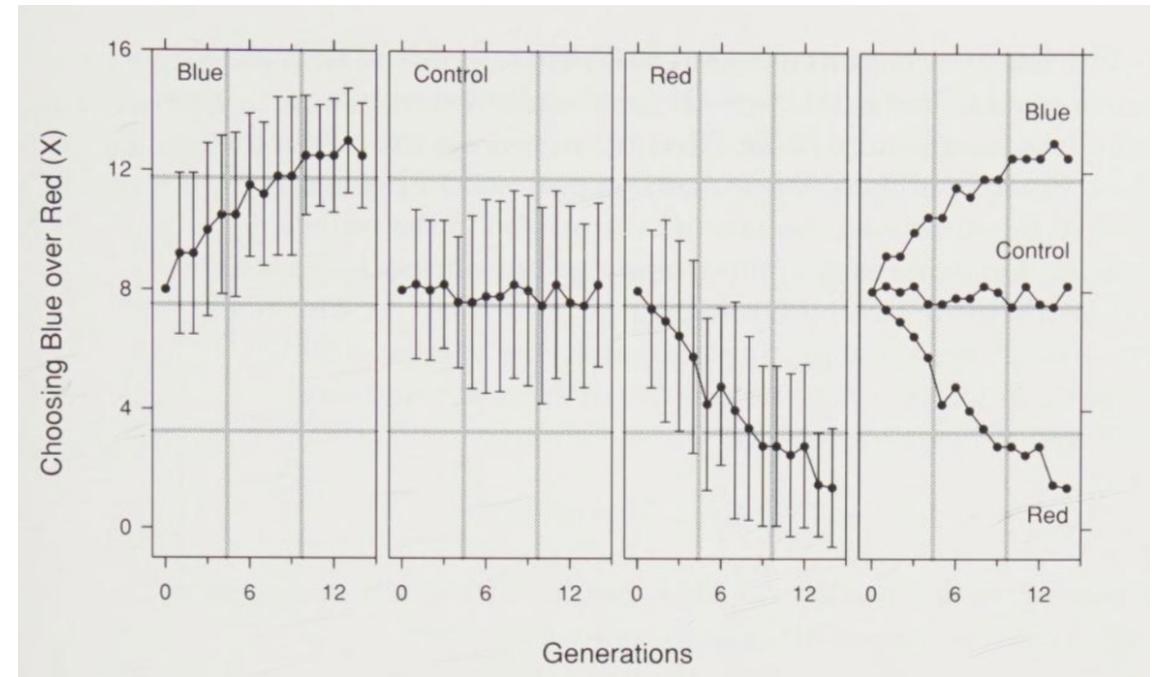
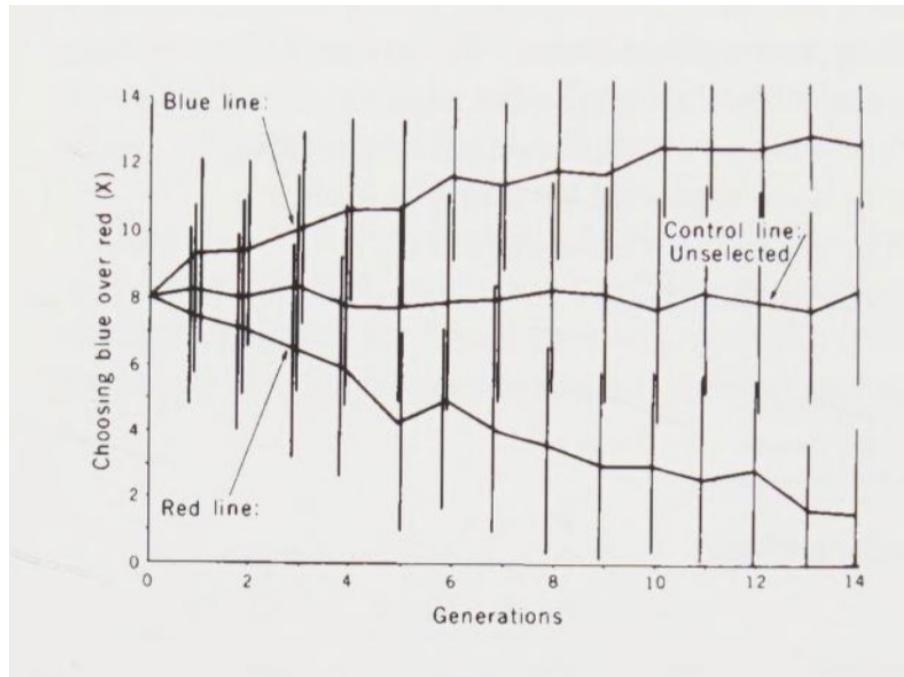
Use a pair of scale lines for each variable; make the data rectangle slightly smaller than the scale-line rectangle; tick marks should point outward



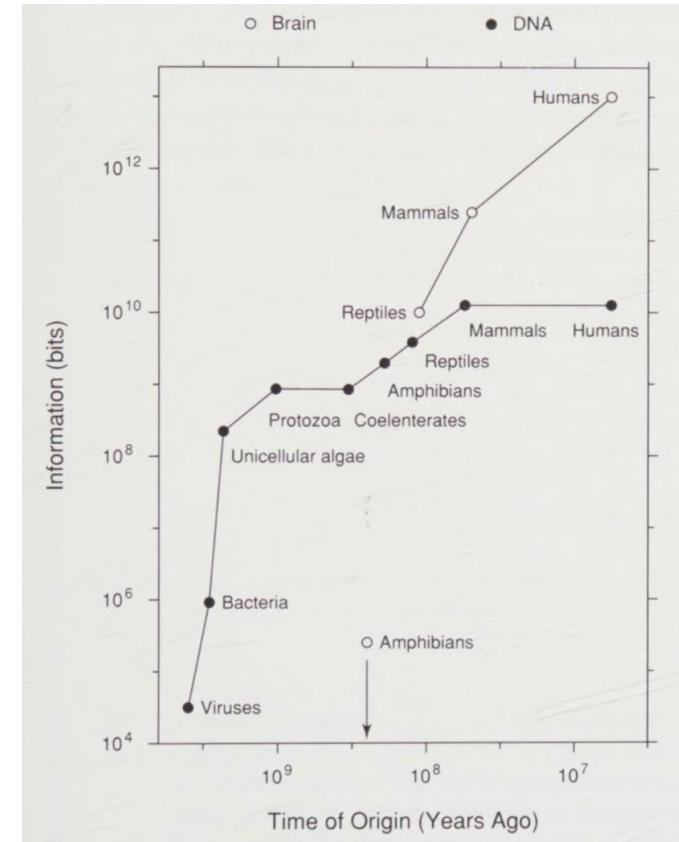
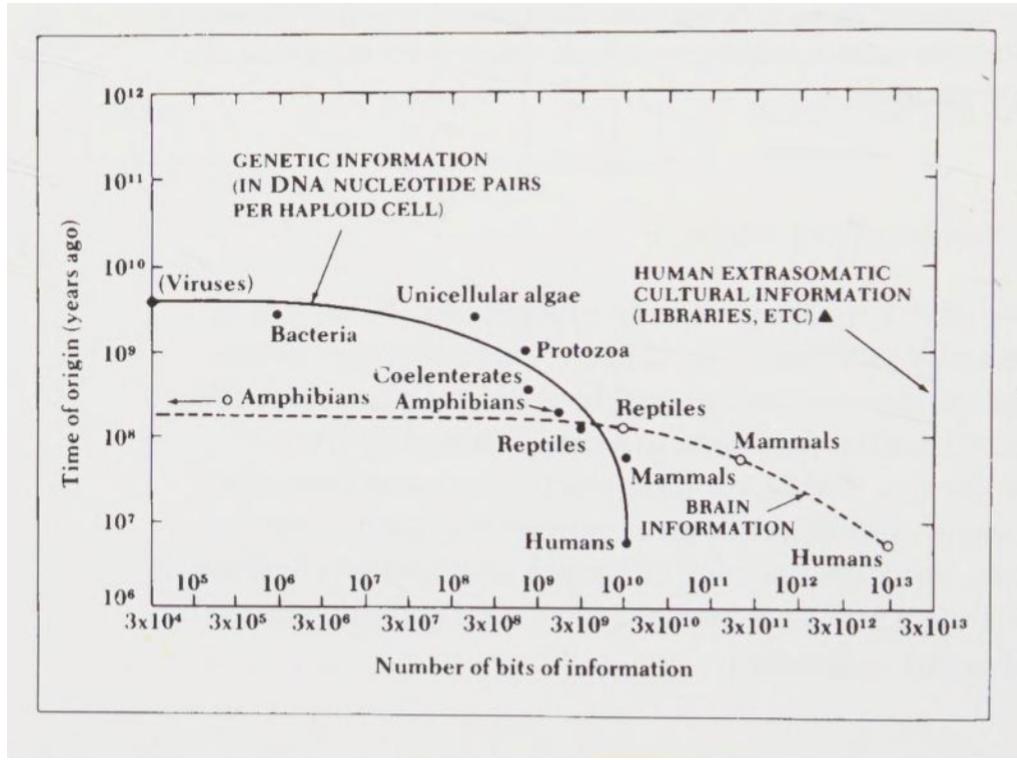
Do not clutter the interior of the scale-line rectangle



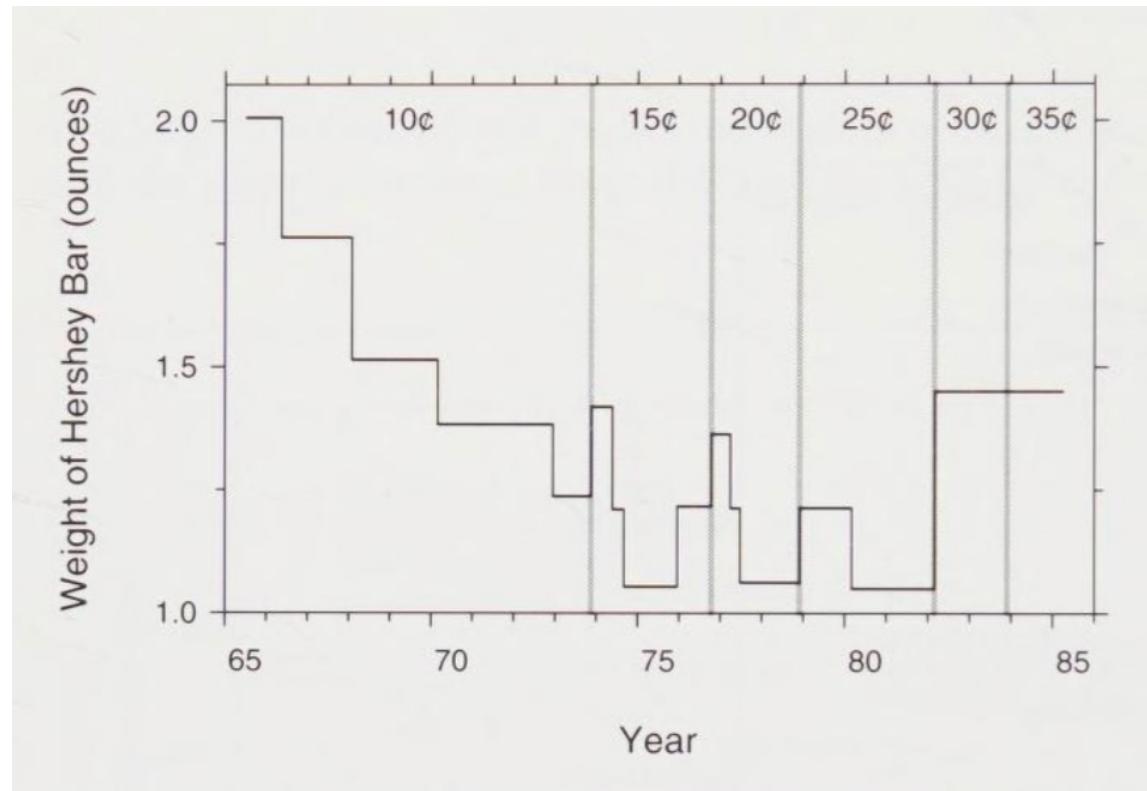
Do not clutter the interior of the scale-line rectangle



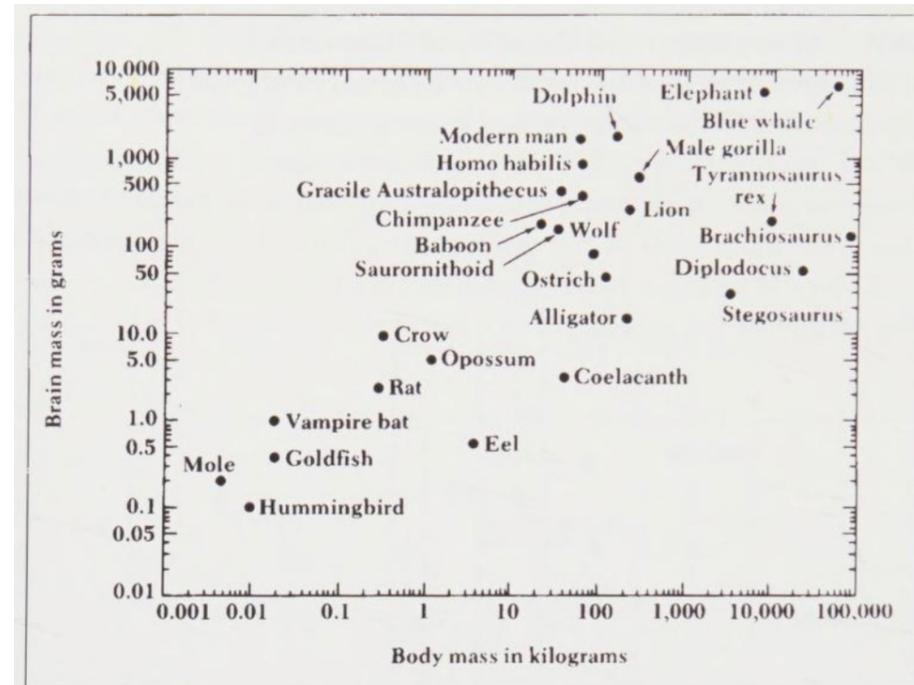
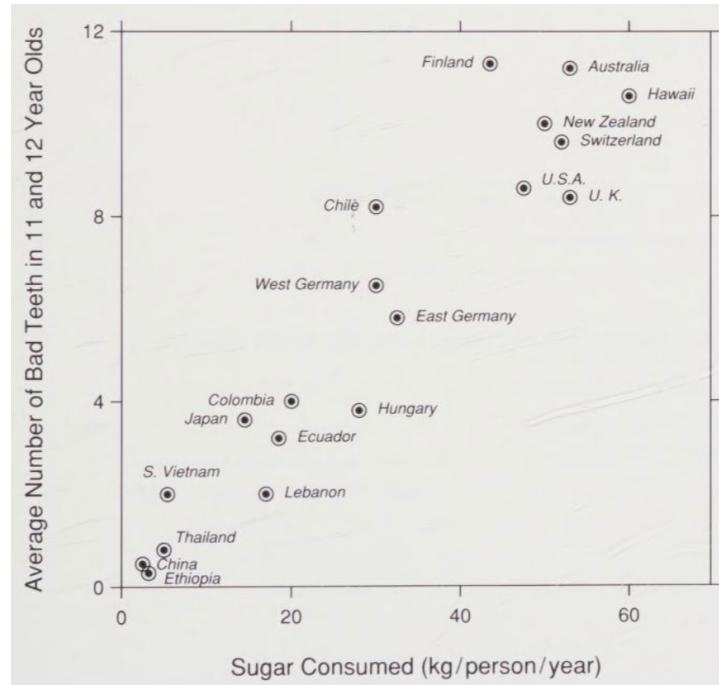
Do not overdo the number of tick marks



Use a reference line when there is an important value that must be seen across the entire graph, but do not let the line interfere with the graph

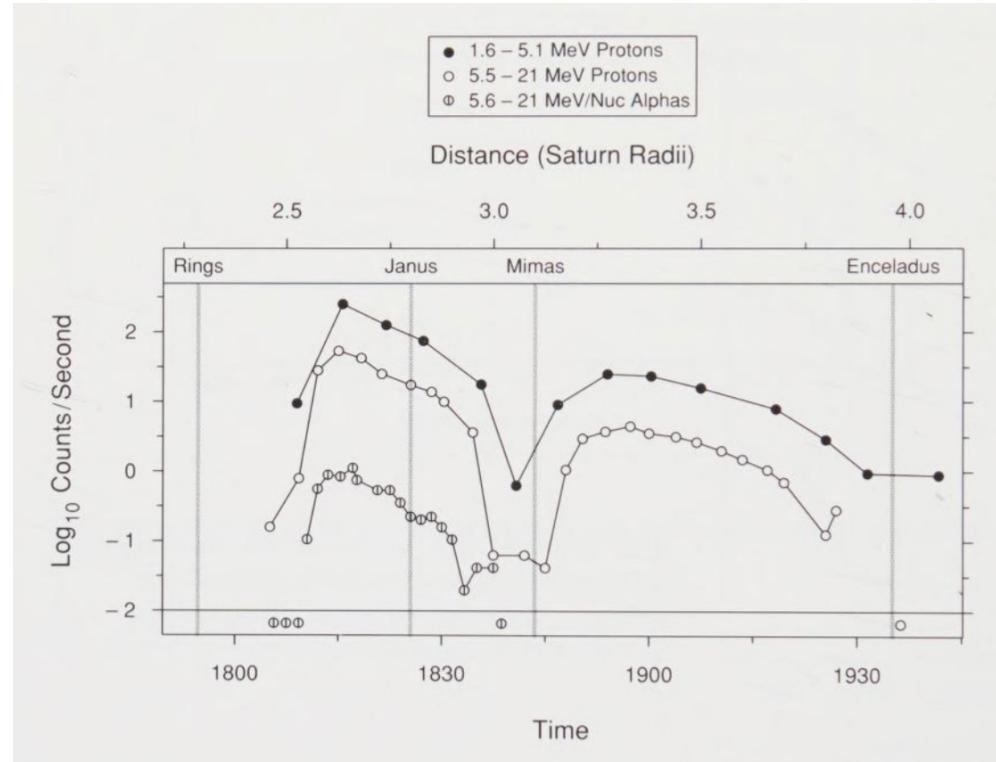
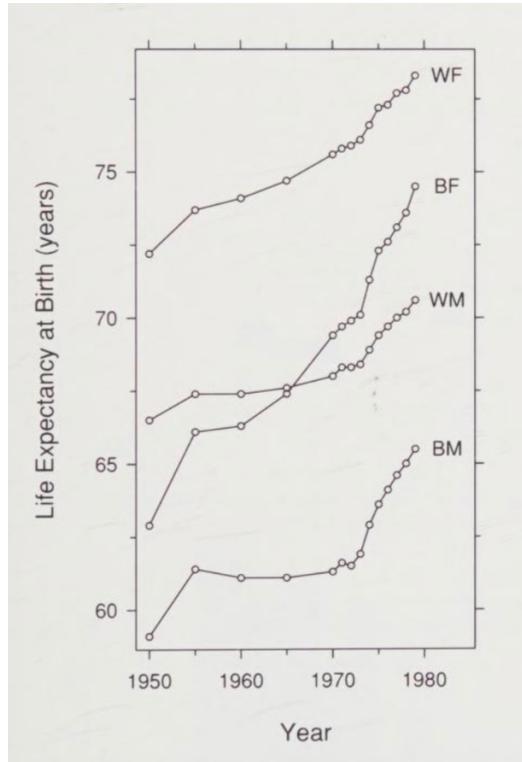


Do not allow labels in the interior of the scale-line rectangle to interfere with the quantitative data or to clutter the graph



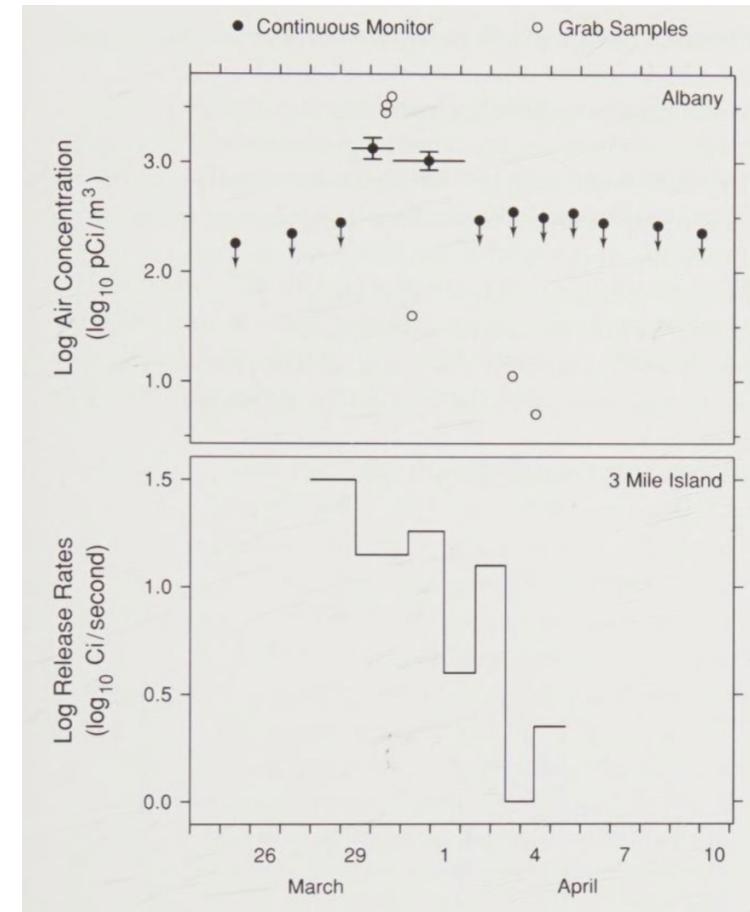
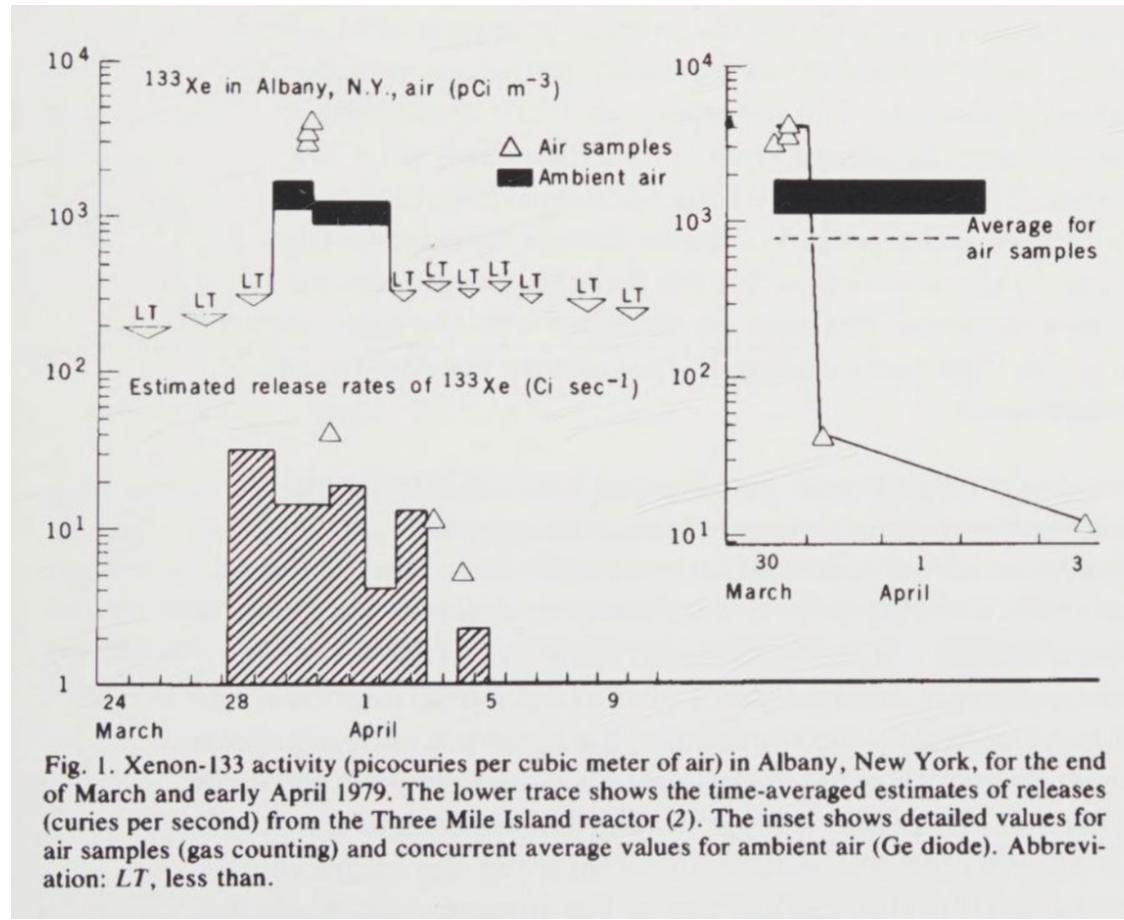
- External key makes identification harder but it reduces clutter

