# Data Science

## Lecture 2: Statistical Analysis / Visualizing Data

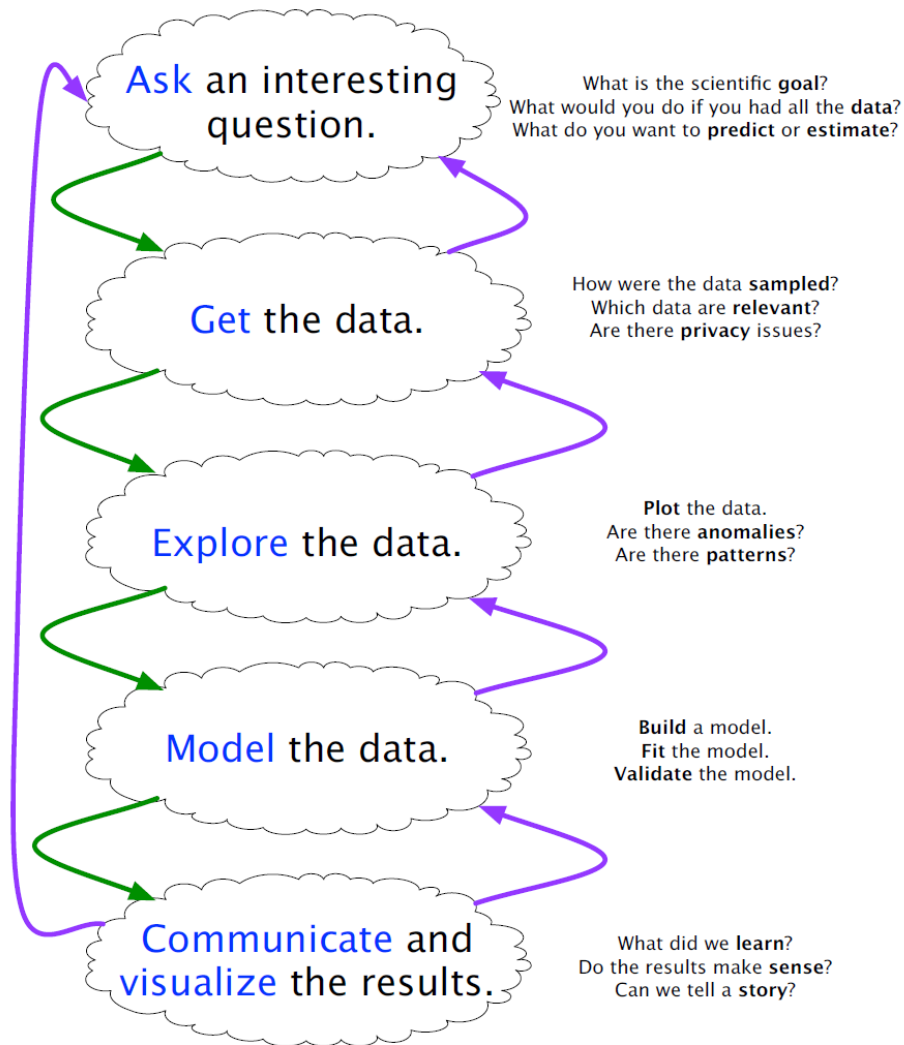# Statistical Analysis

**Typical Data Science Pipeline**

Ask an interesting question.

What is the scientific **goal**?
What would you do if you had all the **data**?
What do you want to **predict** or **estimate**?

Get the data.

How were the data **sampled**?
Which data are **relevant**?
Are there **privacy** issues?

Explore the data.

**Plot** the data.
Are there **anomalies**?
Are there **patterns**?

Model the data.

**Build** a model.
**Fit** the model.
**Validate** the model.

Communicate and visualize the results.

What did we **learn**?
Do the results make **sense**?
Can we tell a **story**?

# Statistical Data Distributions

Every observed random variable has a particular frequency/probability distribution.