Do not allow labels in the interior of the scale-line rectangle to interfere with the quantitative data or to clutter the graph

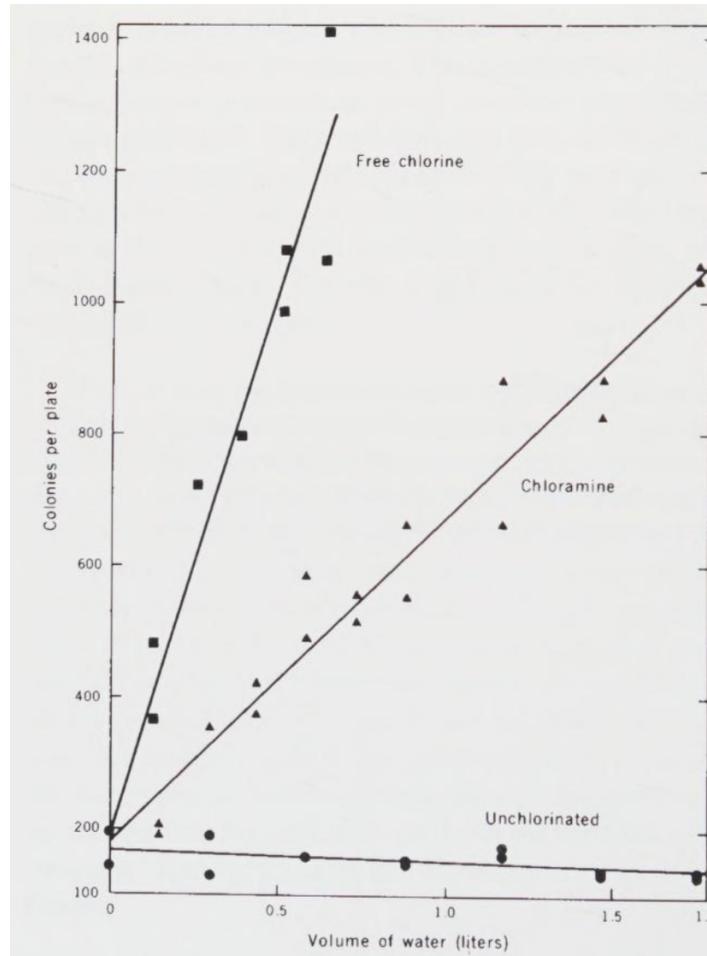


- External key makes identification harder but it reduces clutter

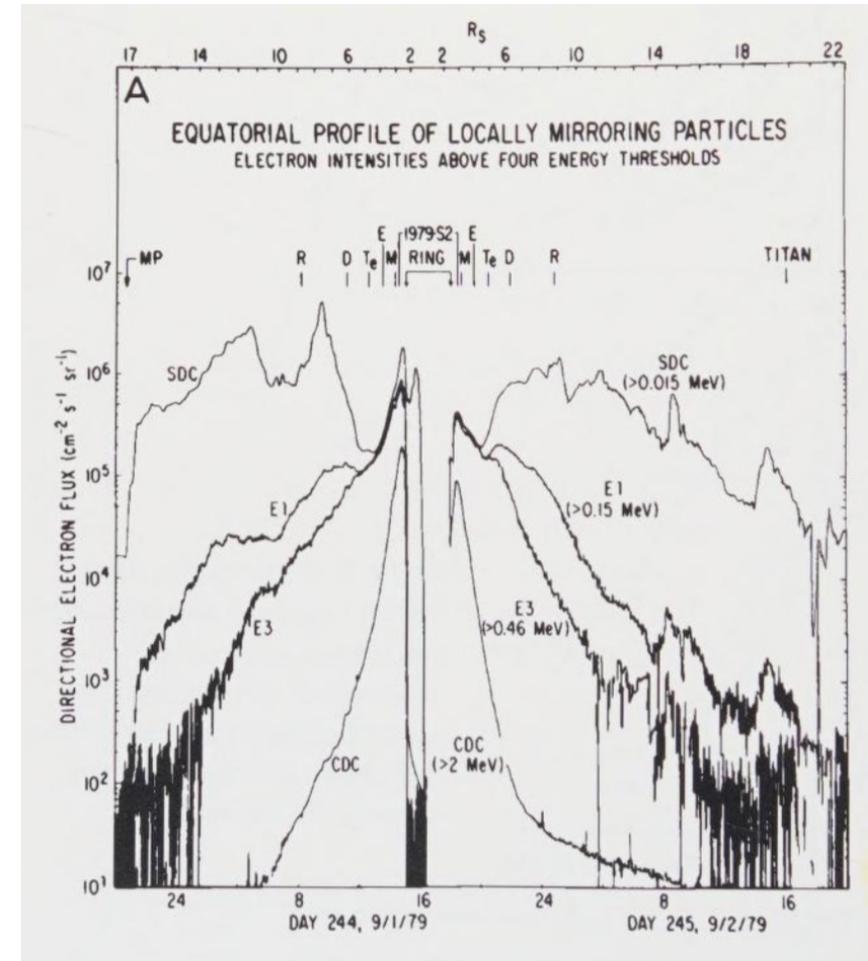
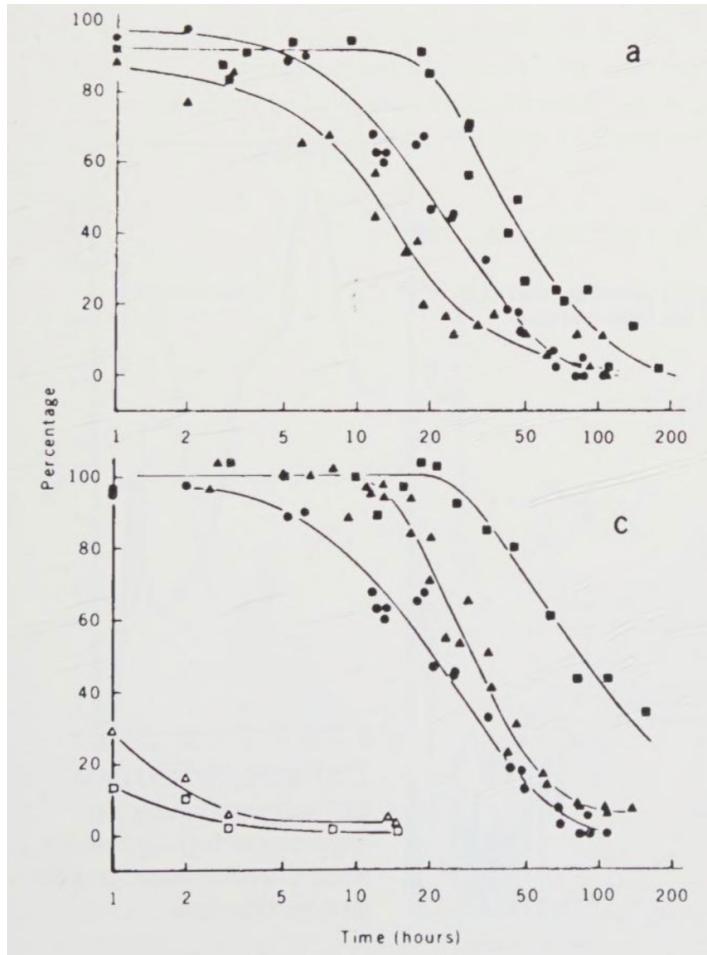
Avoid putting notes and keys inside the scale-line rectangle; put a key outside, and put notes in the caption or in the text



Overlapping plotting symbols must be visually distinguishable

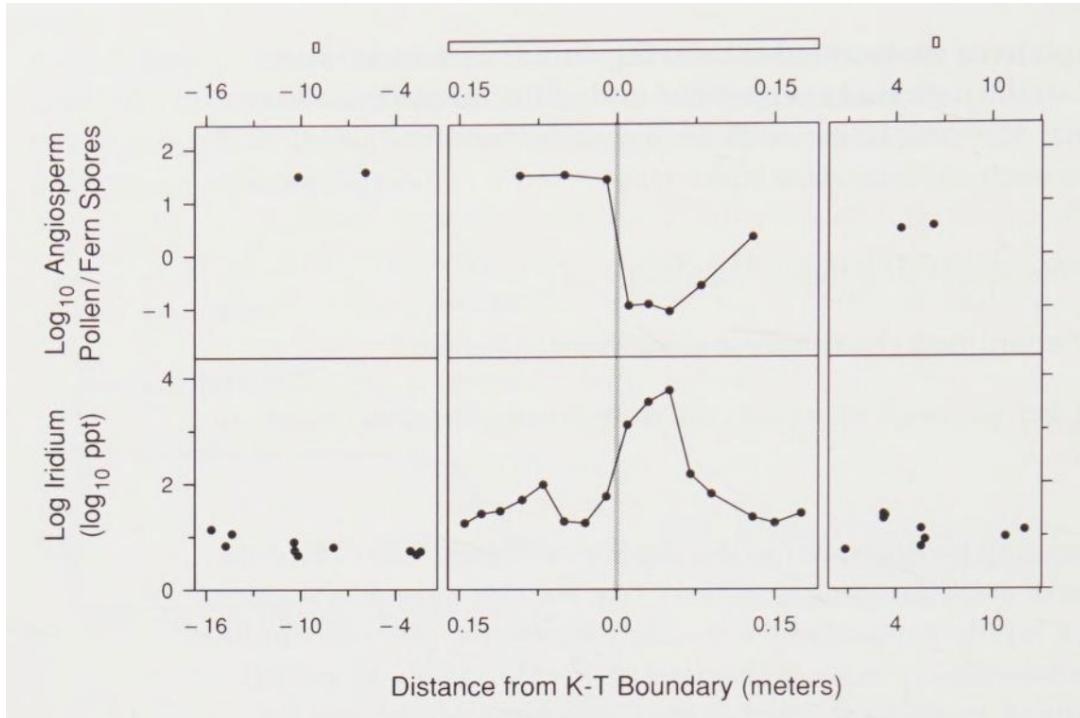


Superposed data sets must be readily visually assembled



Clear Understanding

Put major conclusions into graphical form; make captions comprehensive and informative

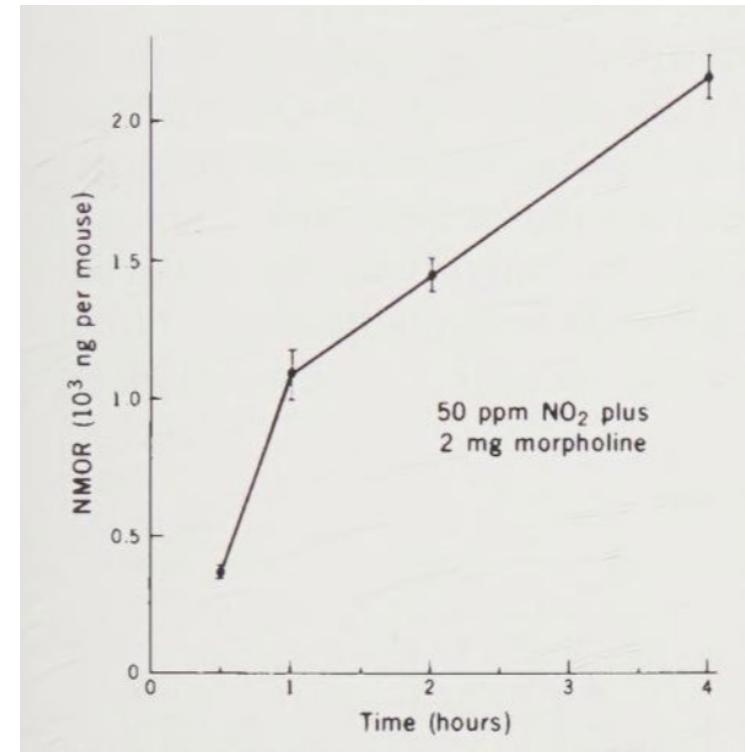
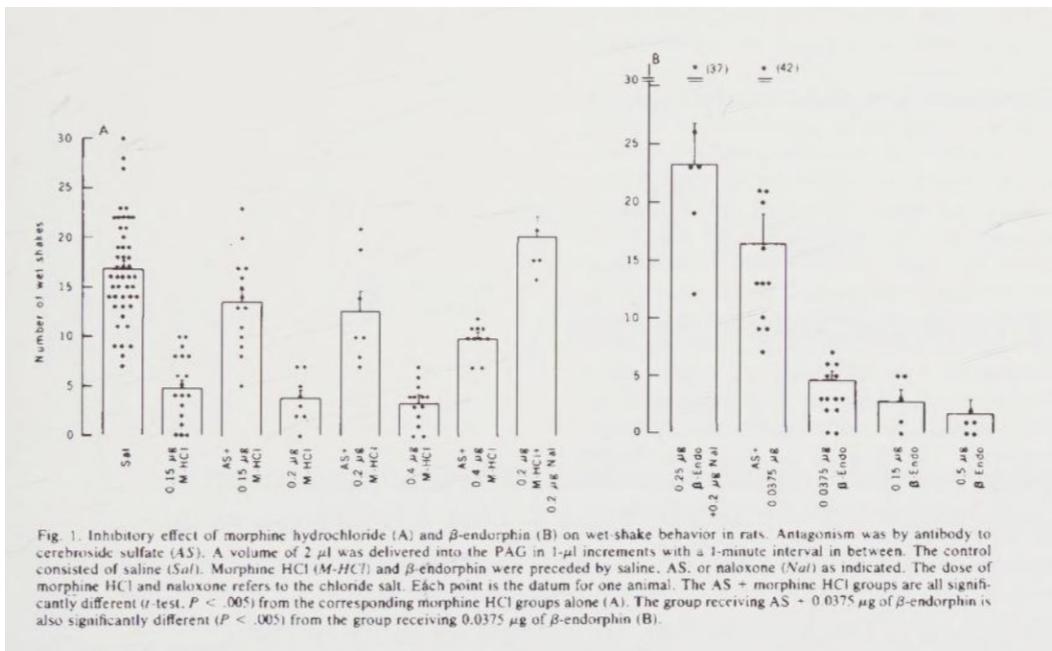


- Describe everything that is graphed
- Draw attention to the important features of the data
- Describe the conclusions that are drawn from the data on the graph

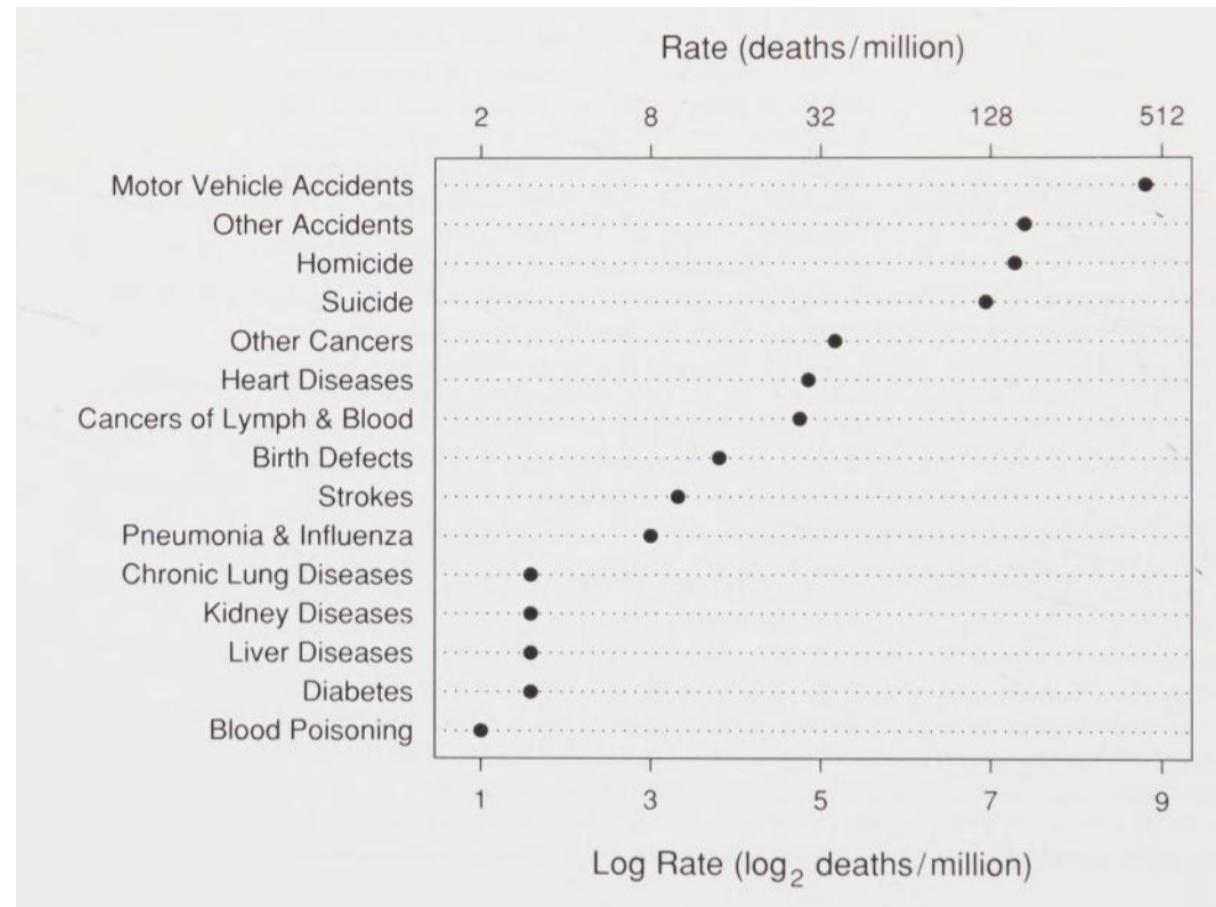
"Angiosperm-Fern Ratio and Iridium Near the K-T Boundary. The graph shows measurements of a core from northeastern New Mexico. The horizontal scale is in meters from the boundary between the Cretaceous and the Tertiary periods; negative values are below the K-T boundary so time goes from earlier to later in going from left to right. The widths of the three rectangles at the top of the graph show the same number of meters on the horizontal scales of the three panels. The top panel shows the ratio of angiosperm pollen to fern spores; the K-T boundary is taken to be the time point at which these values begin to decrease. The bottom panel shows concentrations of iridium; the concentrations begin a dramatic rise and fall at the boundary. Since the principal source of iridium is extraterrestrial, its rise and fall supports the hypothesis that an asteroid struck the earth causing a cloud of dust in the upper atmosphere; this is argued to have caused the large number of extinctions, including the dinosaurs, that occurred at the beginning of the Tertiary period."

If you have error bars, make sure you explain

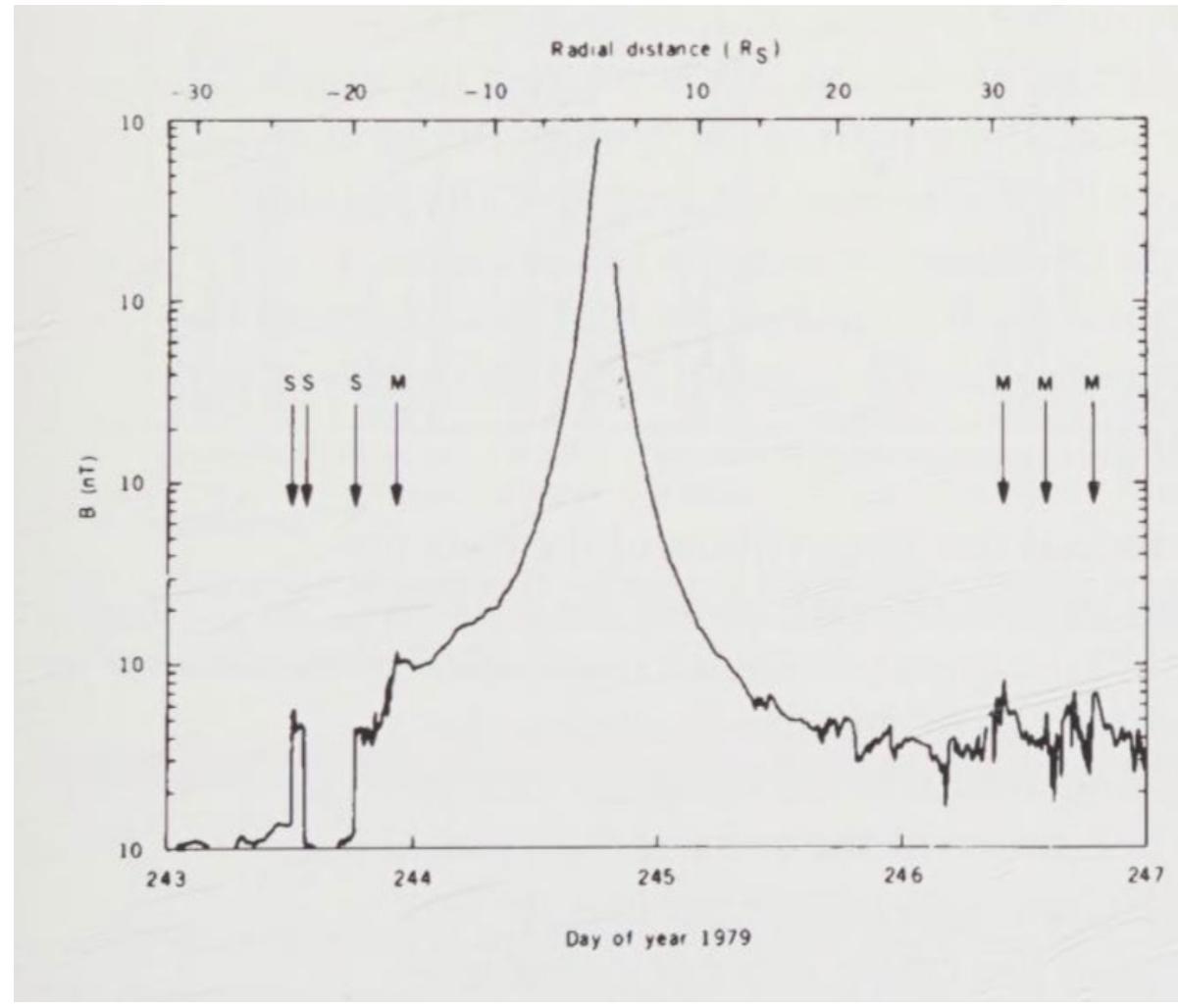
- What “error” are they?
 - Sample standard deviation
 - Standard error (estimate of the standard deviation)
 - Confidence interval



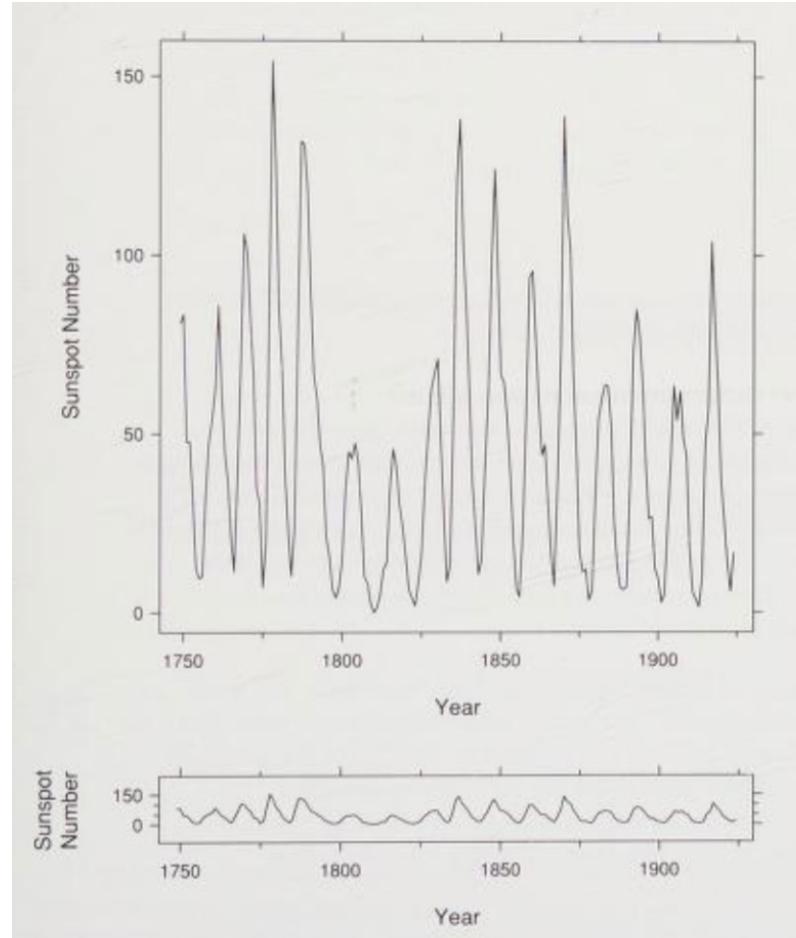
When logarithms of a variable are graphed, the scale label should correspond to the tick mark labels



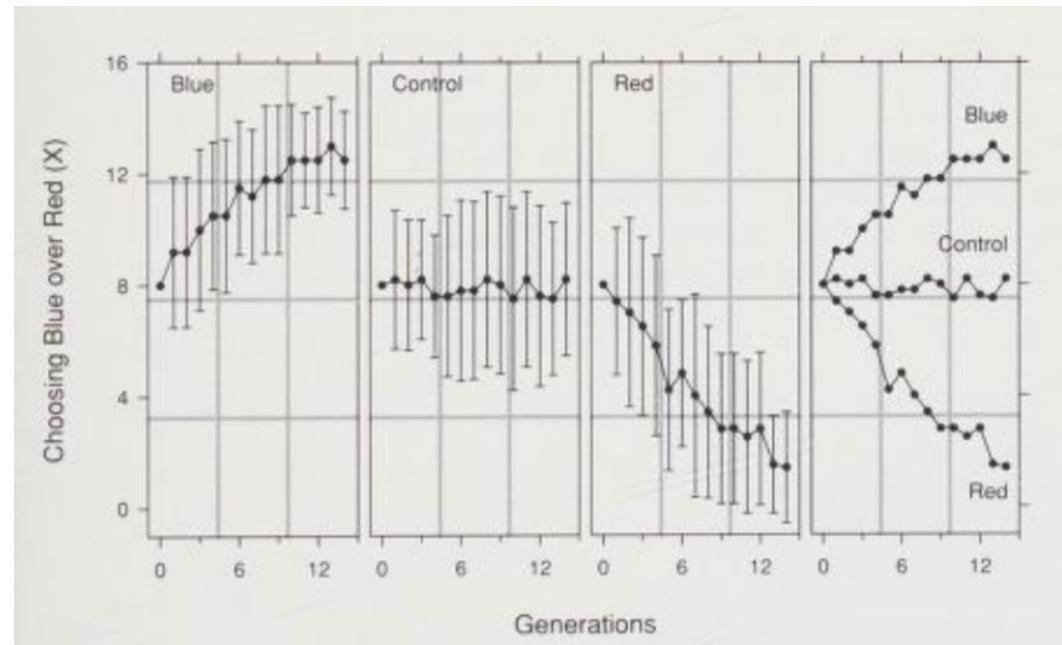
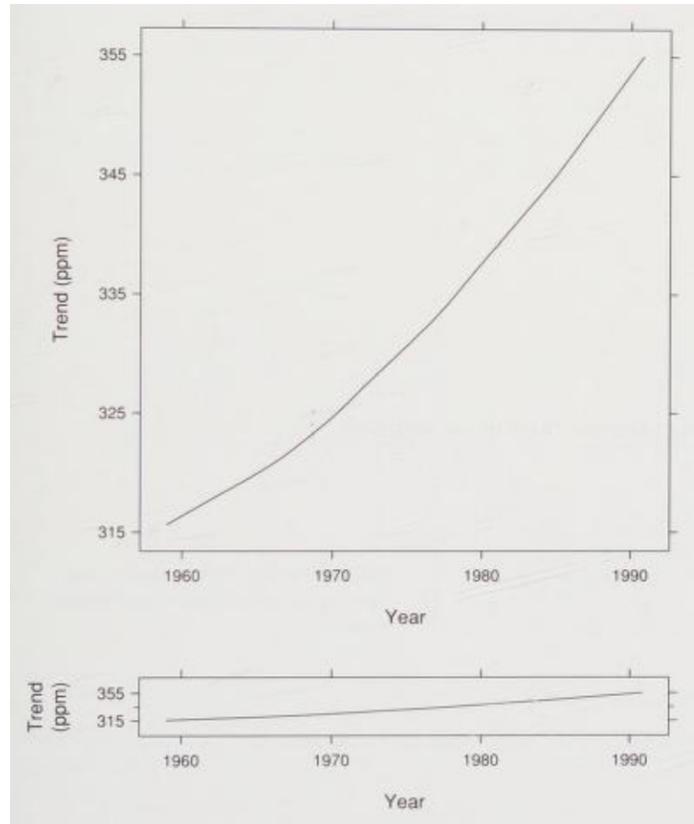
Proofread graphs



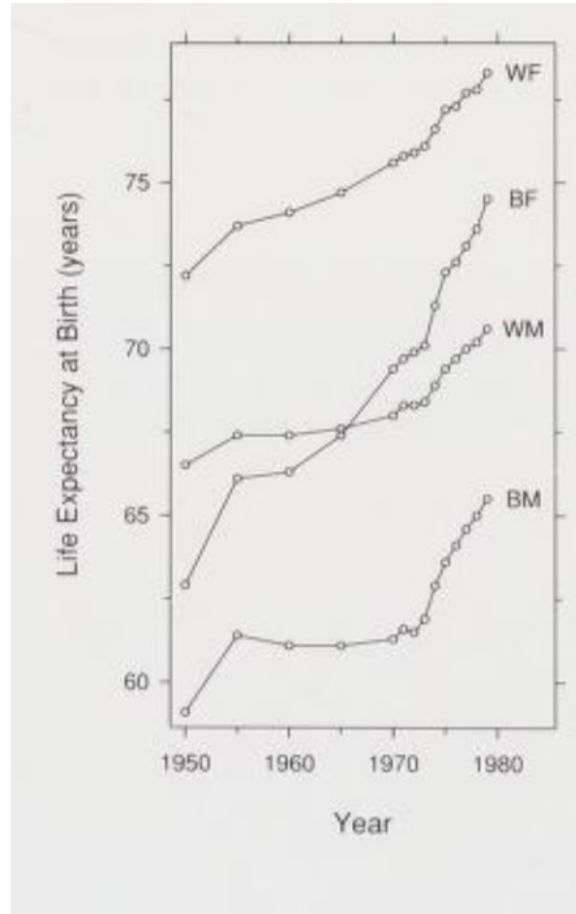
When the orientation of line segments are judged to decode information about rate of change, bank the segments to 45 degrees



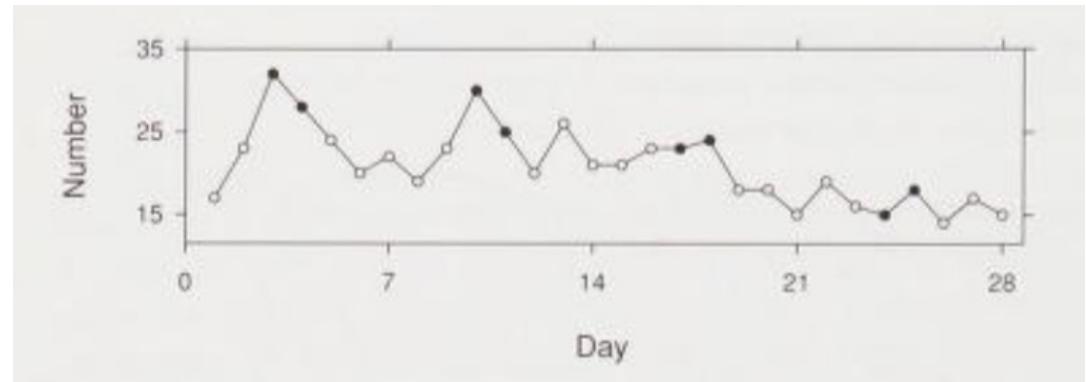
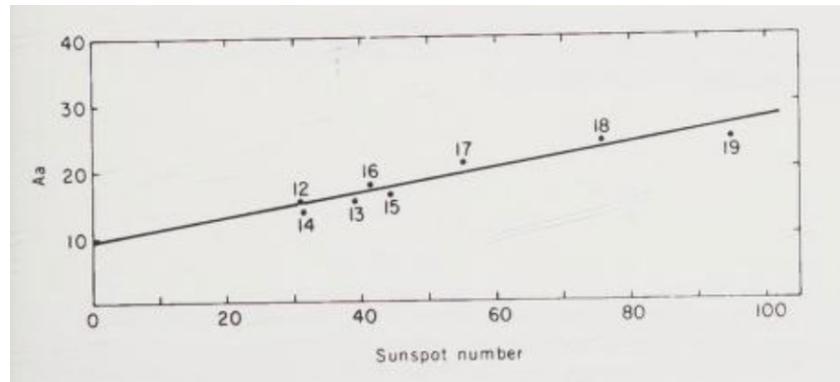
When the orientation of line segments are judged to decode information about rate of change, bank the segments to 45 degrees



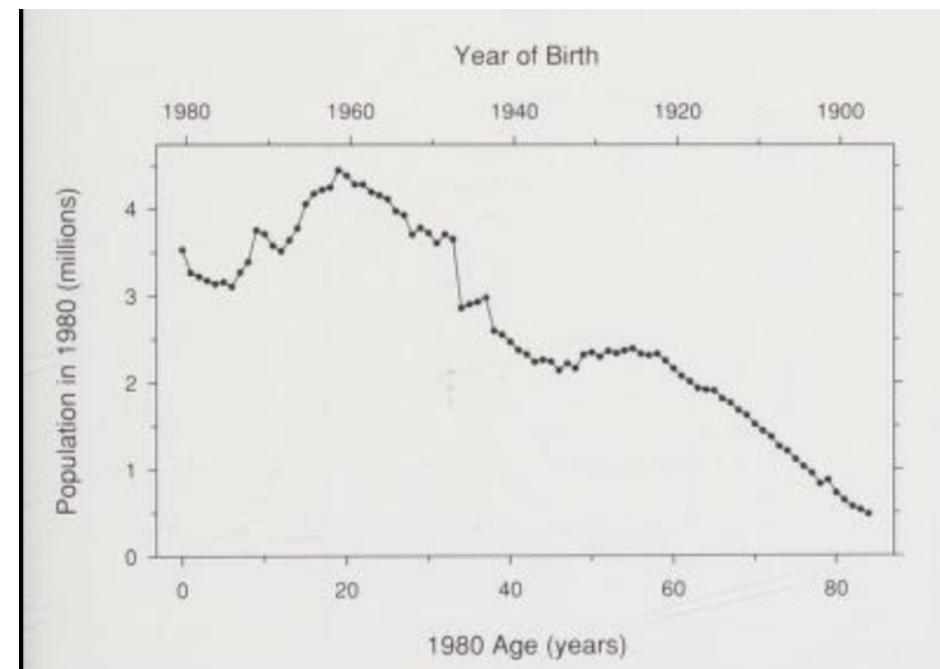
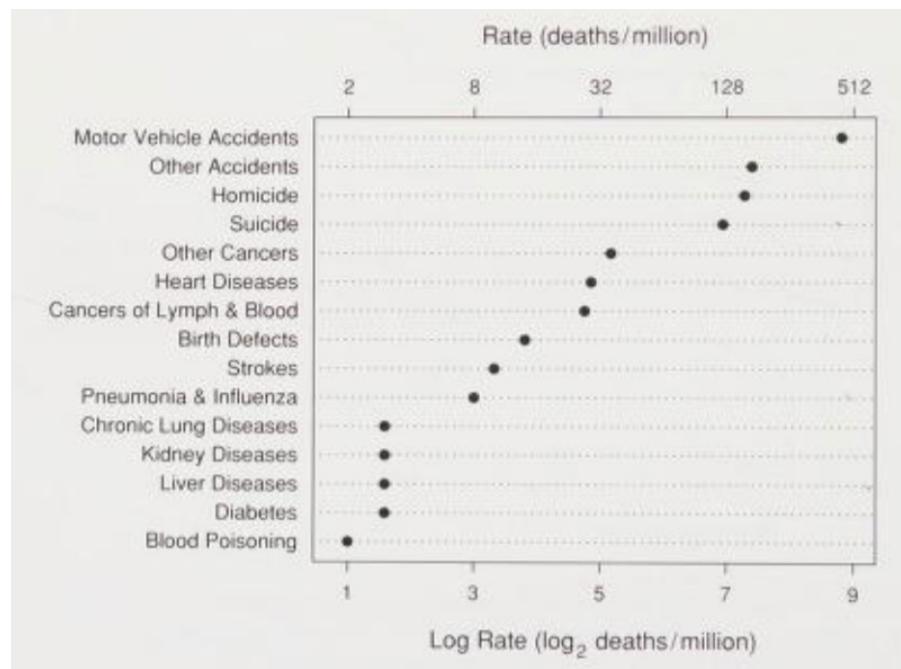
Choose the range of the tick marks to include or nearly include the range of the data



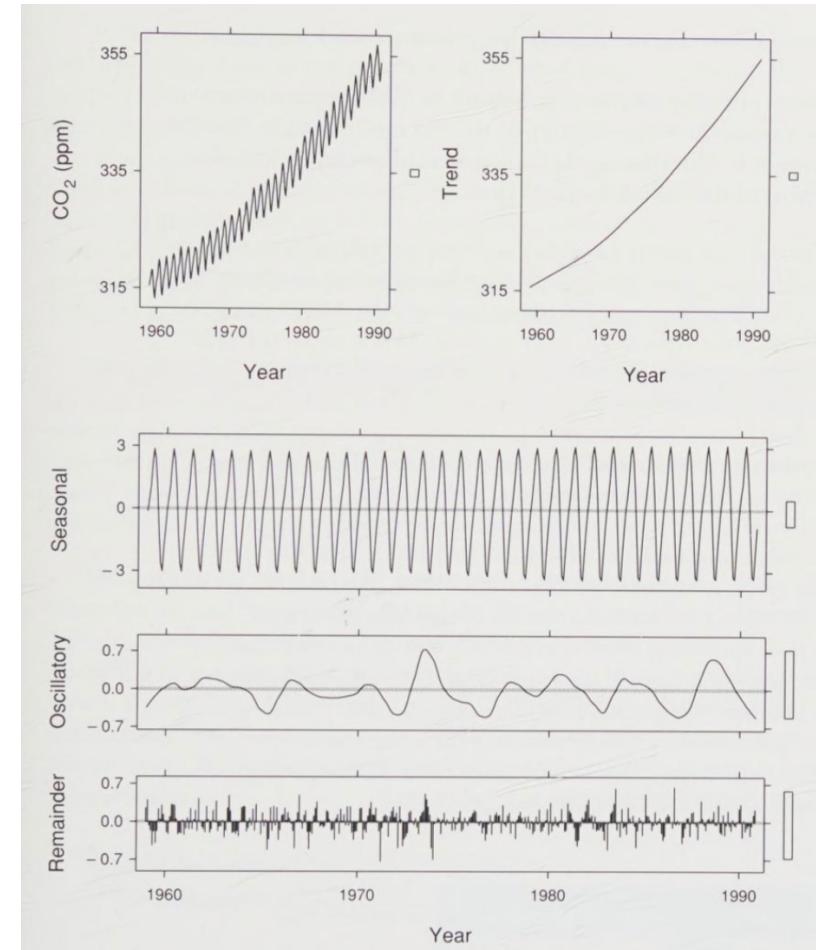
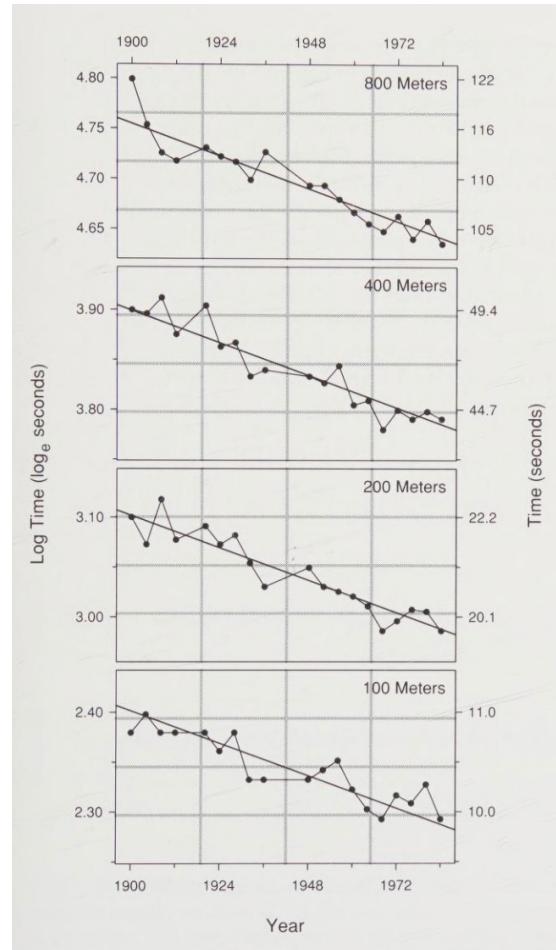
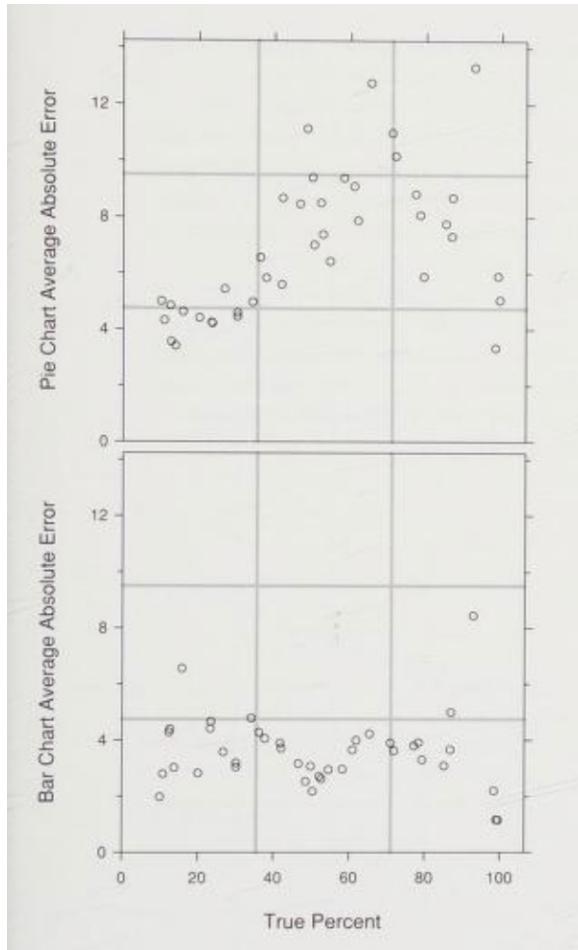
Subject to the constraints that scales have, choose the scales so that the data rectangle fills up as much of the scale-line rectangle as possible



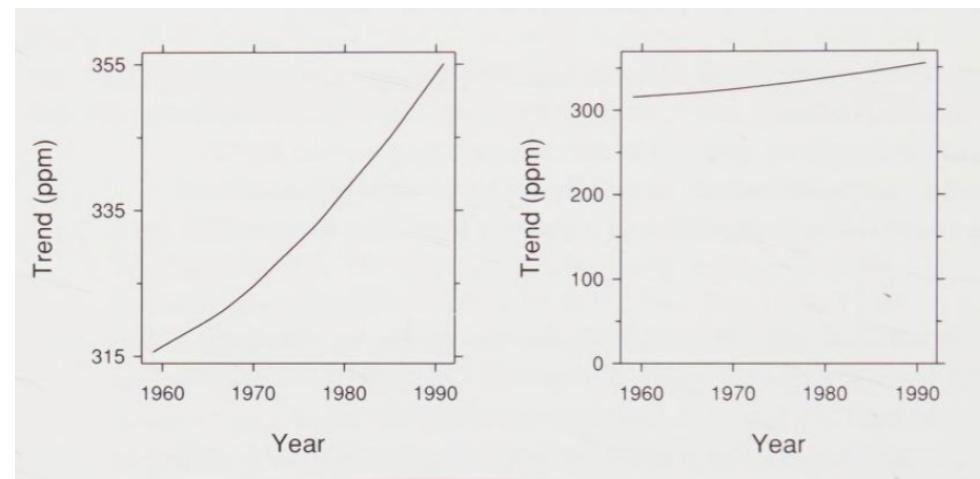
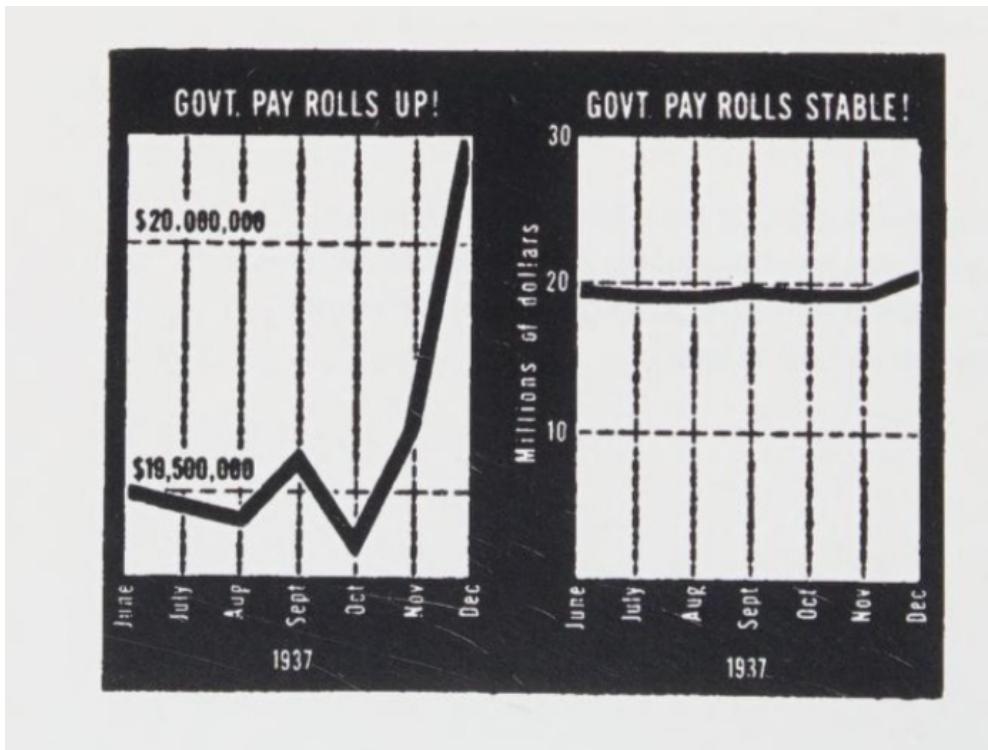
It is sometimes helpful to use the pair of scale lines for a variable to show two different scales