Some distributions occur often in practice/theory:

- The Binomial Distribution
- The Normal Distribution
- The Power Law Distribution

# Binomial Distributions

Experiments consist of *n identical, independent* trials which have two possible outcomes, with probabilities *p* and *(1-p)* like heads or tails.
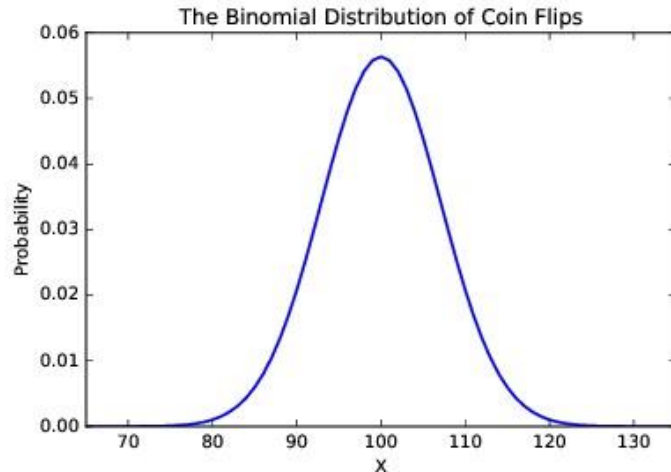
$$P\{X = x\} = \binom{n}{x} p^x (1 - p)^{n-x}$$

# Properties of Binomial Distributions

Discrete, but bell (or half-bell) shaped

Coin flips:  p=0.5     n=200

Lightbulb burnouts: p=0.001 n=1000

The distribution is a function of n and p.

# The Normal Distribution

The bell-shaped distribution of height, IQ, etc.

Completely parameterized by mean and standard deviation:

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

Not all bell-shaped distributions are normal but it is generally a reasonable start.

# Interpreting the Normal Distribution

Tight bounds on probability follow for Z-scores from normally distributed random variables:

IQ is normally distributed, with mean 100 and standard deviation 15.

Thus about 2.5% of people have IQs above 130.

# Power Law Distributions

Power laws are defined $P(X = x) = cx^{-\alpha}$ for exponent $\alpha$ and normalization constant $c$.

They do not cluster around a mean like a normal distribution, instead having very large values rarely but consistently.

They define 80-20 rules: 20% of the $X$ get 80% of the $Y$.

# City Population Yield Power Laws

The average big US city has population 165,719.   Even with a huge standard deviation of 410,730, New York city with 8,008,278 people is too many sigma away from the mean.

Power laws arise when the rich get richer.

# Linear and Log-Log Plots for City Pop

Straight lines on log-log plots say power law.

The biggest values are out of scale on linear plots.

# Many Distributions are Power Laws

- Internet sites with x inlinks.
- Frequency of earthquakes at x on the Richter scale
- Words used with a relative frequency of x
- Wars which kill x people

Power laws show as straight lines on log value, log frequency plots.

# When is an Observation Meaningful?

Computational analysis readily finds patterns and correlations in large data sets.

But when is a pattern significant?

Sufficiently strong correlations on large data sets may seem ``obviously'' significant, but the issues are often quite subtle.

# Comparing Population Means

The T-test evaluates whether the population means of two samples are different.

Sample the IQs of 20 men and 20 women.  Is one group smarter on average?

Certainly the sample means will differ, but is this difference significant?

# Differences in Distributions

It becomes easier to distinguish two distributions as the means move apart...



... or the variance decreases:

# The T-Test

Two means differ significantly if:

- The mean difference is relatively large
- The standard deviations are small enough
- The samples are large enough

Welch's t-statistic is:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

where $s^2$ is the sample variance.

Significance is looked up in a table.



Low variability

Medium variability

High variability

# Permutation Tests and P-values

If your hypothesis is supported by the data, then randomly shuffled data sets should be less likely to support it.

The ranking of the real test statistic among the shuffled test statistics gives a p-value.

You need statistic on your model you believe is interesting, e.g. correlation, std. error, etc.