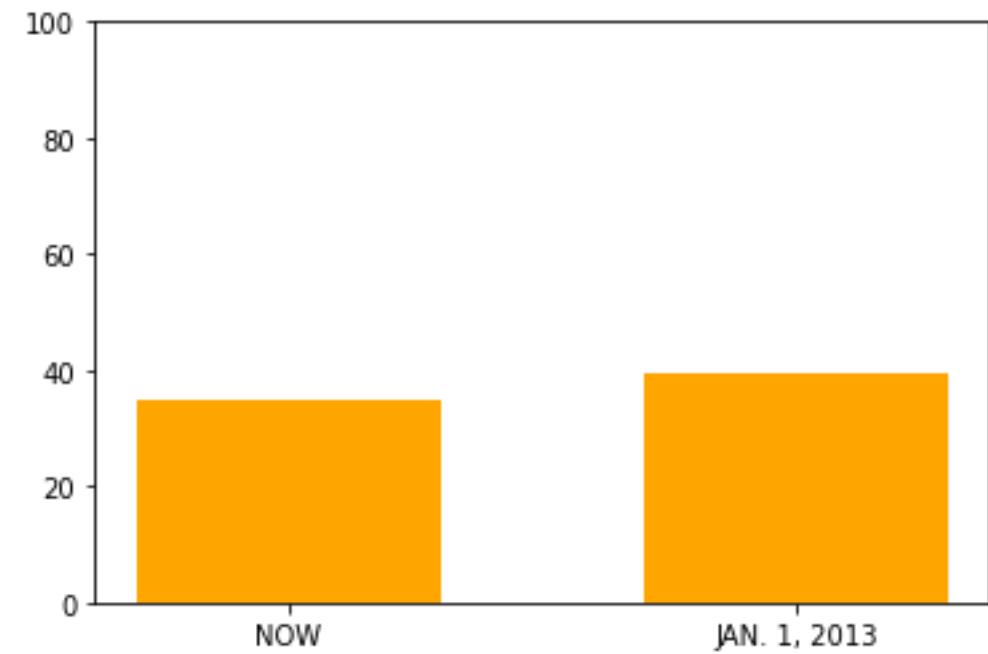
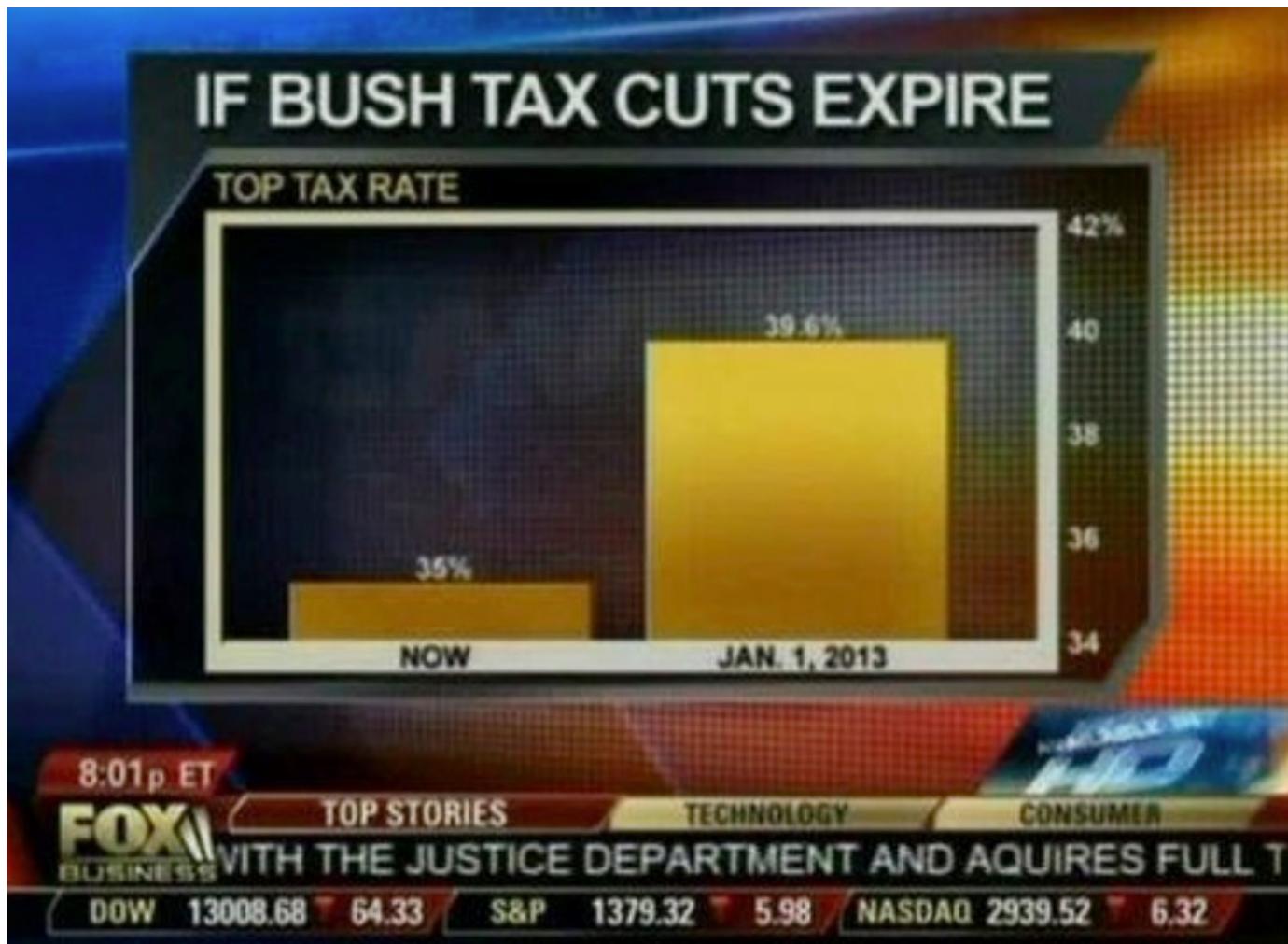
Choose appropriate scales when data on different panels are compared



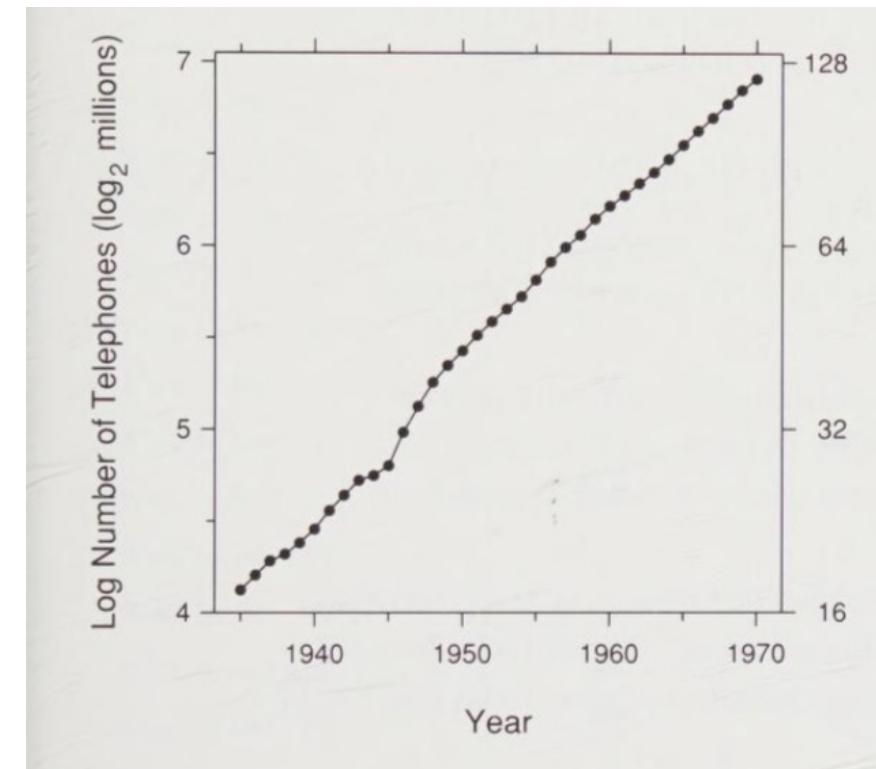
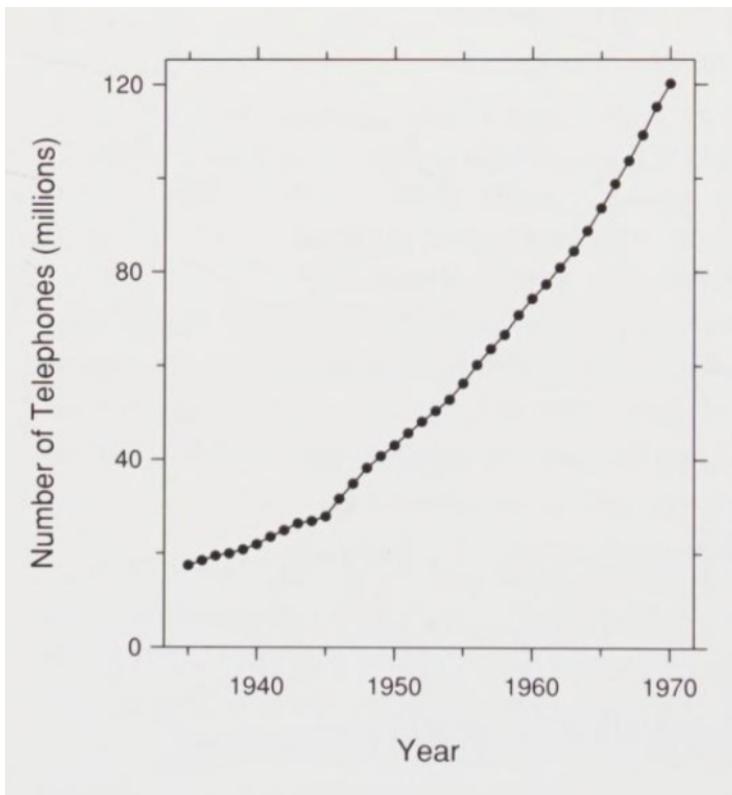
Do not insist that zero always be included on a scale showing magnitude



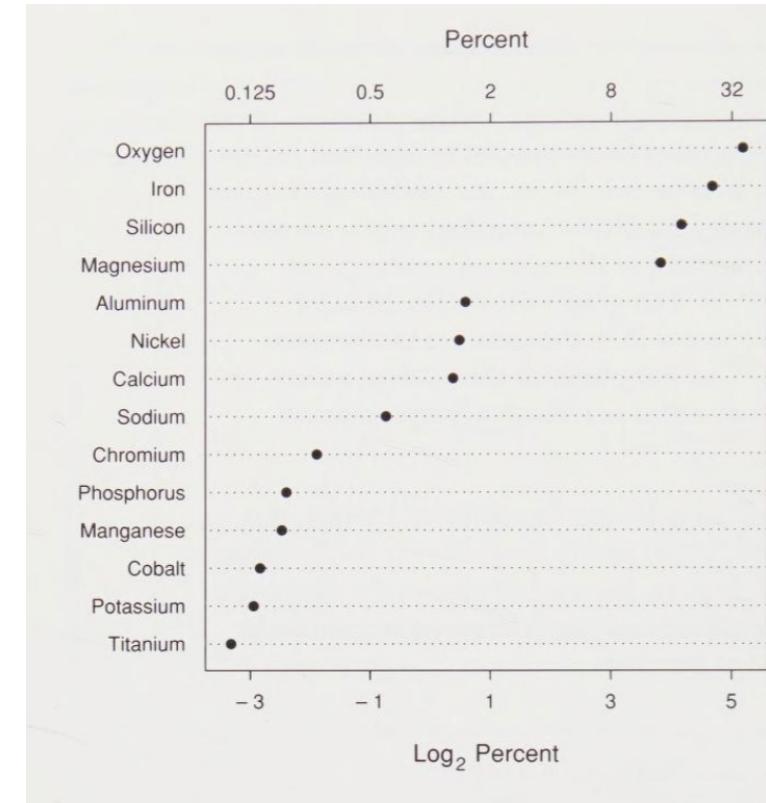
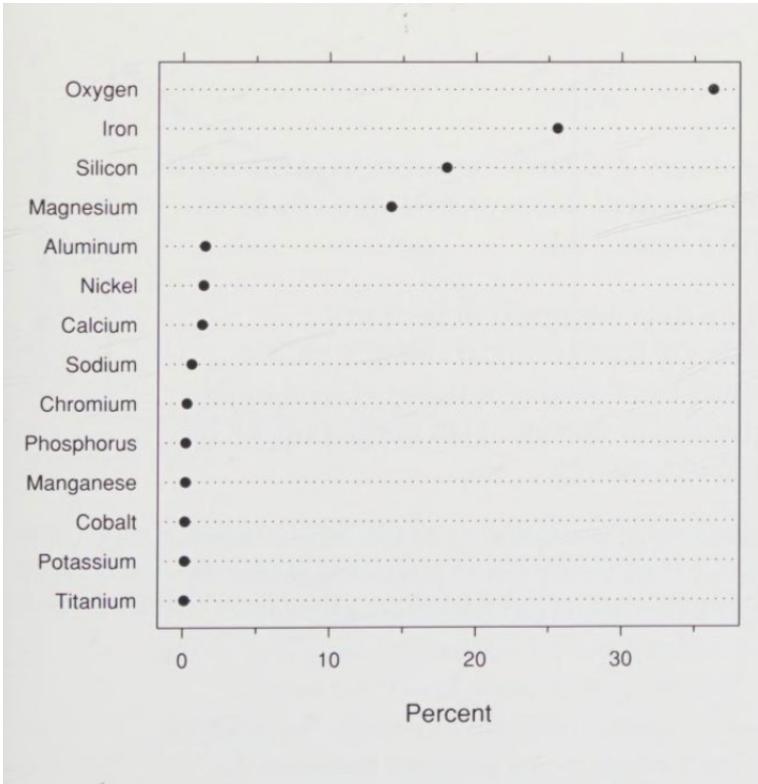
Is there a conflict with our rule here?



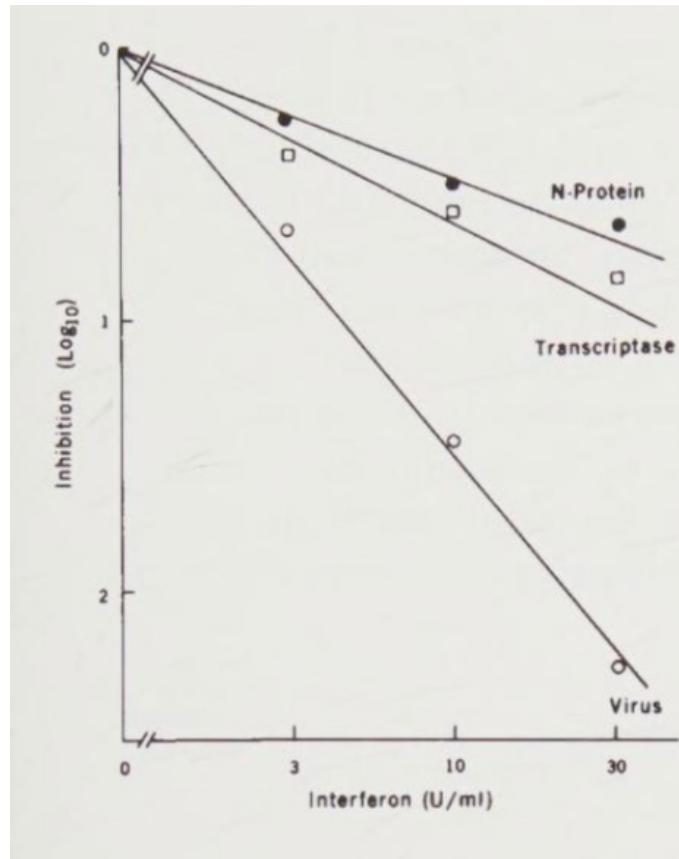
Use a logarithmic scale when it is important to understand percent change or multiplicative factors



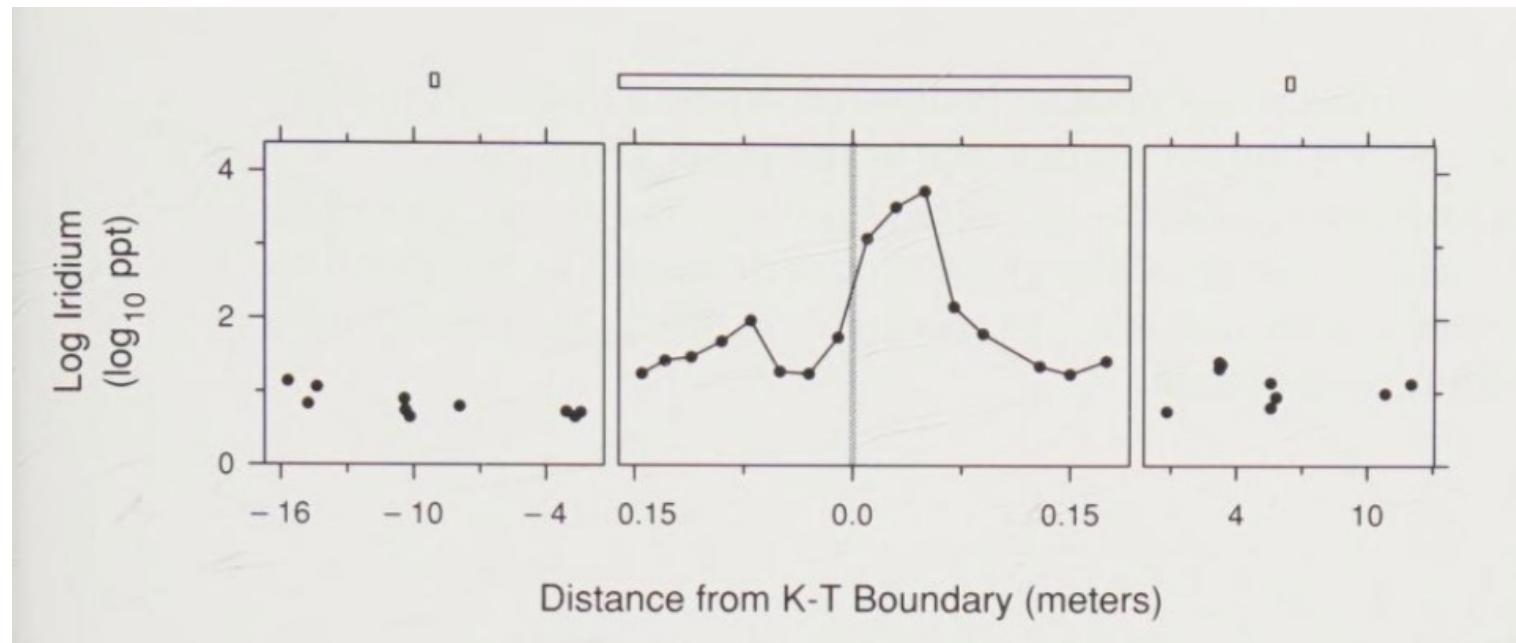
Showing data on a logarithmic scale can cure skewness toward large values



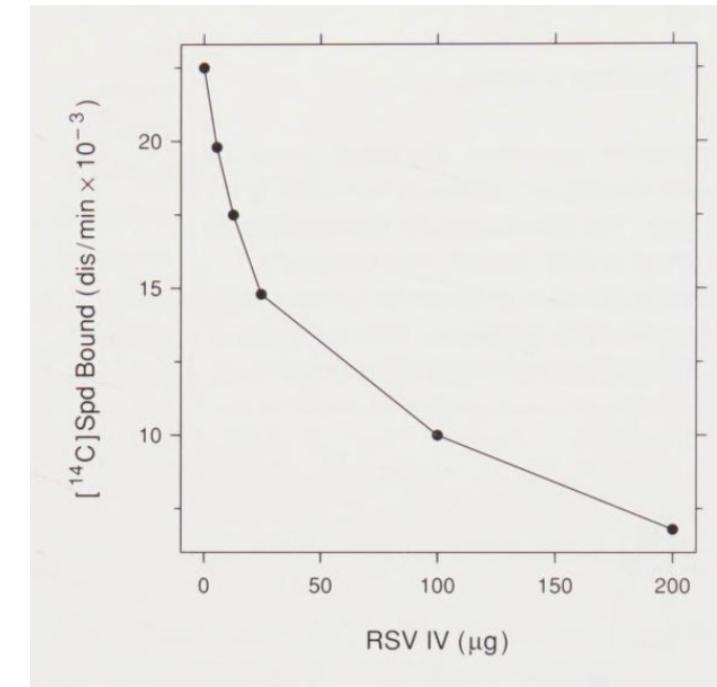
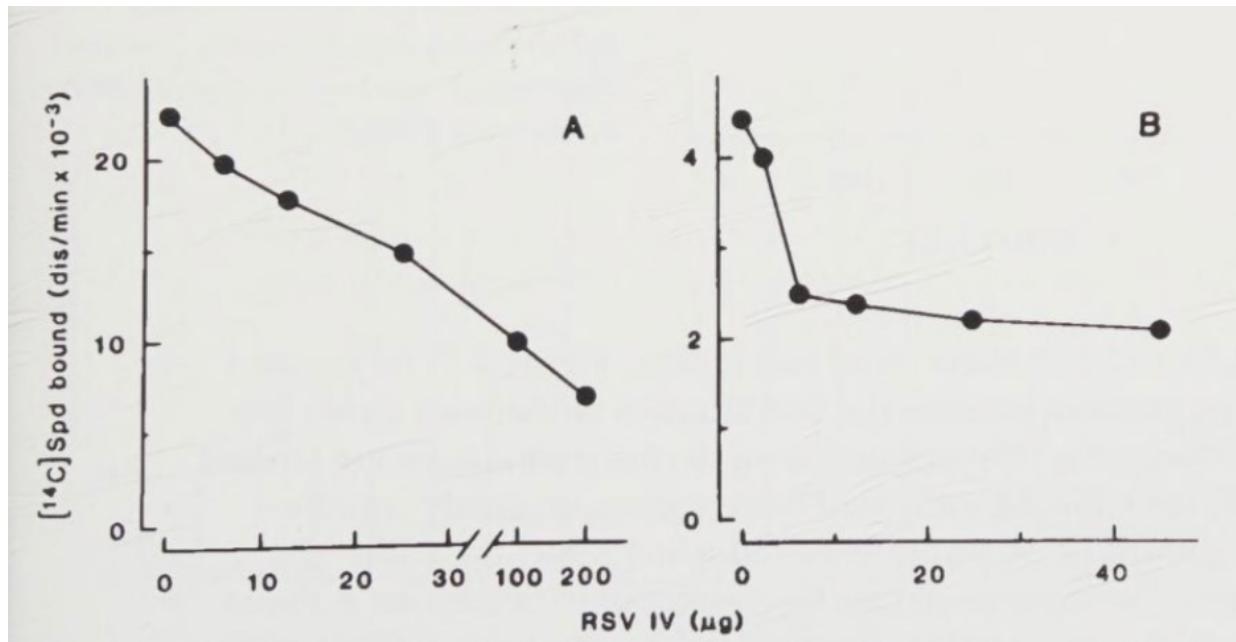
Use a scale break only when necessary; if a break cannot be avoided, use a full scale break; do not connect numerical values on two sides of a break; taking logs can cure the need for a break



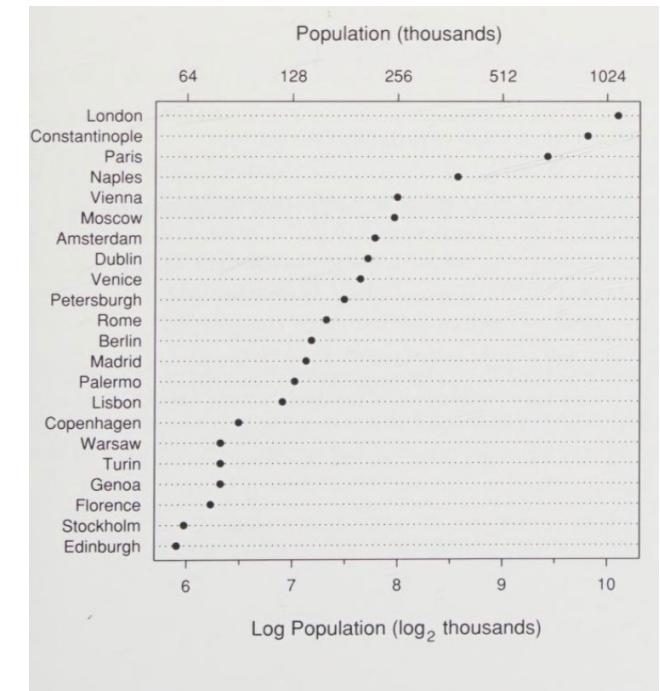
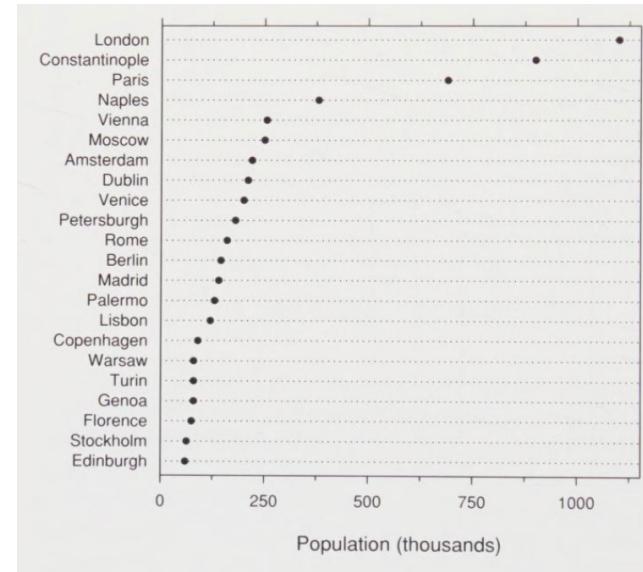
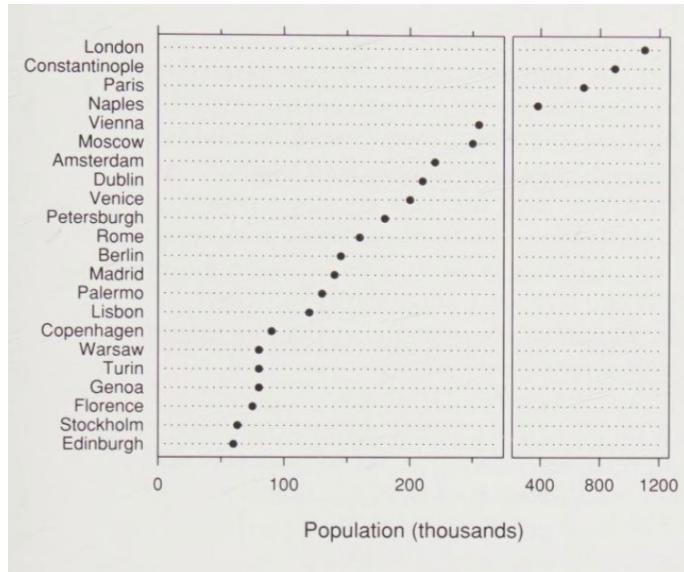
Use a scale break only when necessary; if a break cannot be avoided, use a full scale break; do not connect numerical values on two sides of a break; taking logs can cure the need for a break



Use a scale break only when necessary; if a break cannot be avoided, use a full scale break; do not connect numerical values on two sides of a break; taking logs can cure the need for a break



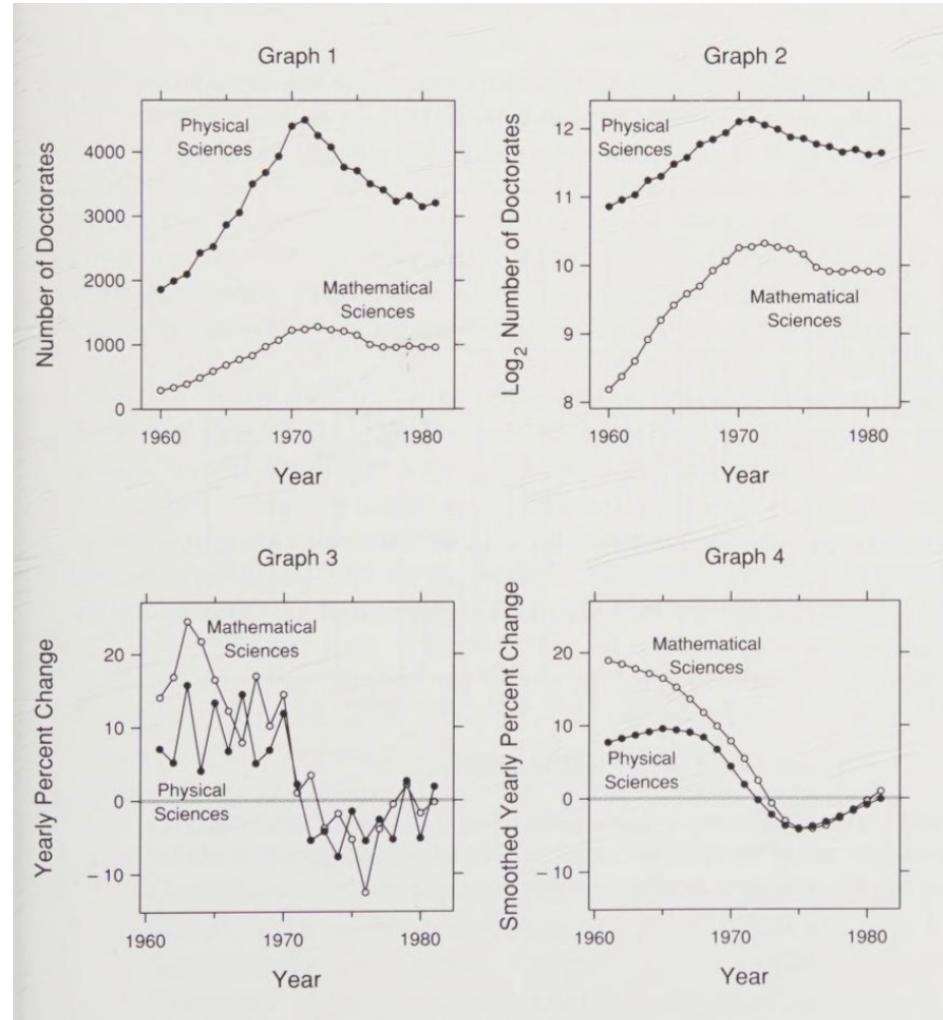
Use a scale break only when necessary; if a break cannot be avoided, use a full scale break; do not connect numerical values on two sides of a break; taking logs can cure the need for a break



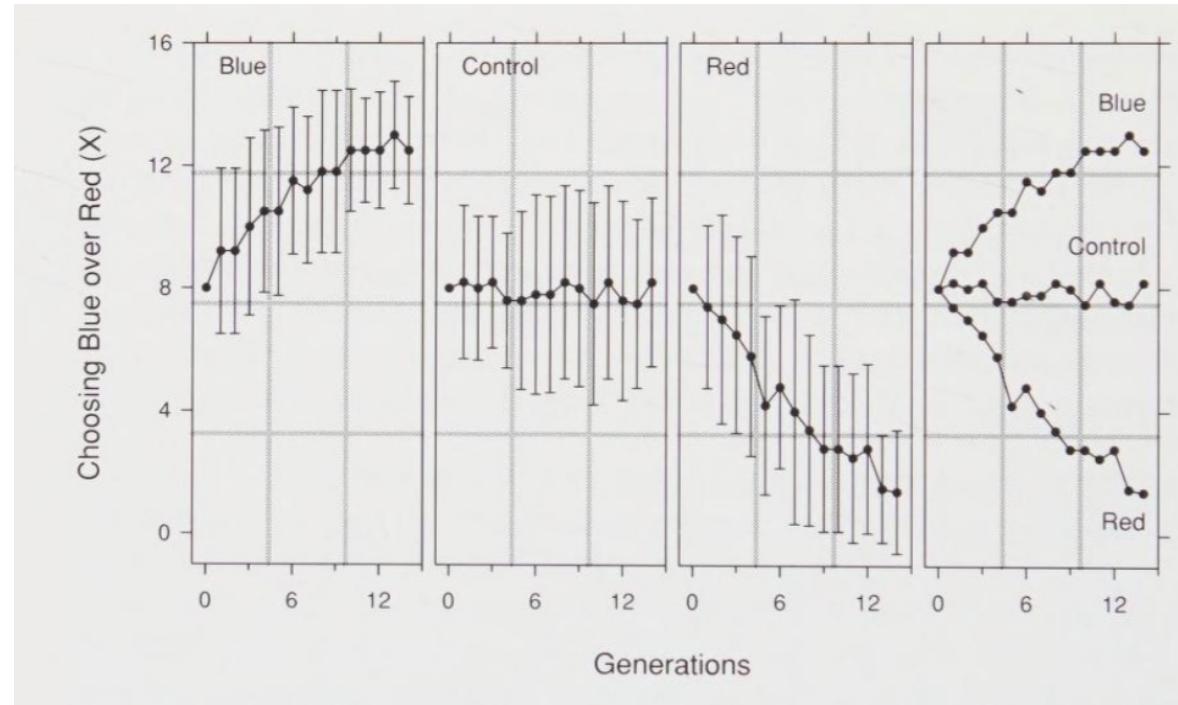
Principles of Graph Construction

- General strategy
 - A large amount of quantitative information can be packed into a small region
 - Graphing data should be an iterative, experimental process
 - Graph data two or more times when it is needed
 - Many useful graphs require careful, detailed study

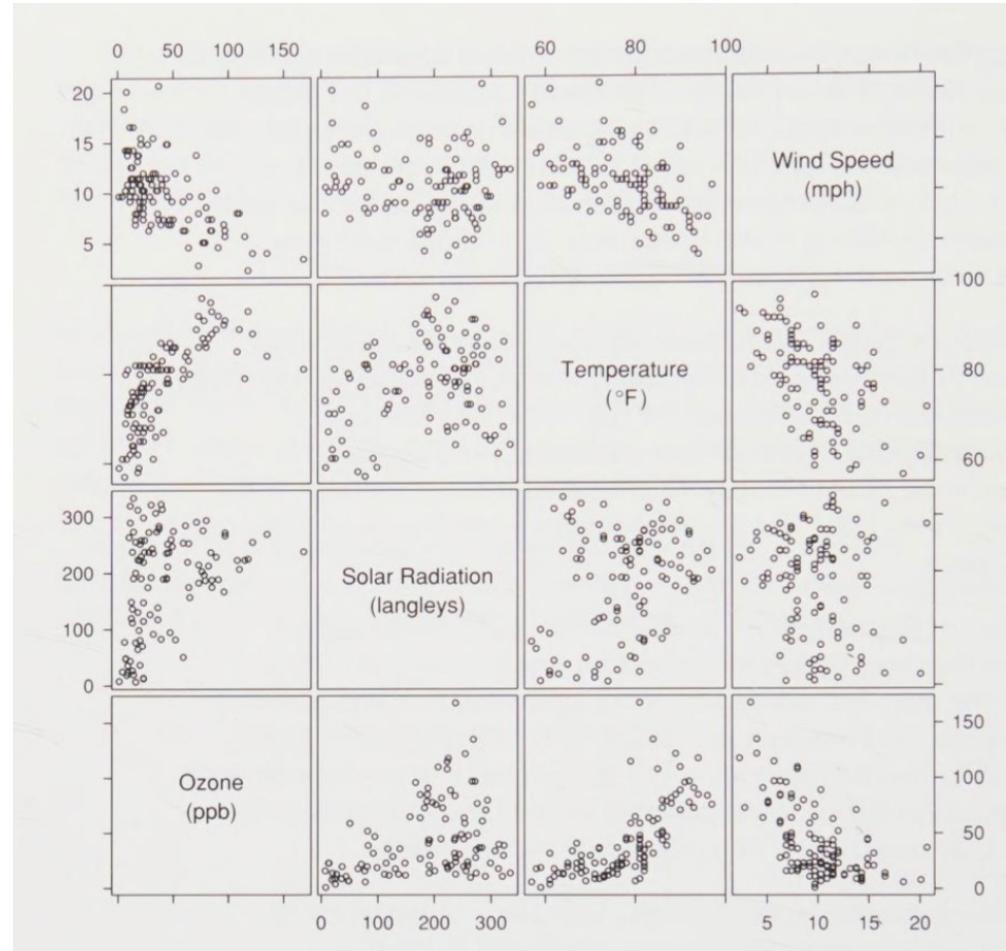
Graphing data should be an iterative, experimental process



Graph data two or more times when it is needed



Many useful graphs require careful, detailed study



Graphical Methods

- The usefulness of various types of plots → this is left to be studied throughout our course
 - Logarithms
 - Residuals
 - Distributions
 - 1D scatterplots and histograms
 - Quantile plots
 - Box plots
 - Q-Q plots (quantile-quantile plots)
 - Dot plots and multi-dot plots
 - Plotting symbols
 - Overlap
 - Superposed plotting symbols
 - Superposed curves
 - Visual reference grids
 - Loess (**locally weight** regression)
 - Time series
 - Scatter plots
 - Statistical variation

Graphical Perception

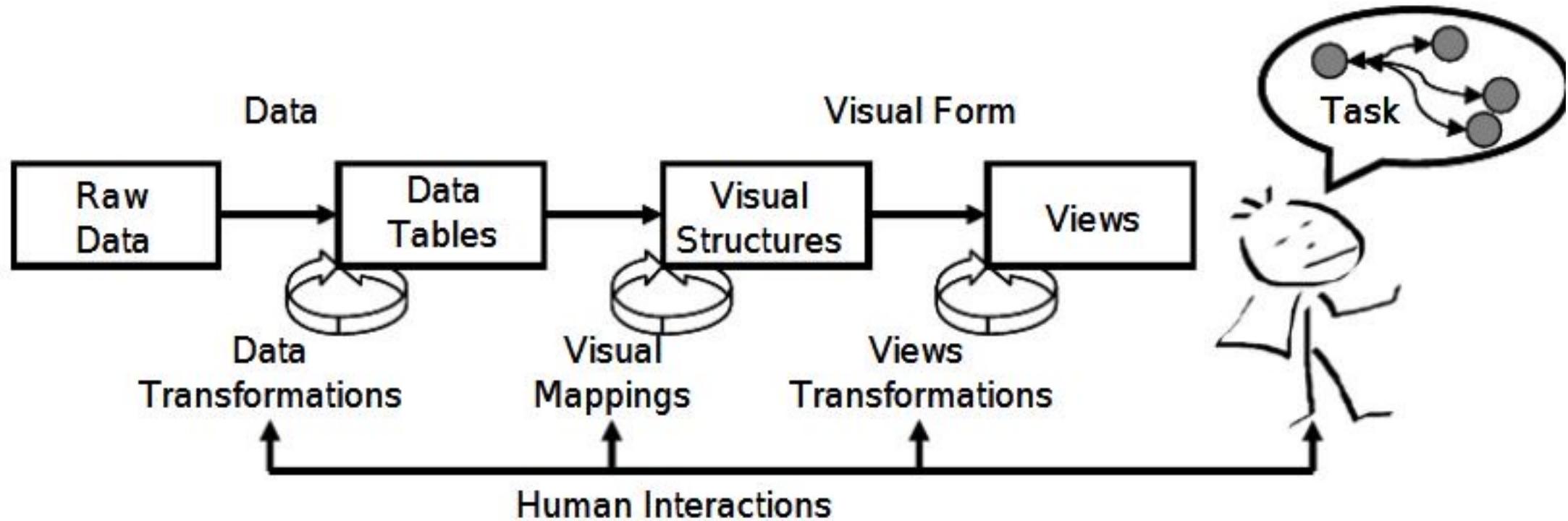
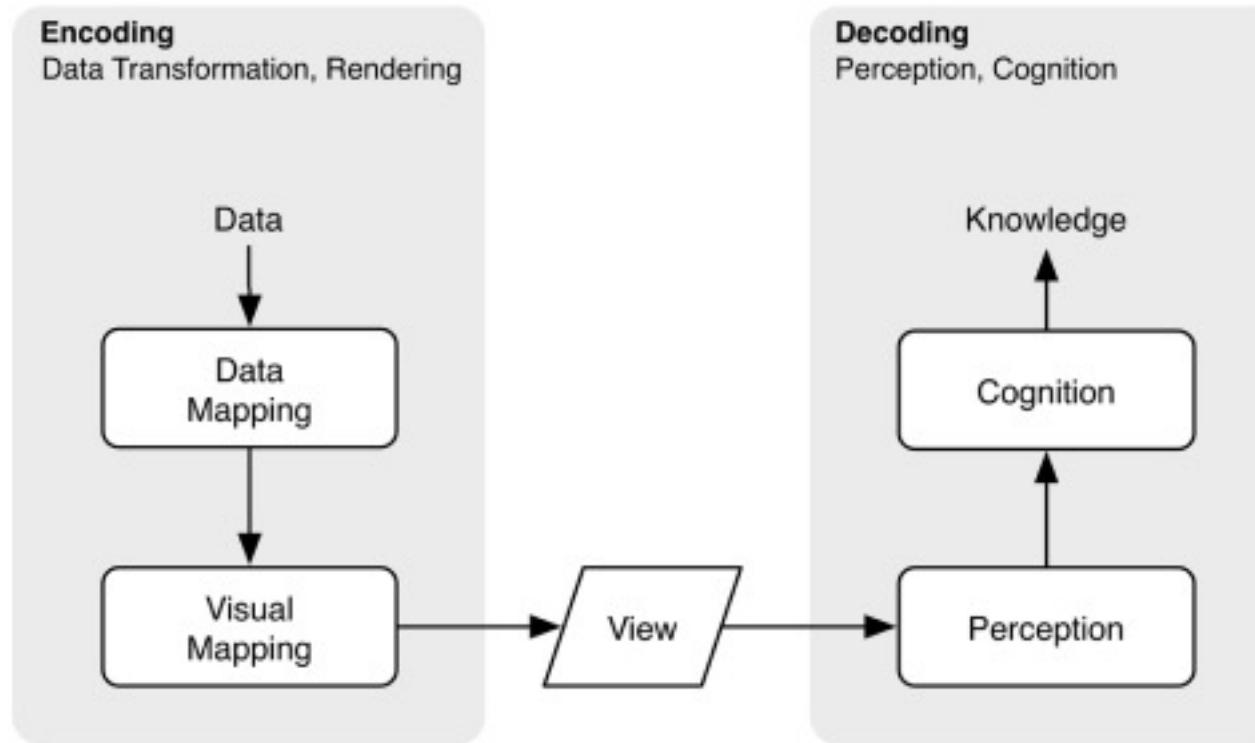


Figure 1 - Reference model for visualization (Card 1999).

Graphical Perception

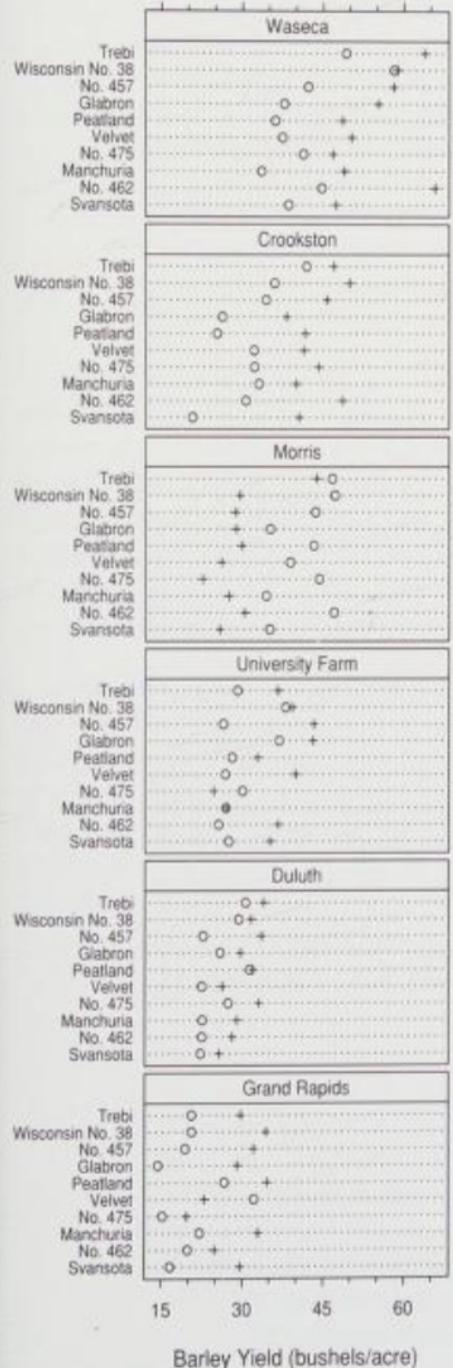


Graphical Perception

- Quantitative and categorical information
- Scale information and physical information
- Table look-up and pattern perception
 - Visual decoding of scale information
 - Visual decoding of physical information

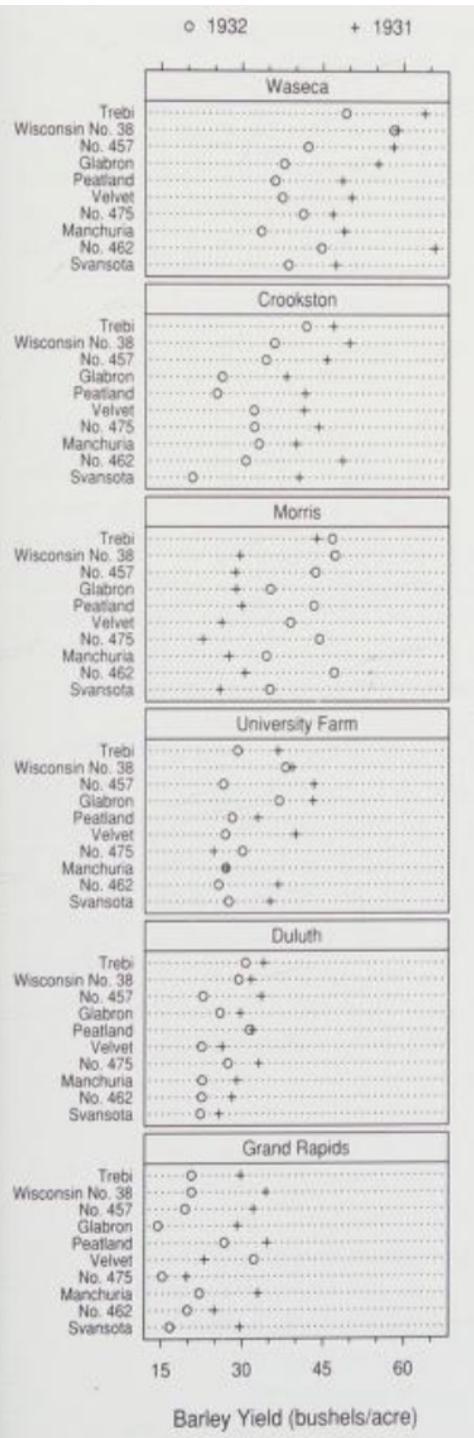
o 1932

+ 1931



Graphical Perception

- Three visual operations of table look-up
 - Scanning
 - Visual inspection of the physical elements
 - Interpolation
 - Visual estimation of values based on a set of discrete known data values
 - Matching
 - Visual identification of a geometric element with a labeled data value