# Significance of a Permutation Test

The rank of the real data among the random permutations determines significance:
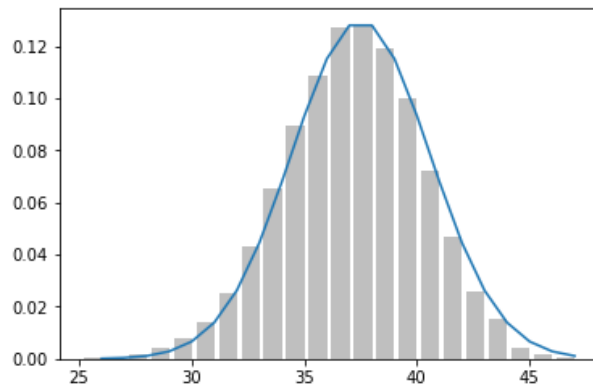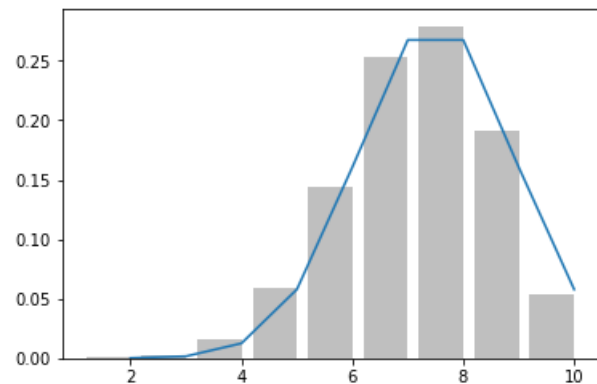
# The Central Limit Theorem

- A random variable defined as the average of a large #of independent and identically distributed(i.i.d.) random variables is itself approximately normally distributed.

- If $x_1, \ldots, x_n$ are r.v. with $\mu$ and $\sigma^2$, and if n is large:

  $Z = 1/n \ (x_1 + \ldots + x_n)$ is approx. normally distributed

# Significance of central limit theorem

- If n gets large, Binomial(n, p) ~ Normal(np, np(1-p))

# Significance of central limit theorem

- It implies that probabilistic and statistical methods that work for normal distributions can be applicable to many problems involving other types of distributions.

# Statistical Hypothesis Testing

- Example: Flipping a Coin – when speculating the coin is not fair.
  - **null hypothesis ($H_0$):** coin is fair, i.e., p = 0.5
  - **alternative hypothesis ($H_1$):** coin is not fair.

- We use statistics to decide whether we can reject $H_0$ as false or not.
  - In particular, flipping the coin $n$ times and counting the #of heads $X$.
  - Each coin flip is a **Bernoulli** trial, meaning X is a **Binomial**(n,p).
  - Due to CLT, X can be approximated by **Normal**(np, np(1-p)).
  - Choose *significance* level– how willing to make a *type I error* (FP)
  - *Typical choices: 5% or 1%*

```
normal_two_sided_bounds(0.95, 500, 15.81) (469.01026640487555, 530.9897335951244)
normal_two_sided_bounds(0.99, 500, 15.81) (459.27260472187146, 540.7273952781286)
```

pdf of Normal(500, 15.81)

# Types of errors

| Table of error types | | Null hypothesis ($H_0$) is | |
|---|---|---|---|
| | | True | False |
| Decision About Null Hypothesis ($H_0$) | Reject | Type I error (False Positive) | Correct inference (True Positive) |
| | Fail to reject | Correct inference (True Negative) | Type II error (False Negative) |

Type I error is detecting an effect that is not present, while a type II error is failing to detect an effect that is present.

# Visualizing Data

# Typical Data Science Pipeline



Ask an interesting question.

What is the scientific **goal**?
What would you do if you had all the **data**?
What do you want to **predict** or **estimate**?

Get the data.

How were the data **sampled**?
Which data are **relevant**?
Are there **privacy** issues?

Explore the data.

**Plot** the data.
Are there **anomalies**?
Are there **patterns**?

Model the data.

**Build** a model.
**Fit** the model.
**Validate** the model.

Communicate and visualize the results.

What did we **learn**?
Do the results make **sense**?
Can we tell a **story**?

# Exploratory Data Analysis

"The greatest value of a picture is when it forces us to notice what we never expected to see."

John Tukey

# Exploratory Data Analysis

Looking carefully at your data is important:

- to identify mistakes in collection/processing
- to find violations of statistical assumptions
- to observe patterns in the data
- to make hypothesis.

Feeding unvisualized data to a machine learning algorithm is asking for trouble.