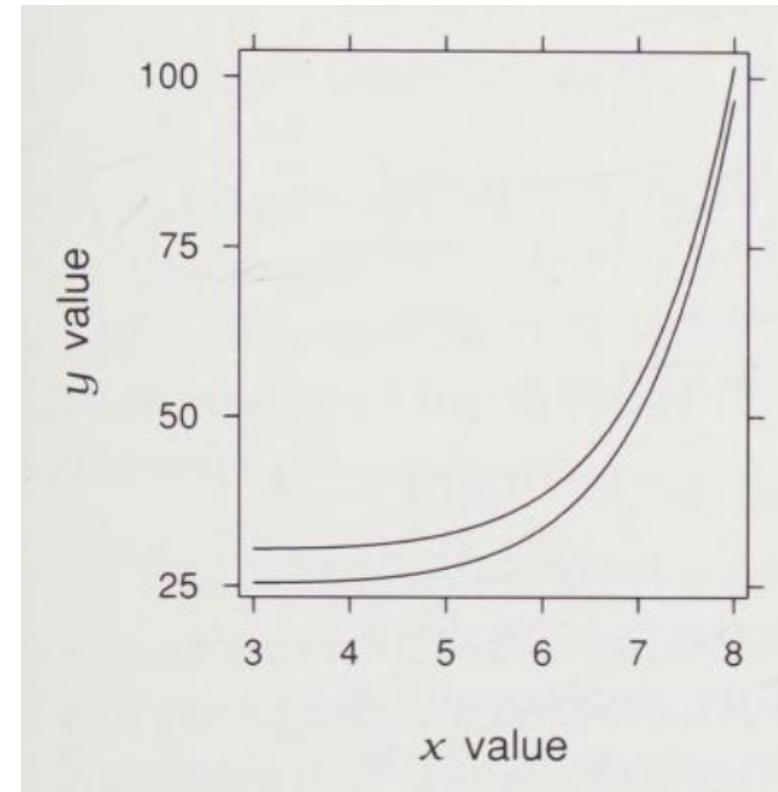


Graphical Perception

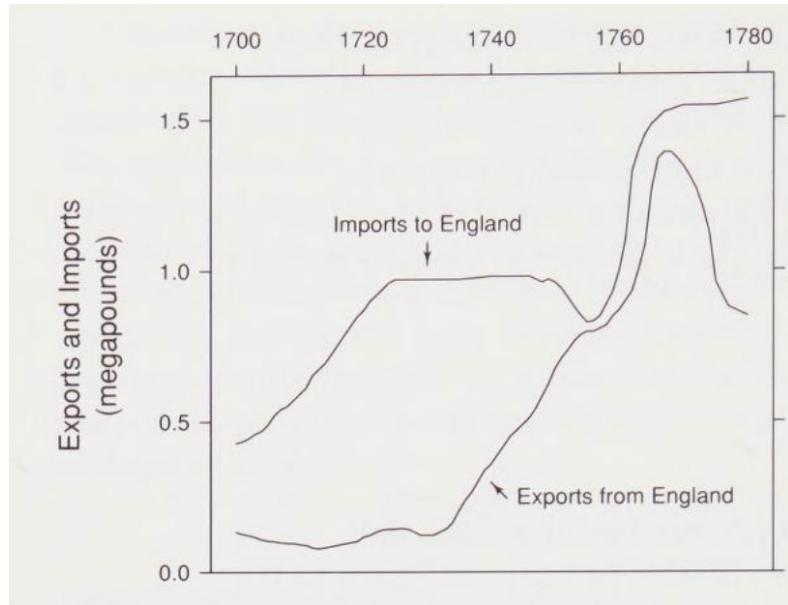
- Three visual operations of pattern perception
 - Detection
 - Visual recognition of a geometric aspect that encodes physical values
 - Assembly
 - Visual grouping of detected physical elements; discerning overall patterns in data
 - Estimation
 - Visual assessment of relative magnitudes of two or more physical values.
 - Discrimination (equal or not equal)
 - Ranking (less than or greater than)
 - Ratioing (a/b)

Pattern perception for curves can be inaccurate

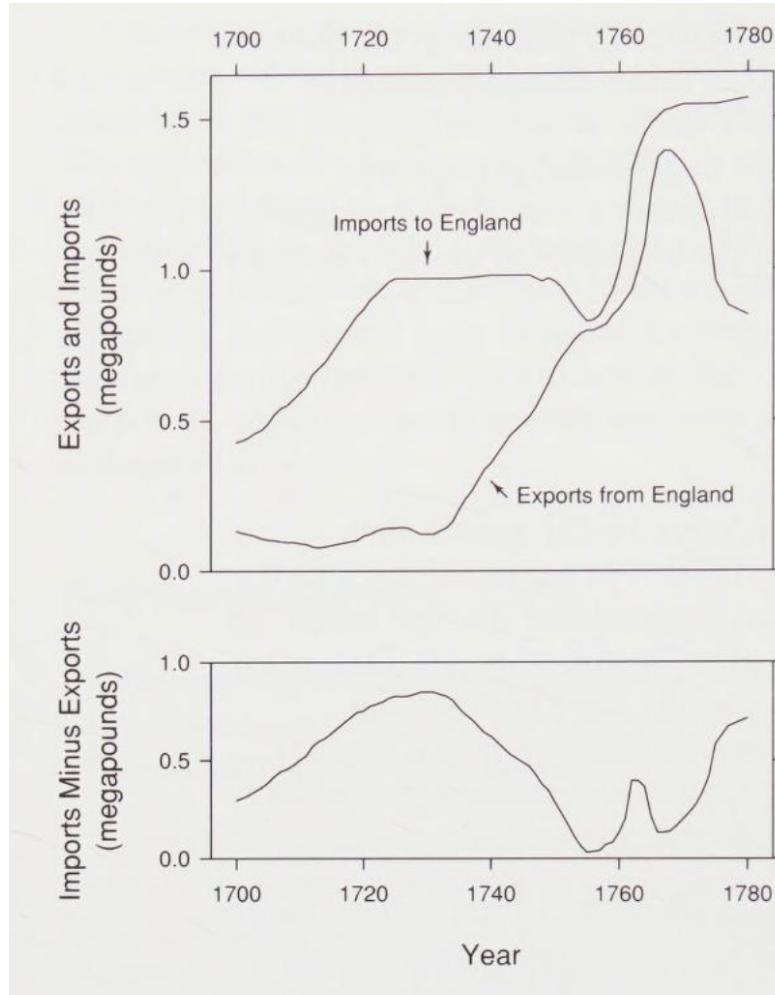
Not because we are inaccurate at estimating differences,
but because we're inaccurate at detecting differences



Pattern perception for curves can be inaccurate

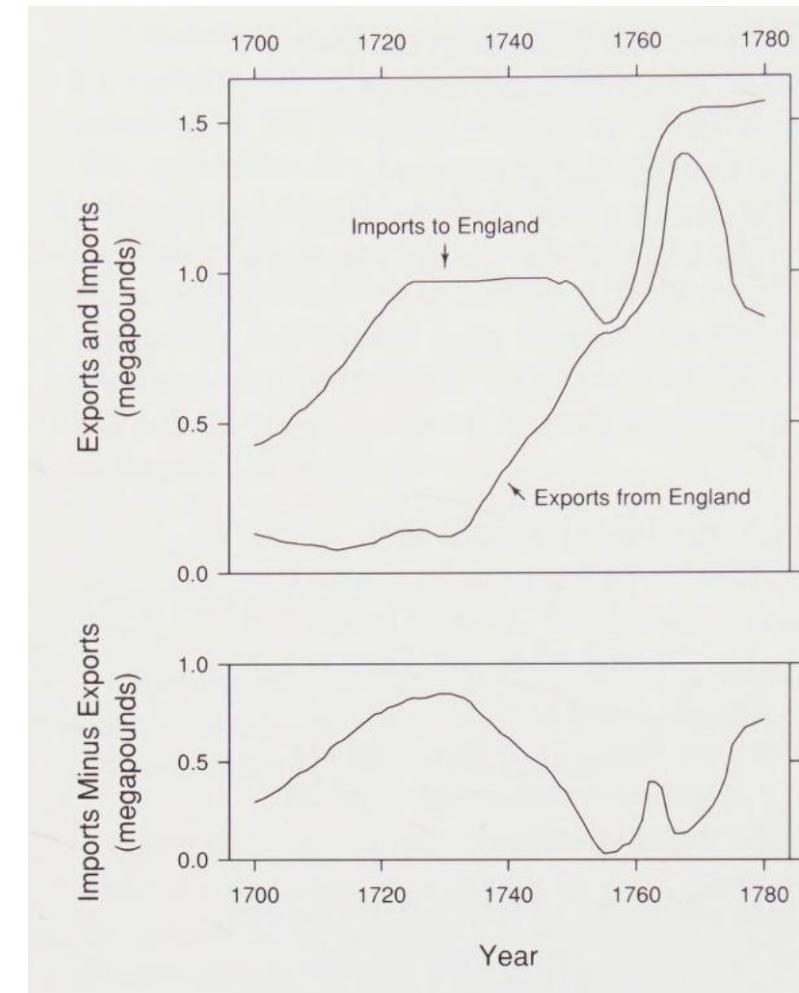


Pattern perception for curves can be inaccurate



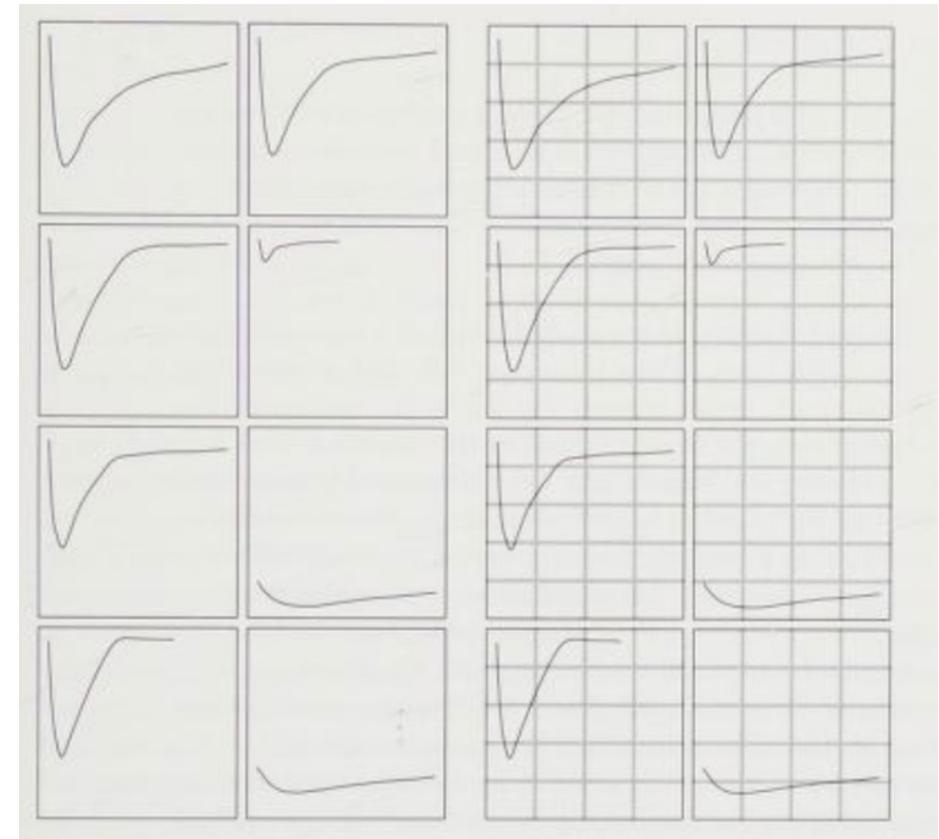
Pattern perception for curves can be inaccurate

- Table look-up is more efficient for the bottom panel than the top when calculating differences
- Remedy: superpose two curves on one panel and graph the differences on another



Pattern perception for curves can be inaccurate

- Table look-up is more efficient for the bottom panel than the top when calculating differences
- Remedy: superpose two curves on one panel and graph the differences on another
- Another is to juxtapose curves on separate panels – eliminates the distortion caused by superposition



Color

- There are different approaches to color
 - HSL (Hue, Saturation, Lightness); CMYK (Cyan, Magenta, Yellow, Key (Black)), RGB (Red, Green, Blue)
- Changing hue can provide efficient discrimination of colors
 - For fixed hue, we can perceive ordering as either lightness changes or saturation changes
 - We cannot effortlessly perceive an ordering to changing hue.
 - Hue can be effective at perceiving boundaries between adjacent levels of color encoding
 - Two hues provide more color variation than just one
- We will show more on color in another session

- Marker identification
 - Assembly (the grouping of elements) depends on effectively discriminating the texture of chosen markers

*	*	*	+	○	○	○
*	+	*	+	○	○	○
+						
*	*	*	+	○	○	○
*	+	*	+	○	○	○
*	*	+	*	○	○	○
*	*	+	*	○	○	○
*	*	+	*	○	○	○

(+ ○)

(+ □)

V J J ^ + x *
 J ^ V ^ x x + x
 ^ ^ J ^ + x x x
 ^ ^ < L x x + x
 V ^ V > x x x x
 V V V ^ + + x x
 L J V V + x + x
 V V ^ J x x + x
 (L +)

$\begin{array}{cccccc} x & x & x & + & - & \times & \times \\ x & + & x & + & x & y & T \\ + & x & x & x & \perp & T & \times \\ x & + & x & + & y & y & -x \\ x & + & x & + & y & y & y \\ x & x & + & x & \perp & - & x \\ x & x & + & x & T & x & T \\ x & x & + & x & - & \times & y \end{array}$

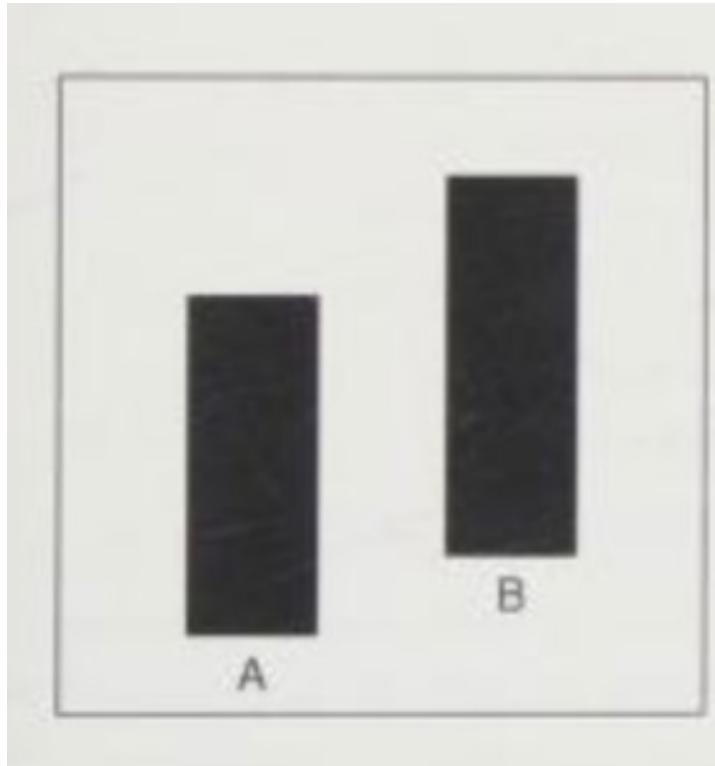
(+ X)

۴۷۸ ۳۷۸
۴۷۹ ۳۷۹
۴۸۰ ۳۸۰
۴۸۱ ۳۸۱
۴۸۲ ۳۸۲
۴۸۳ ۳۸۳
۴۸۴ ۳۸۴
۴۸۵ ۳۸۵
۴۸۶ ۳۸۶
۴۸۷ ۳۸۷
۴۸۸ ۳۸۸
۴۸۹ ۳۸۹
۴۹۰ ۳۸۹

$(L_L M_L)$

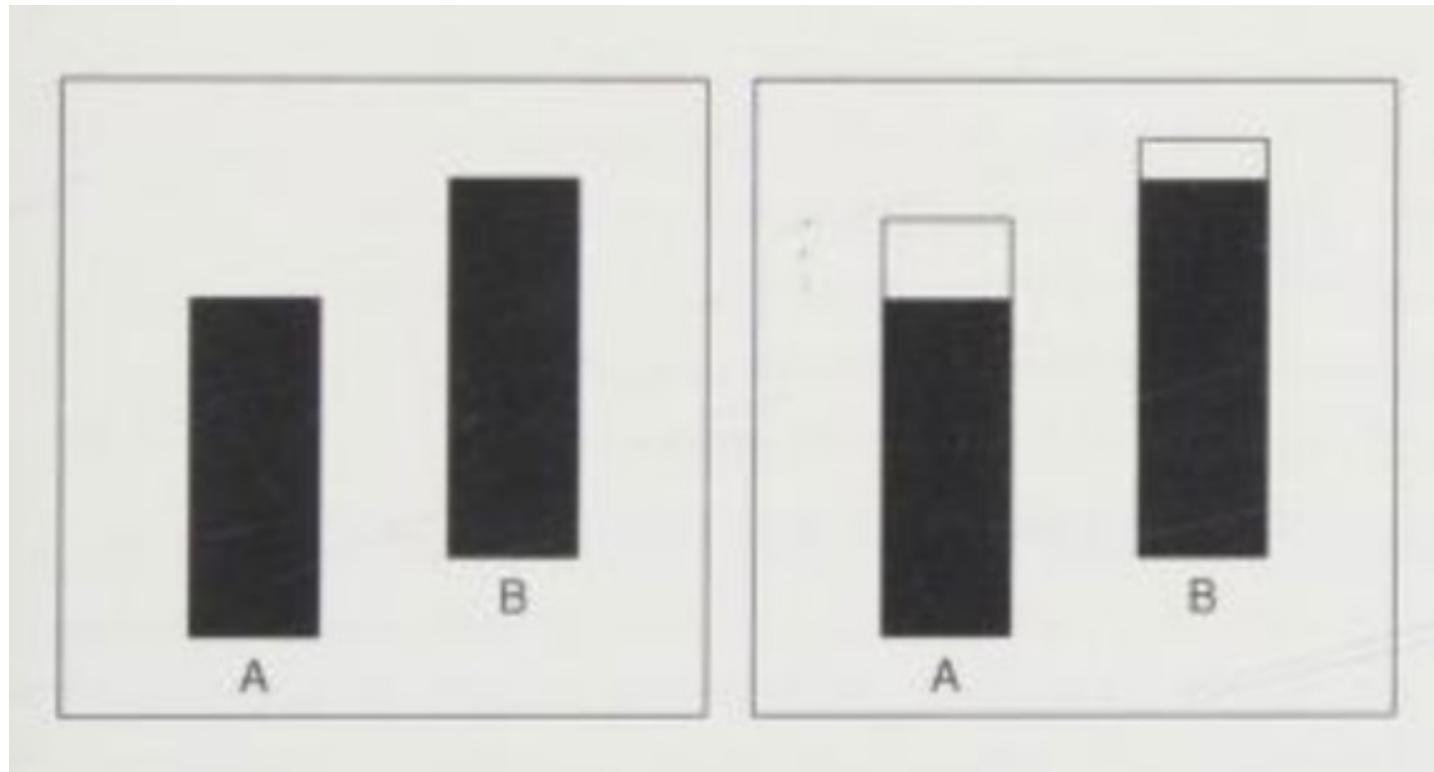

 (R-mirror-R)

Detecting length differences is more accurate for large percent change

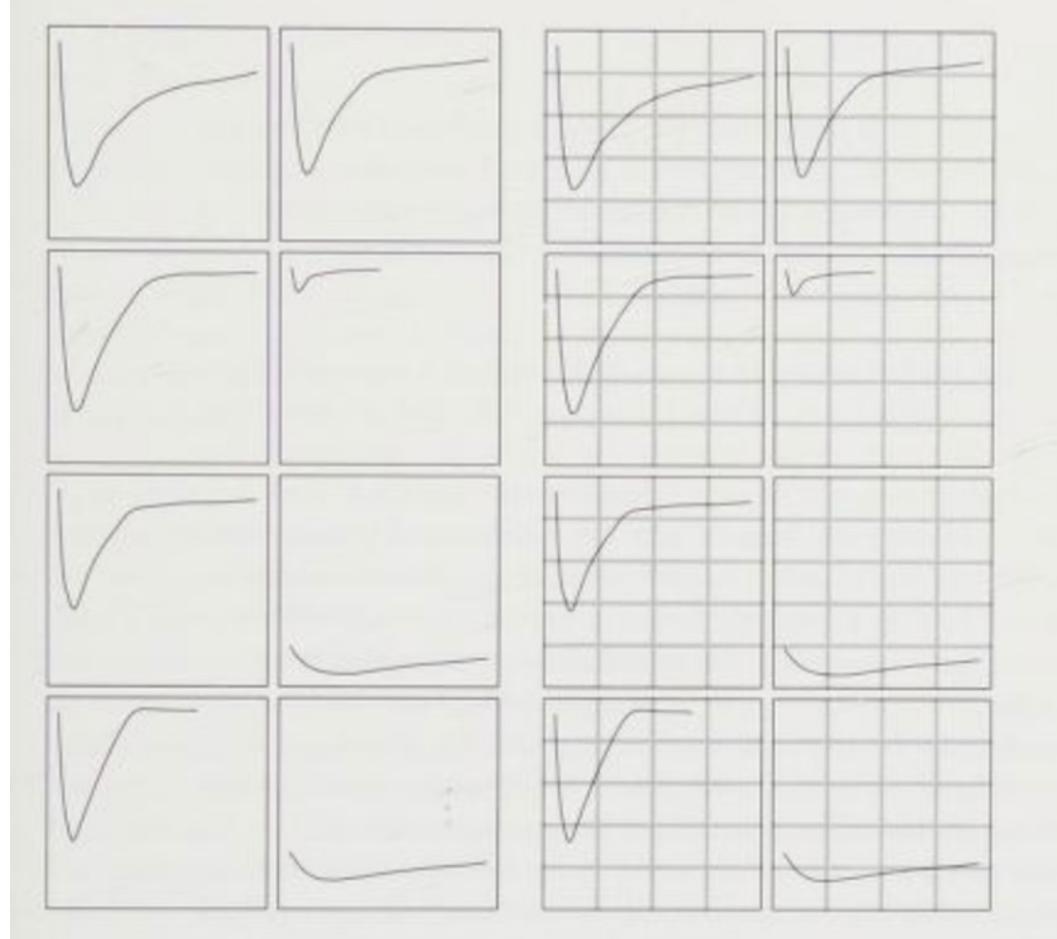


- Which bar is longer?