# Anscombe's Quartet

All four data sets have exactly the same mean, variance, correlation, and regression line:

| | I | | II | | III | | IV | |
|---|---|---|---|---|---|---|---|---|
| | x | y | x | y | x | y | x | y |
| | 10.0 | 8.04 | 10.0 | 9.14 | 10.0 | 7.46 | 8.0 | 6.58 |
| | 8.0 | 6.95 | 8.0 | 8.14 | 8.0 | 6.77 | 8.0 | 5.76 |
| | 13.0 | 7.58 | 13.0 | 8.74 | 13.0 | 12.74 | 8.0 | 7.71 |
| | 9.0 | 8.81 | 9.0 | 8.77 | 9.0 | 7.11 | 8.0 | 8.84 |
| | 11.0 | 8.33 | 11.0 | 9.26 | 11.0 | 7.81 | 8.0 | 8.47 |
| | 14.0 | 9.96 | 14.0 | 8.10 | 14.0 | 8.84 | 8.0 | 7.04 |
| | 6.0 | 7.24 | 6.0 | 6.13 | 6.0 | 6.08 | 8.0 | 5.25 |
| | 4.0 | 4.26 | 4.0 | 3.10 | 4.0 | 5.39 | 19.0 | 12.50 |
| | 12.0 | 10.84 | 12.0 | 9.13 | 12.0 | 8.15 | 8.0 | 5.56 |
| | 7.0 | 4.82 | 7.0 | 7.26 | 7.0 | 6.42 | 8.0 | 7.91 |
| | 5.0 | 5.68 | 5.0 | 4.74 | 5.0 | 5.73 | 8.0 | 6.89 |
| mean | 9.0 | 7.5 | 9.0 | 7.5 | 9.0 | 7.5 | 9.0 | 7.5 |
| var. | 10.0 | 3.75 | 10.0 | 3.75 | 10.0 | 3.75 | 10.0 | 3.75 |
| corr. | 0.816 | | 0.816 | | 0.816 | | 0.816 | |

# Plotting Anscombe's Quartet

All four data sets have exactly the same mean, variance, correlation, and regression line:

# Mapping Data to Image

# Most Effective

# Less Effective

# Least Effective



Estimated fraction of precipitation lost to evapotranspiration 1971-2000

| | | |
|---|---|---|
| 0.0 - 0.09 | 0.5 - 0.59 | 1.0 - 1.09 |
| 0.1 - 0.19 | 0.6 - 0.69 | 1.1 - 1.19 |
| 0.2 - 0.29 | 0.7 - 0.79 | 1.2 - 1.29 |
| 0.3 - 0.39 | 0.8 - 0.89 | |
| 0.4 - 0.49 | 0.9 - 0.99 | |

# Order These Values

# Perceived as Ordered

| Brightness | Saturation | Hue: not as much |
|:---:|:---:|:---:|

# Examples: Not Effective



Sources: US Treasury and WHO reports

# Examples: Not Effective

# Examples: Much Better

# **Tufte's Design Principle**

Distinguishing good/bad visualizations requires a design aesthetic, and a vocabulary to talk about data representations:

- Maximize data ink-ratio
- Minimize lie factor
- Minimize chartjunk
- Use proper scales and clear labeling

# Maximize Data-Ink Ratio

$$\text{Data-Ink Ratio} = \frac{\text{Data ink}}{\text{Total ink used in graphic}}$$

# Maximize Data-Ink Ratio

$$\text{Data-Ink Ratio} = \frac{\text{Data ink}}{\text{Total ink used in graphic}}$$

# The Lie Factor

$$\frac{\text{Size of effect shown in graphic}}{\text{Size of effect in data}}$$



THE SHRINKING FAMILY DOCTOR
In California

Percentage of Doctors Devoted Solely to Family Practice

| 1964 | 1975 | 1990 |
|------|------|------|
| 27% | 16.0% | 12.0% |

1: 4,232
6,212

1: 3,167
6,694

1: 2,247 RATIO TO POPULATION
8,023 Doctors



IN THE BARREL...
Price per bbl. of
light crude, leaving
Saudi Arabia
on Jan. 1

April 1
$14.55

$13.34

$12.70

$12.09

$11.51

$10.46

$10.95

$2.41

'73  '74  '75  '76  1977  1978  1979

# Reduce Chartjunk

Extraneous visual elements distract from the message the data is trying to tell.