Detecting length differences is more accurate for large percent change

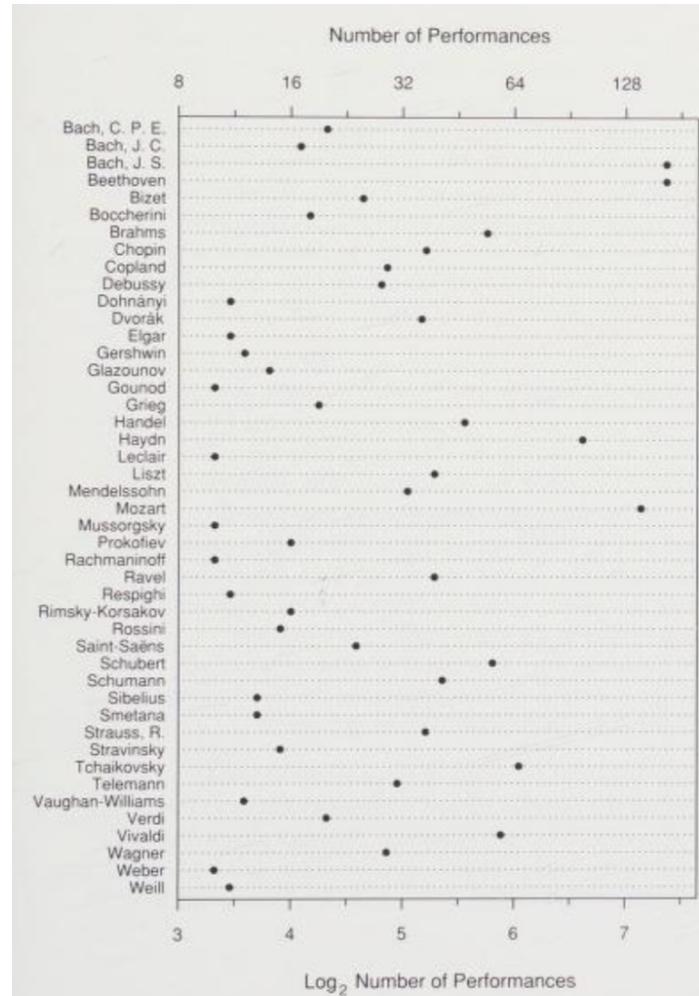


Grids can enhance pattern perception



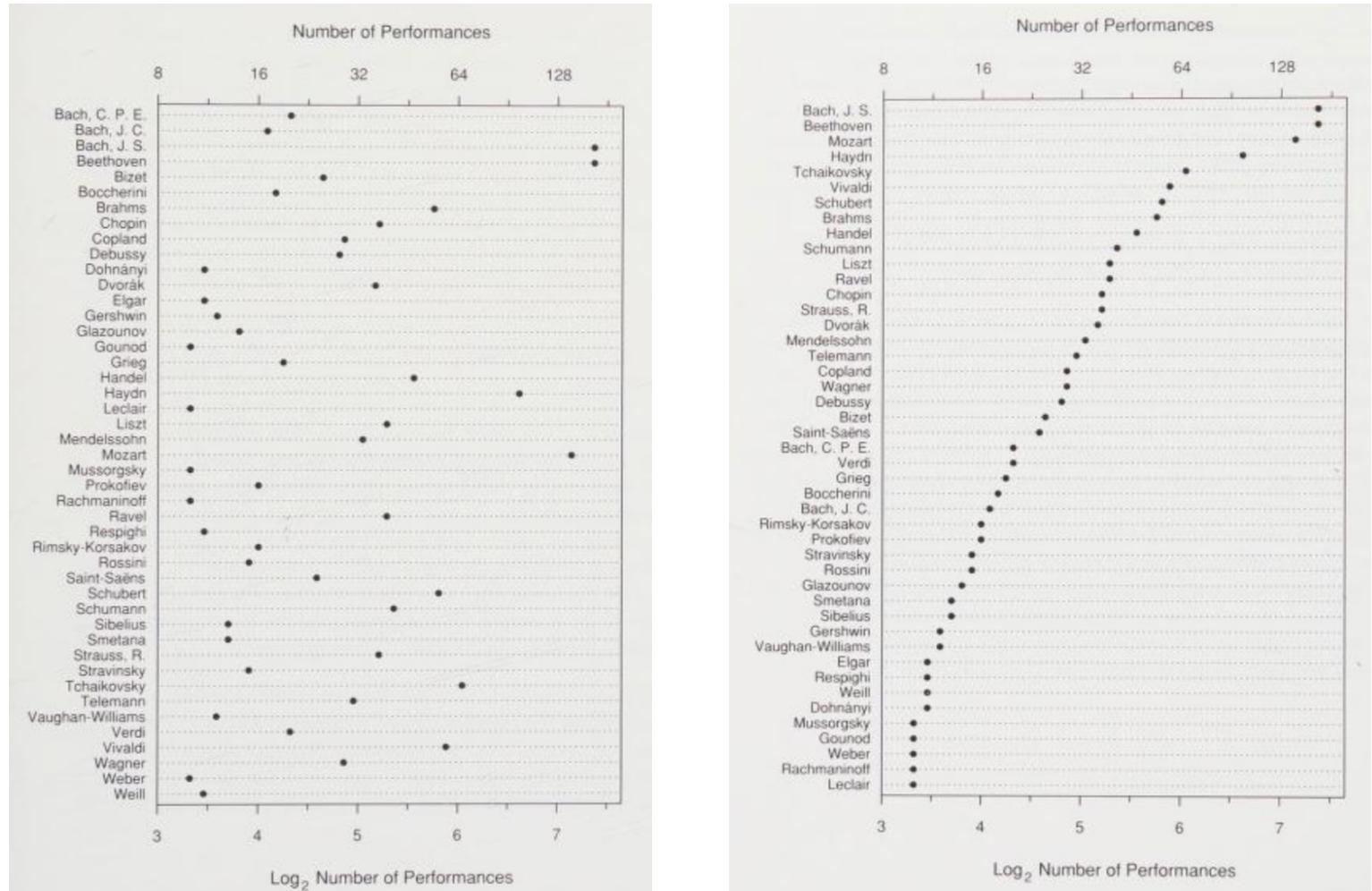
Ordering can make perception more effective

- Who are the top five composers?



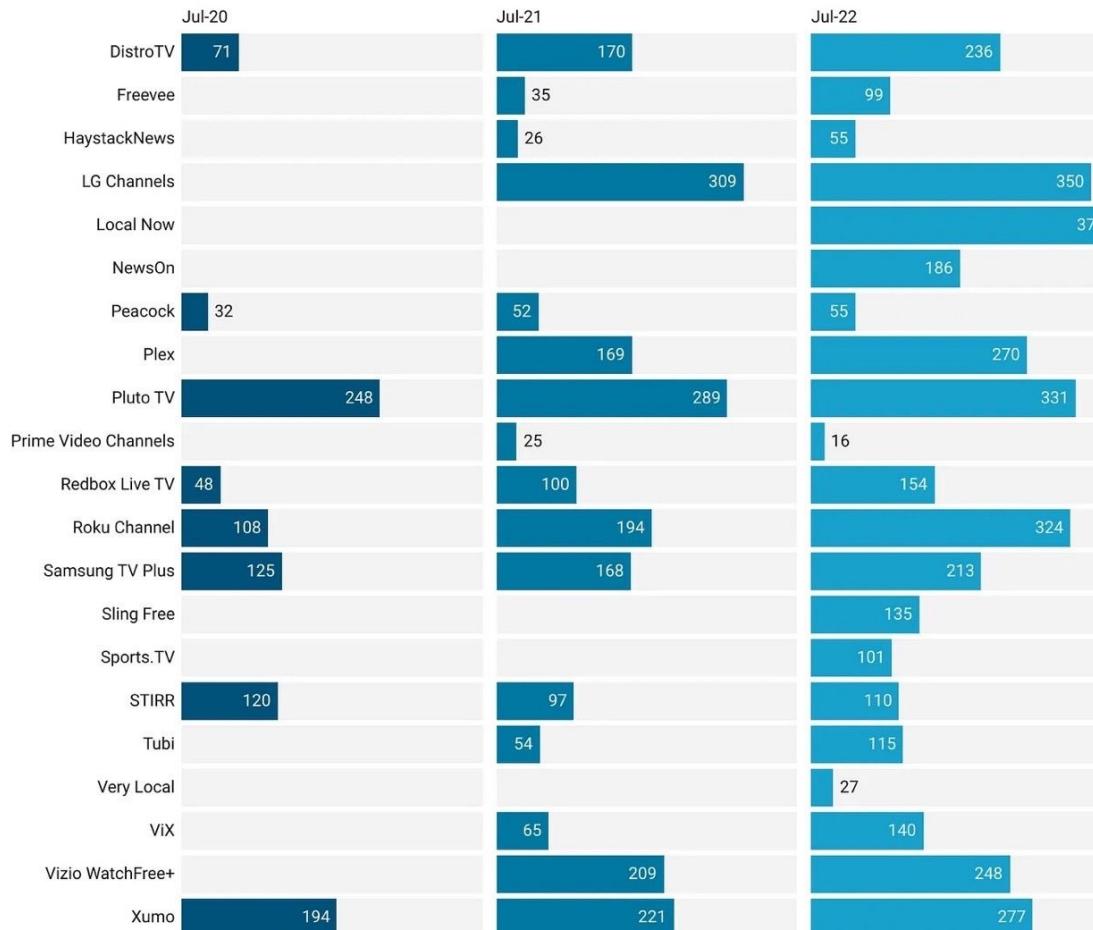
Ordering can make perception more effective

- Who are the top five composers?



Ordering can make perception more effective

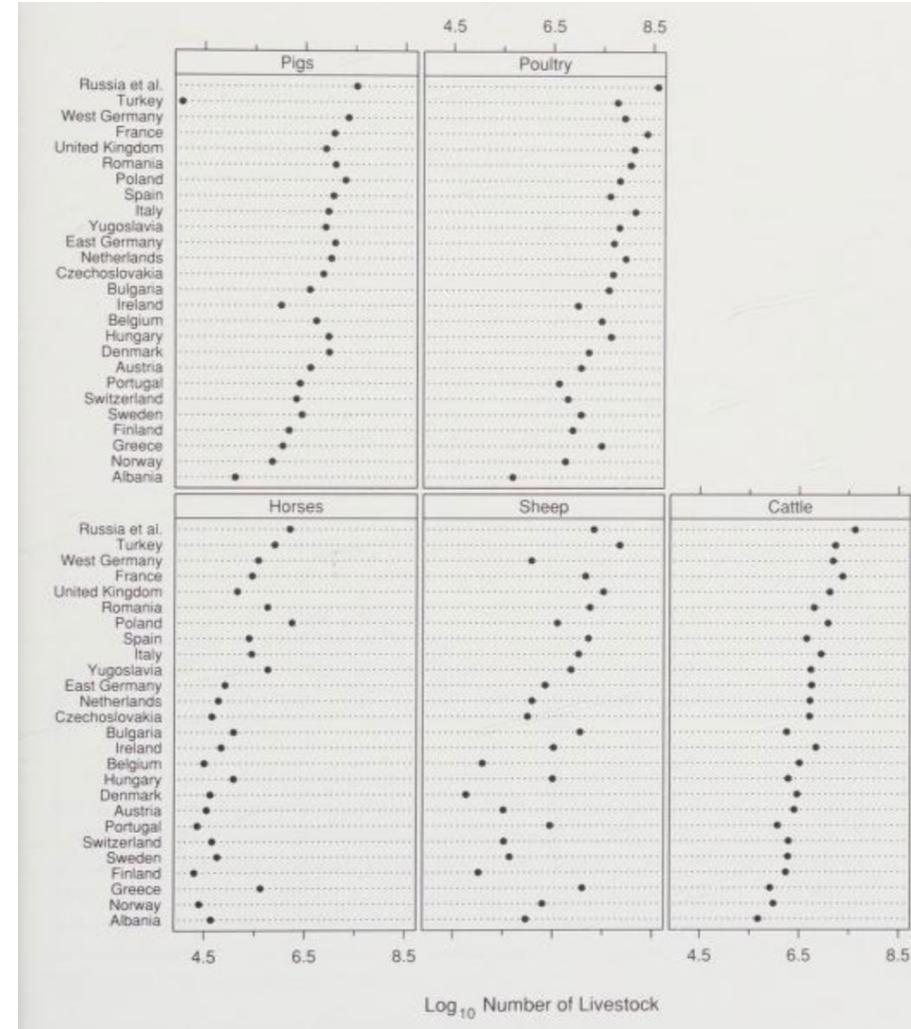
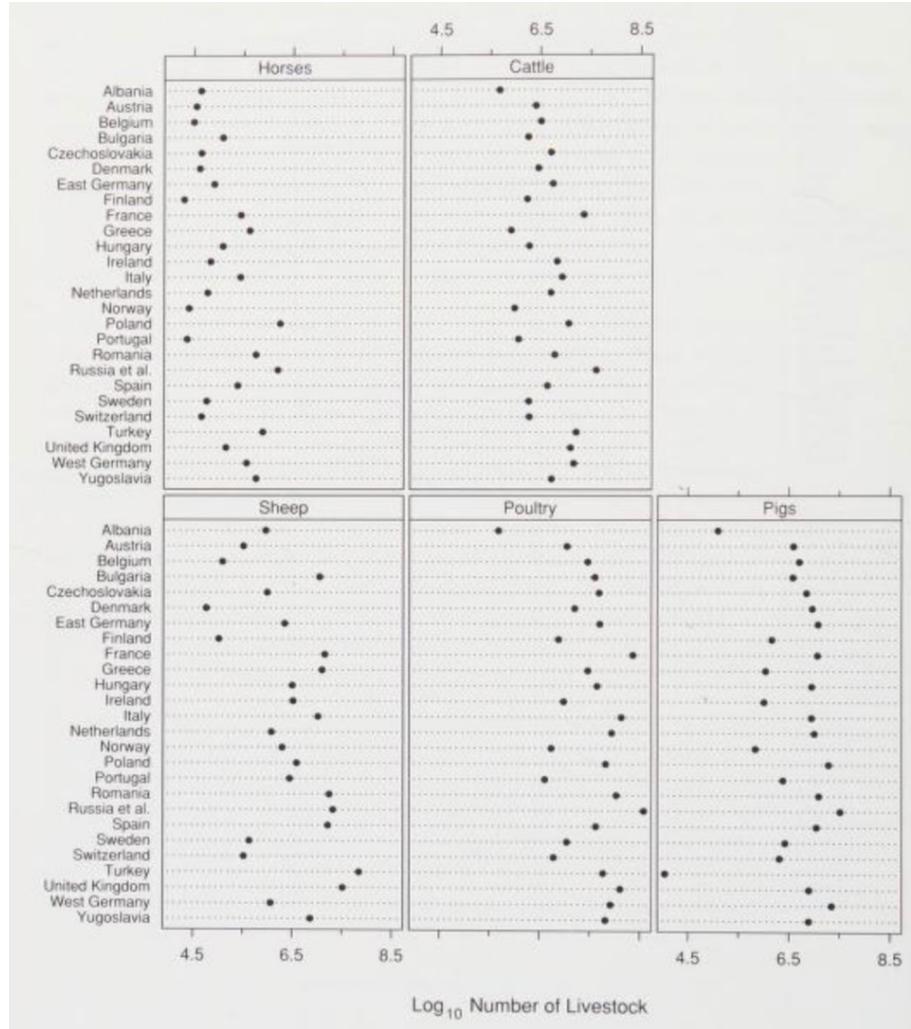
Number of FAST Channels on FAST Services in July 2020, 2021 and 2022



Figures for News by Fire TV are under revision and not included for July

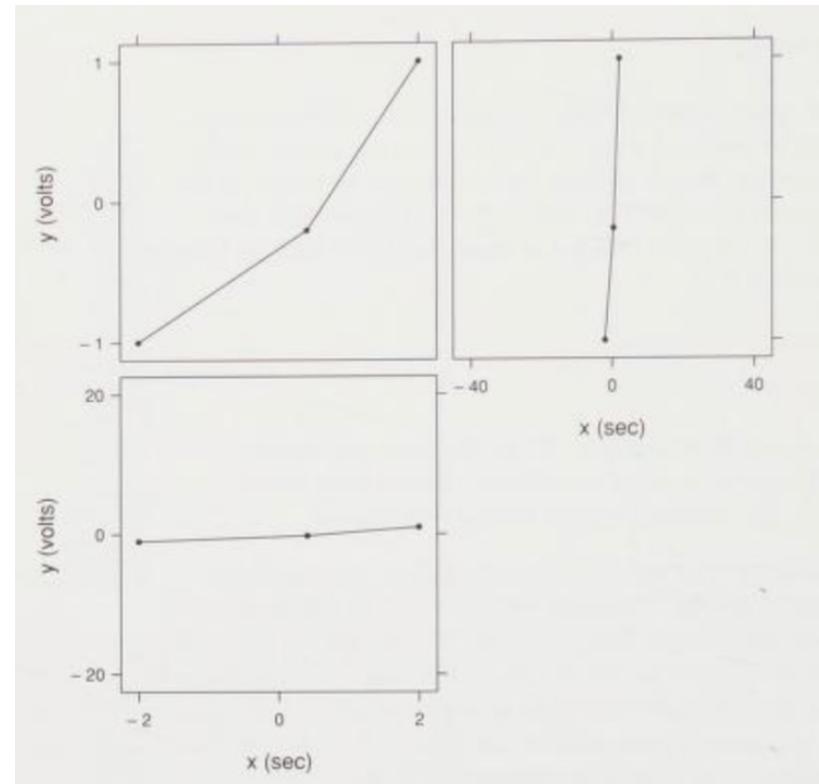
Chart: Gavin Bridge (c/o The FASTMaster) • Source: fastmaster.substack.com • Created with Datawrapper

The same follows for multivariate data

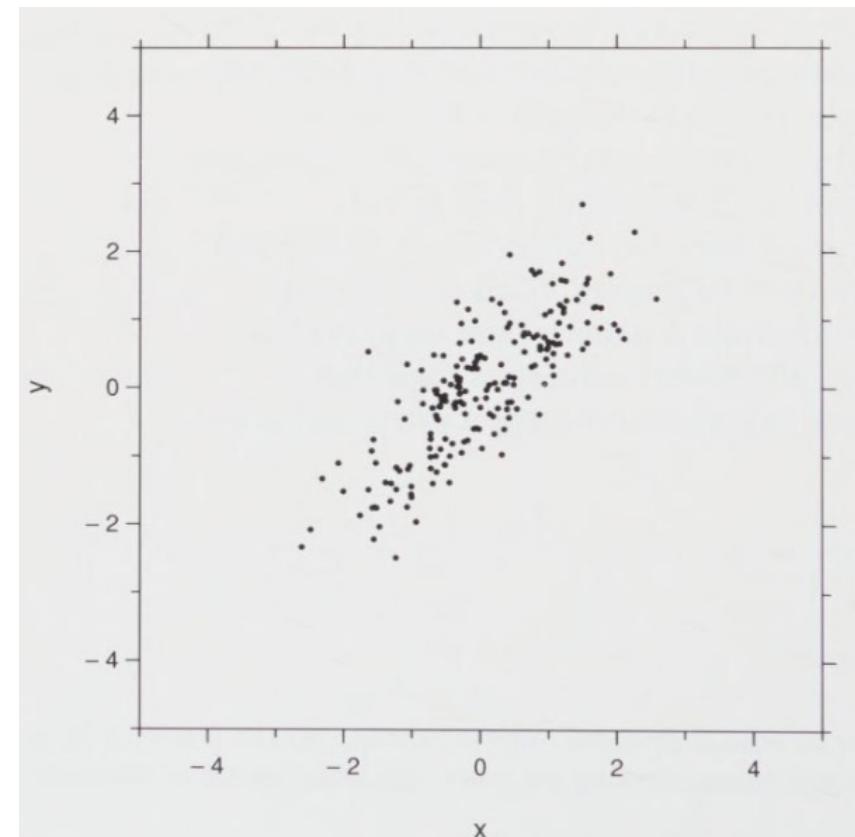
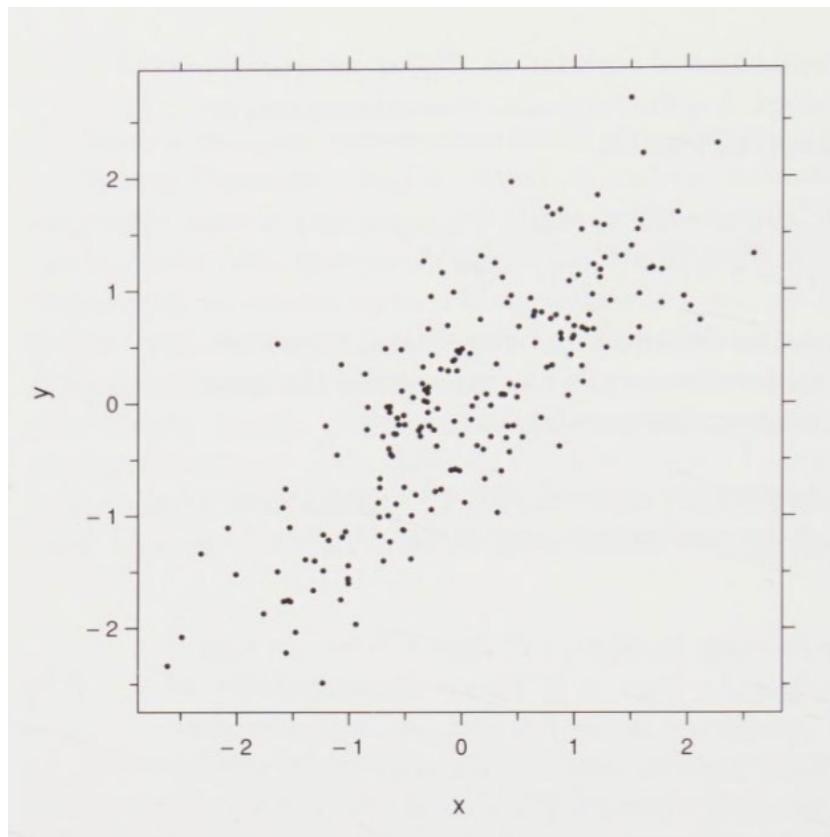


Aspect ratios should be set so that slopes average to be 45 degrees

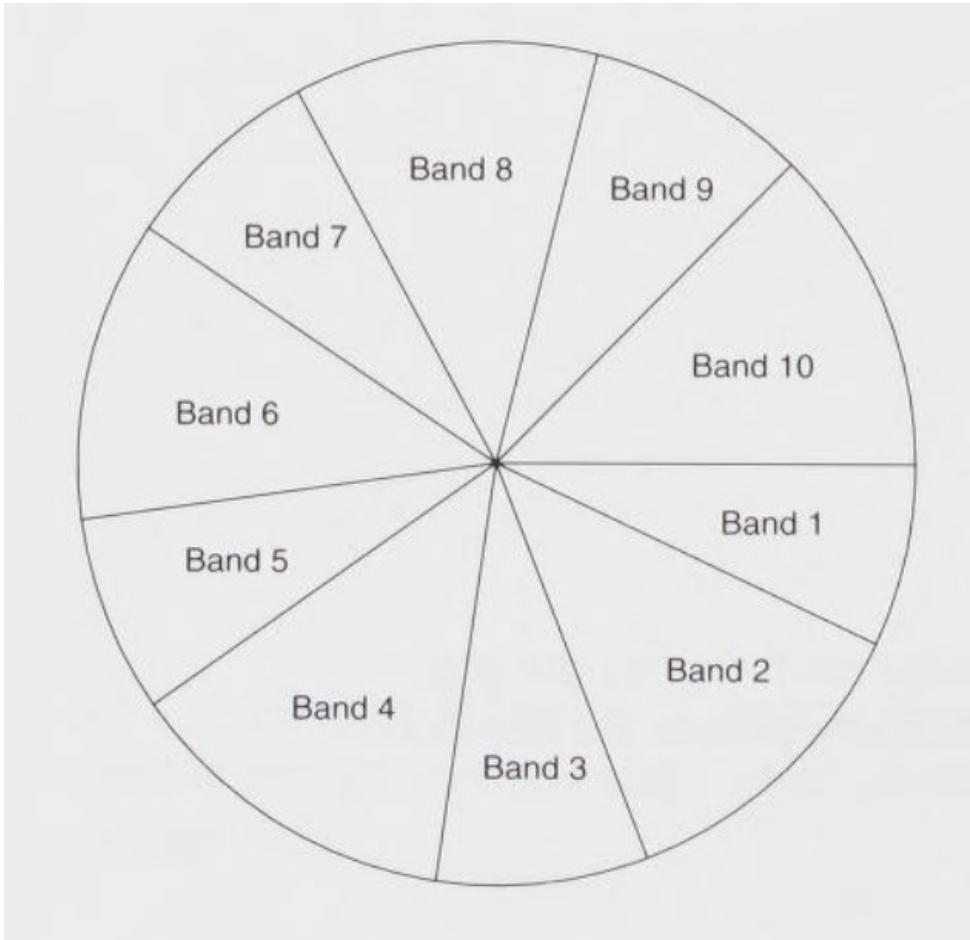
- This allows us to optimally distinguish angular differences
- If you need to efficiently table look-up the slope, consider deriving and graphing the slope directly
 - Steps to find the slope:
 - scan horizontally for the two endpoints
 - interpolate to get two scale values
 - do mental arithmetic to get difference
 - do the same for the horizontal difference
 - mentally divide to get slope



Correlation – when comparing, make sure that the ratio of the data rectangles to the scale-line rectangles are equal

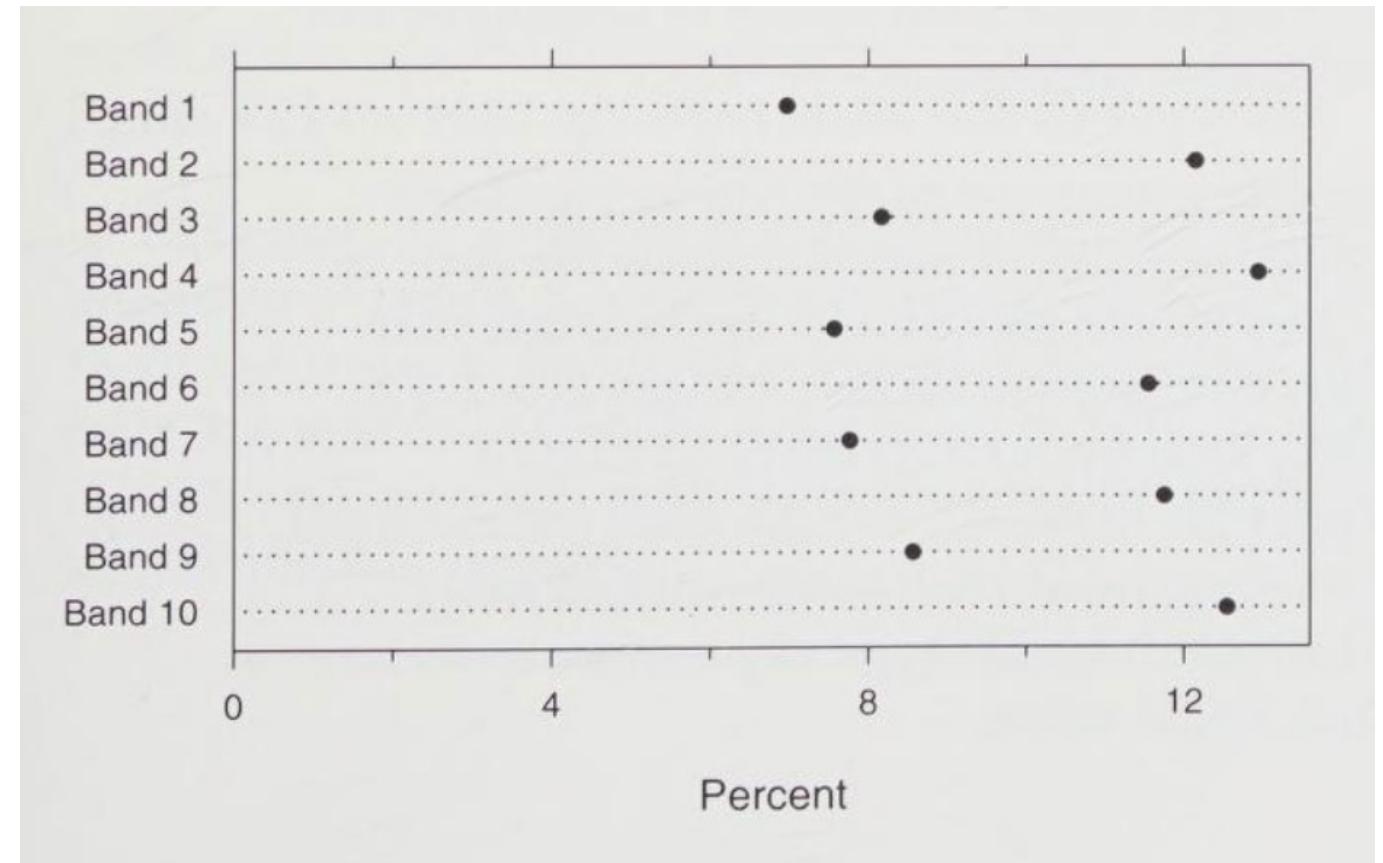
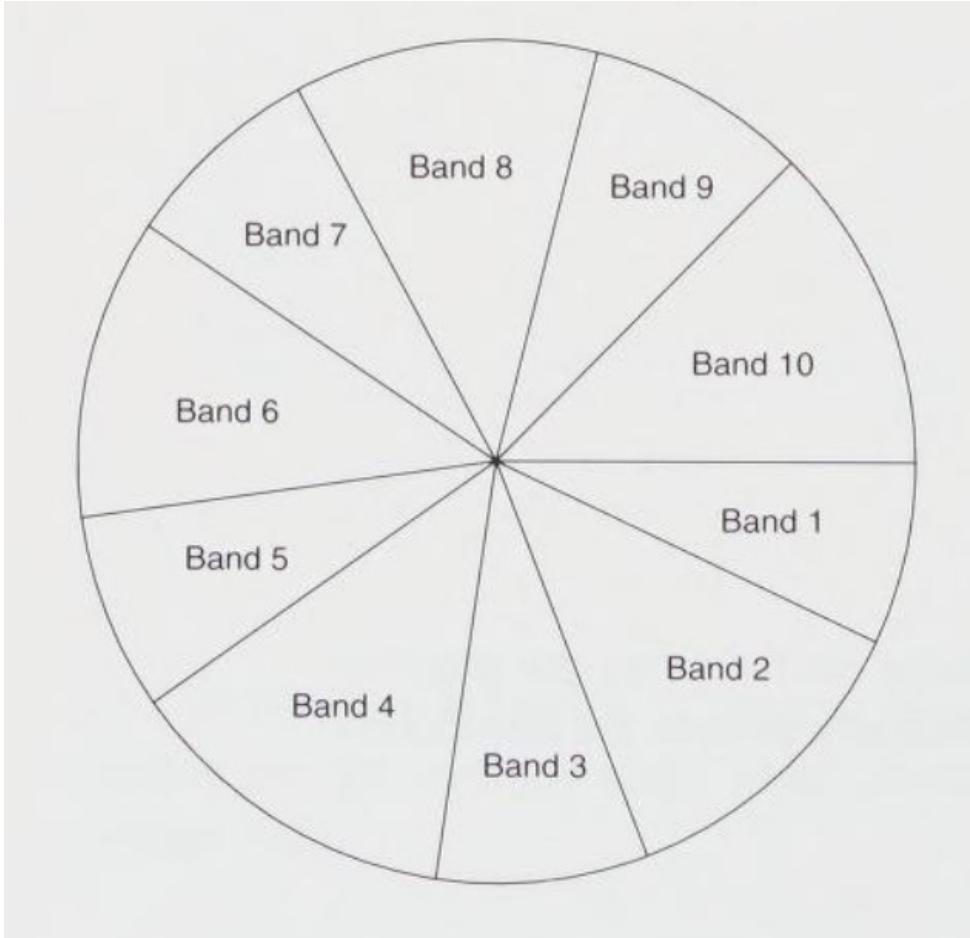


We just aren't as good at estimating areas

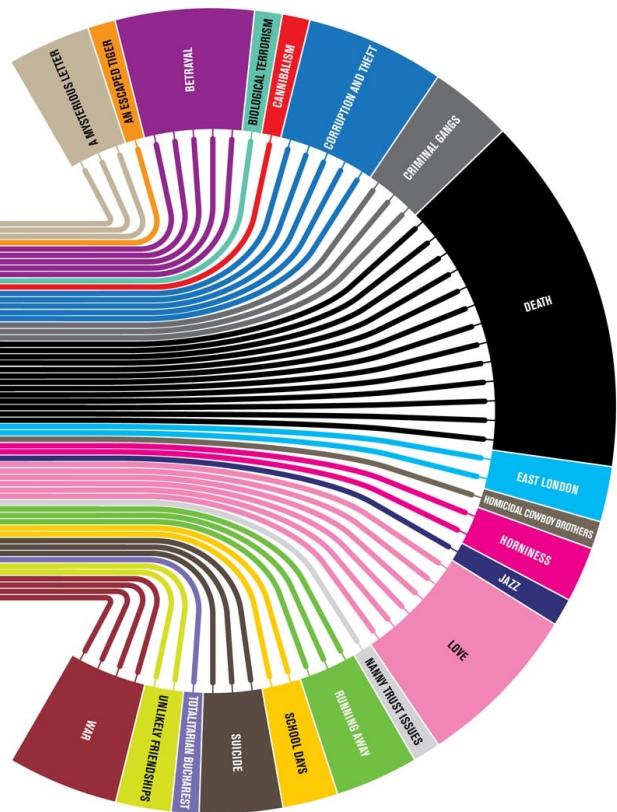
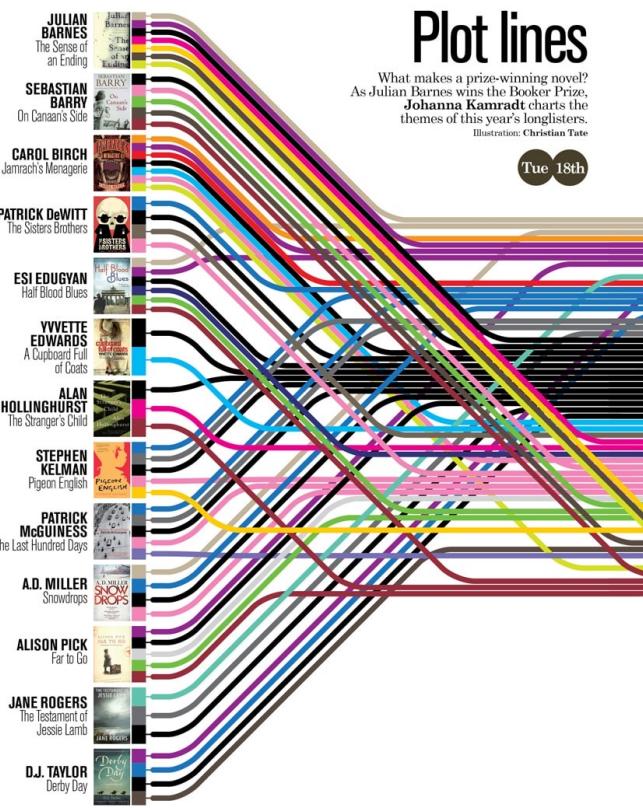


- Who are the top five bands?
- How much larger is #1 than #5?

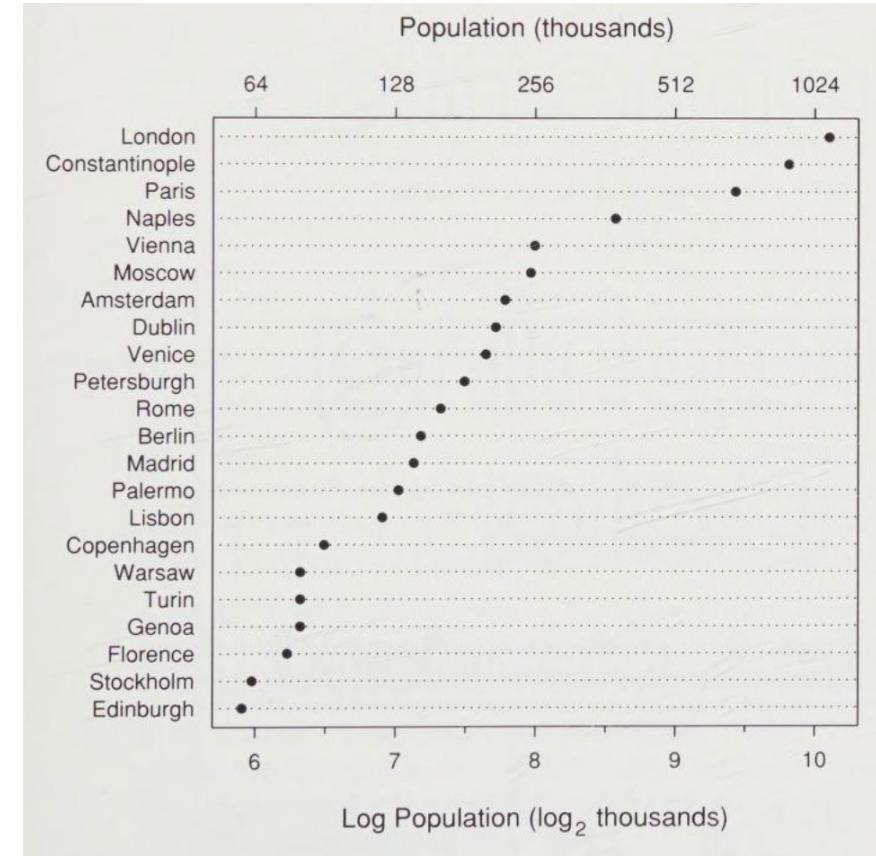
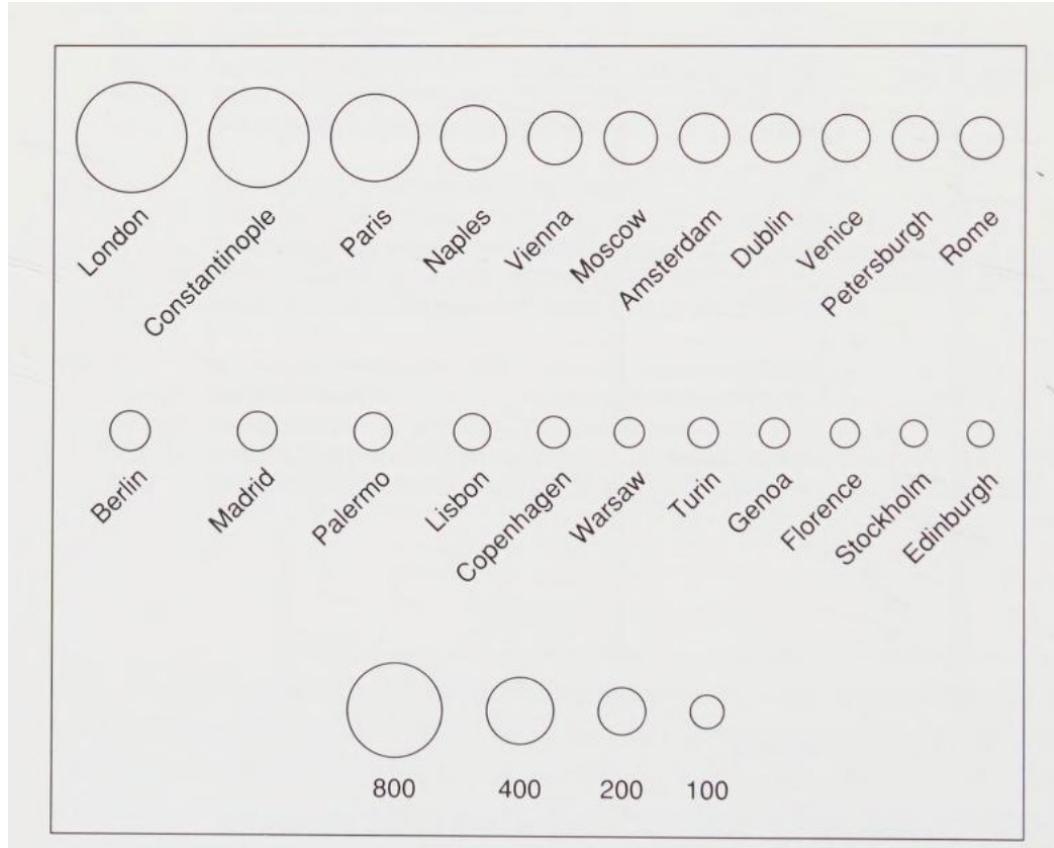
We just aren't as good at estimating areas



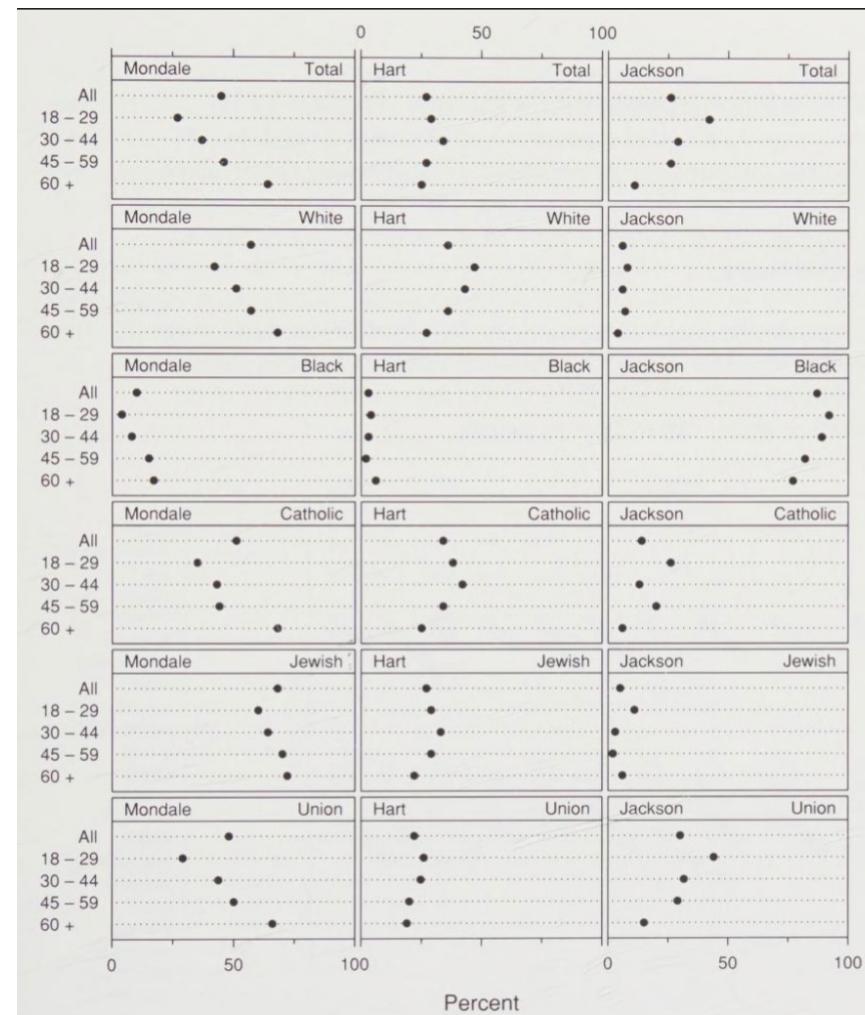
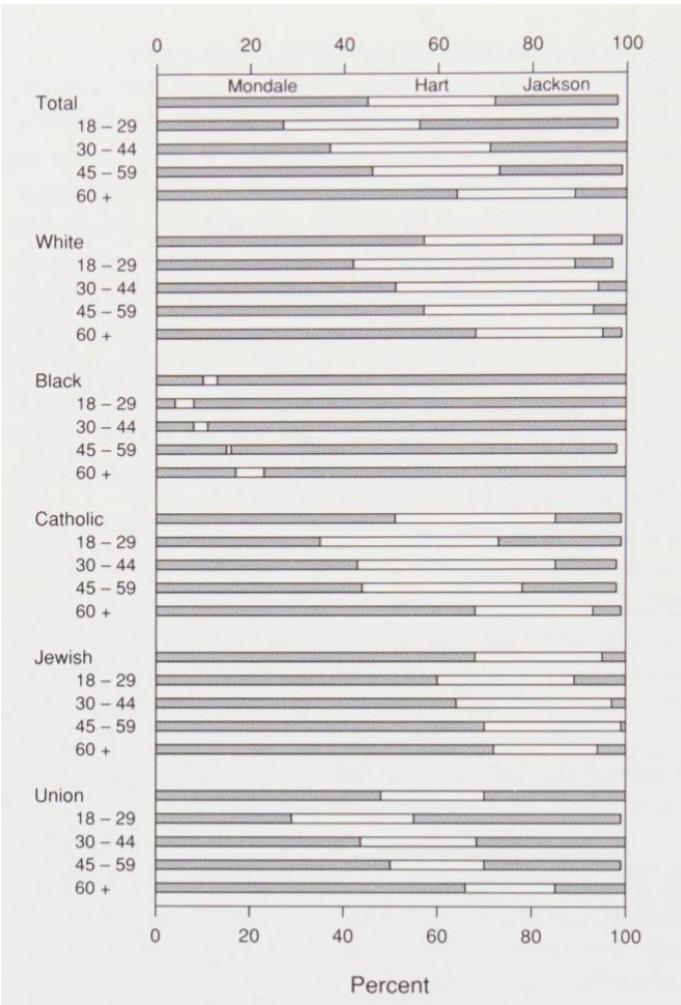
We just aren't as good at estimating areas



We just aren't as good at estimating areas

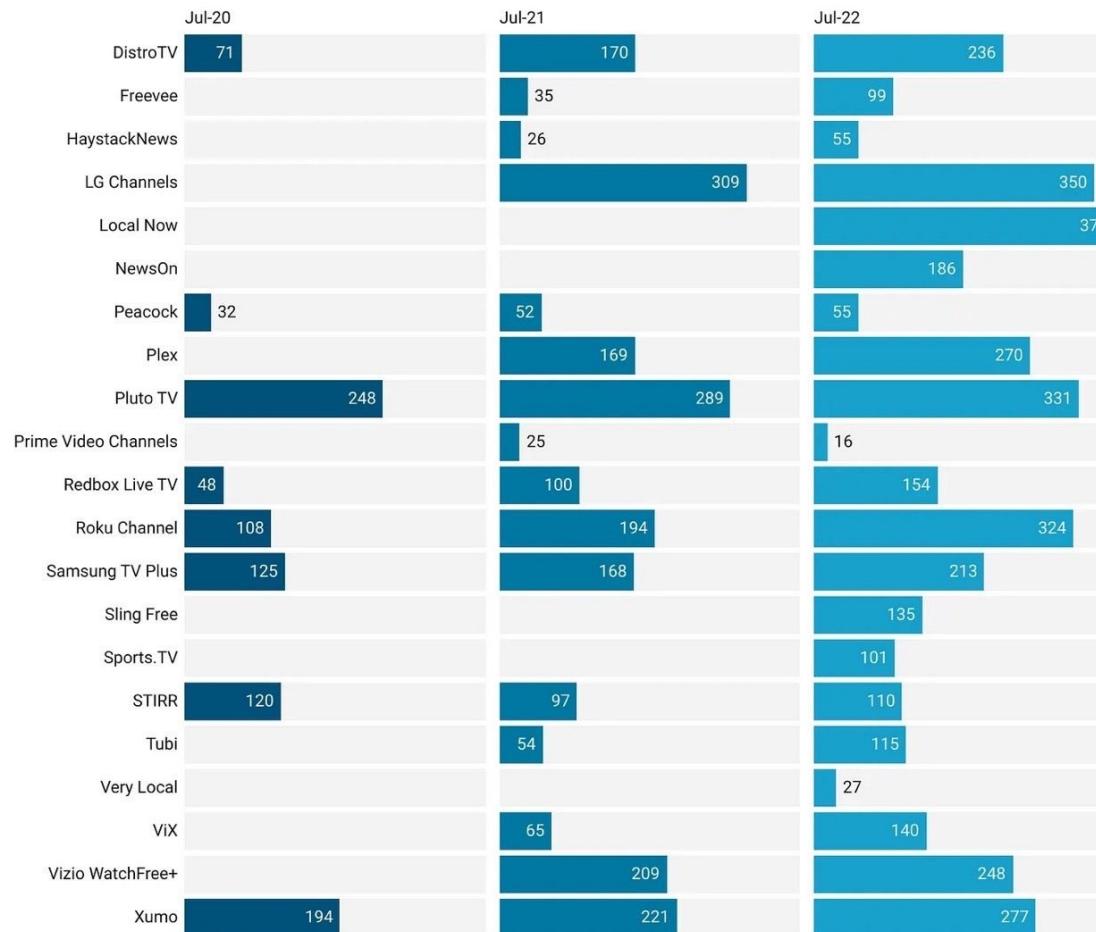


We aren't as perceptive with stacked bar charts as we are with plots that have positions along common scales.



We aren't as perceptive with stacked bar charts as we are with plots that have positions along common scales.

Number of FAST Channels on FAST Services in July 2020, 2021 and 2022



Figures for News by Fire TV are under revision and not included for July

Chart: Gavin Bridge (c/o The FASTMaster) • Source: fastmaster.substack.com • Created with Datawrapper

Next time:

We explore these analysis and visualization points
for
univariate data