- Extra dimensionality
- Uninformative coloring
- Excessive grids and figurative decoration

In an exciting graphic, the data tells the story, not the chartjunk.

# Can you Simplify this Plot?



Tim Brey

# Can You Further Simplify?

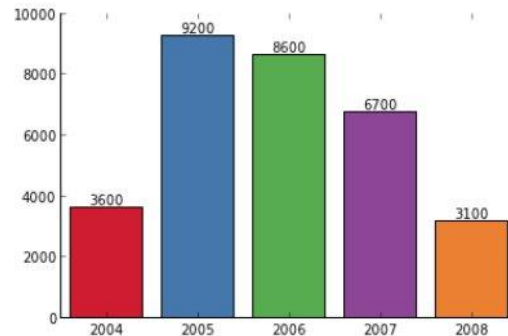# Better, but can you Further Simplify?

# Anything Else that Can Go?

# ``Less is More"

# MatPlotLib Supports Nice Plots



http://nbviewer.ipython.org/5357268

# Graphical Integrity: Scale Distortion

# Graphical Integrity: Scale Distortion

Always start bar graphs at zero.

# Scale Distortions



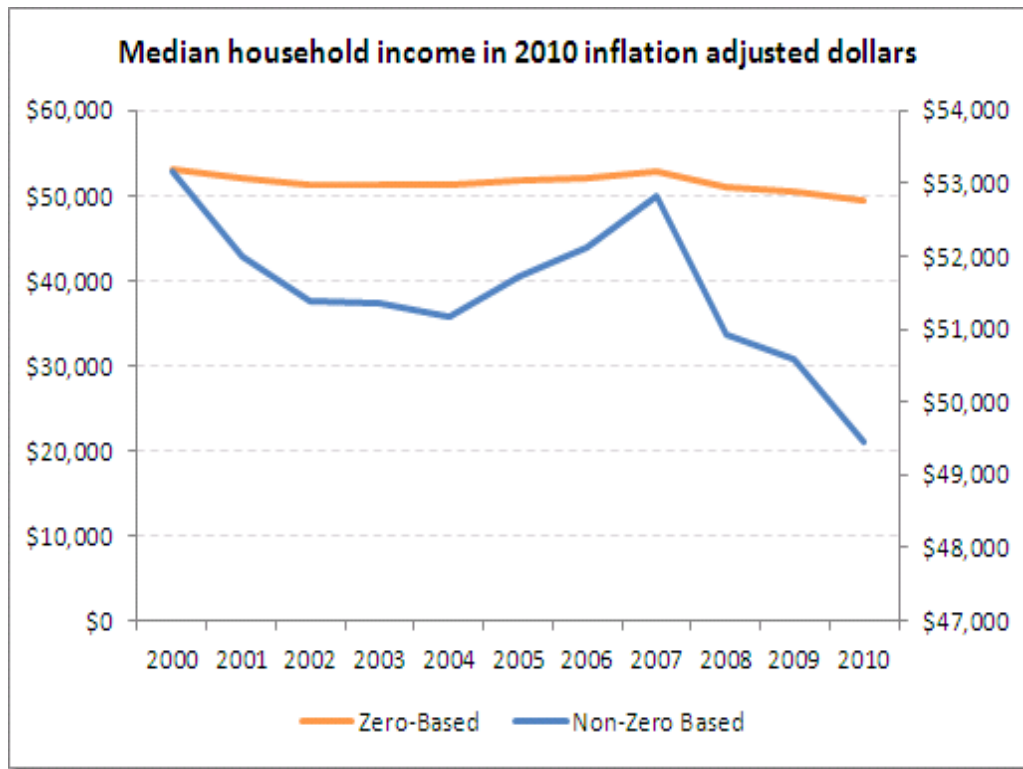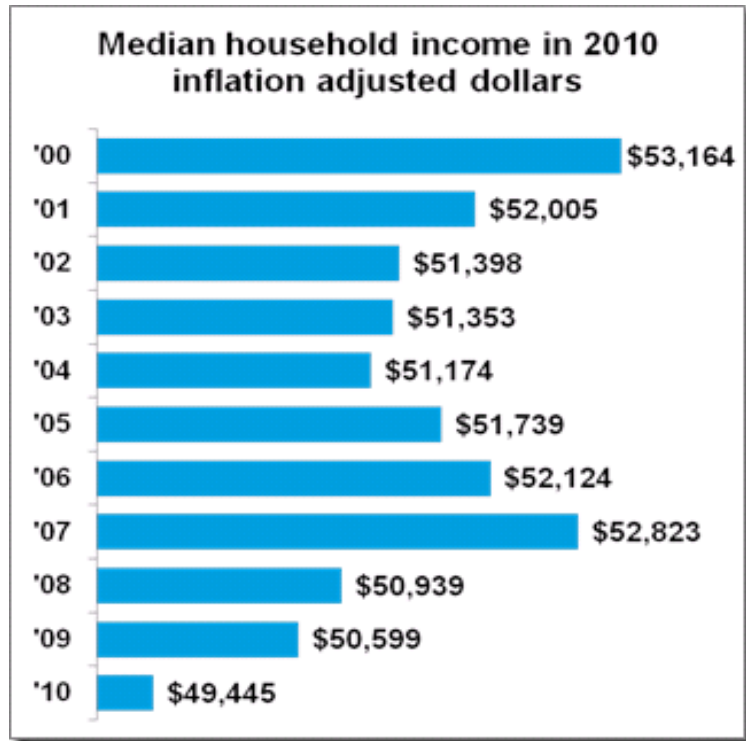HOW 2012 STACKS UP
THE WARMEST YEARS ON RECORD
CONTIGUOUS U.S.

Source: NOAA's National Climatic Data Center - State of the Climate National Overview
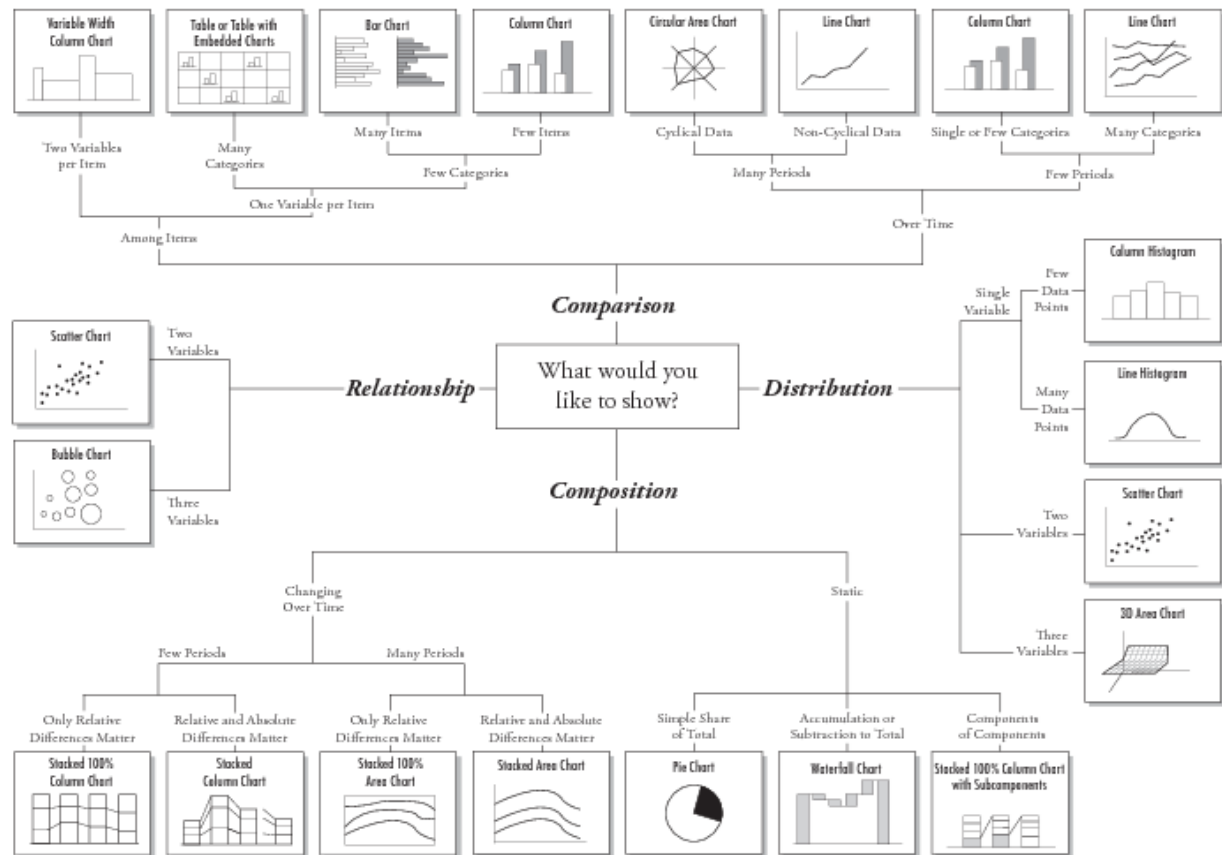CLIMATE CENTRAL

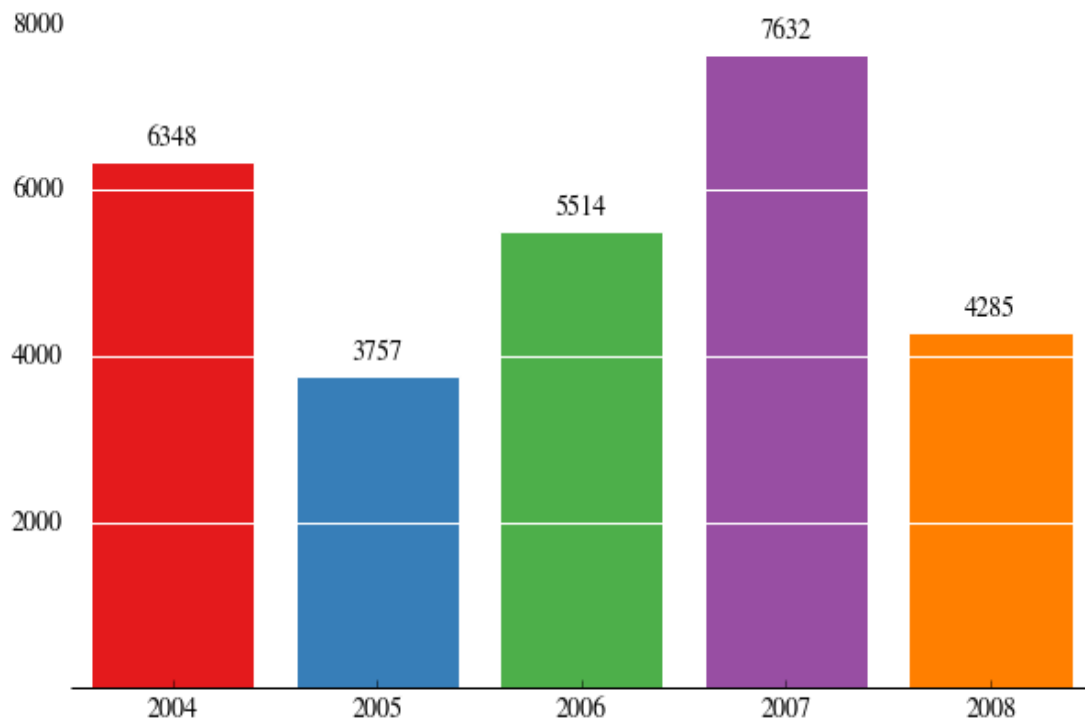# Scale Distortions

*Always* start your bar graphs at zero!



Median household income in 2010 inflation adjusted dollars

'00 — $53,164
'01 — $52,005
'02 — $51,398
'03 — $51,353
'04 — $51,174
'05 — $51,739
'06 — $52,124
'07 — $52,823
'08 — $50,939
'09 — $50,599
'10 — $49,445

Median household income in 2010 inflation adjusted dollars

Zero-Based | Non-Zero Based

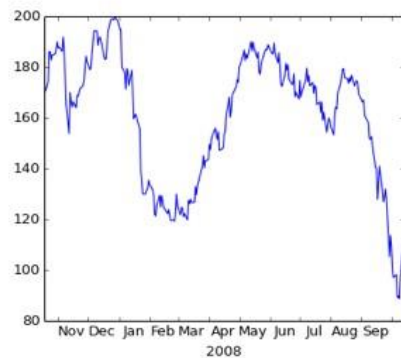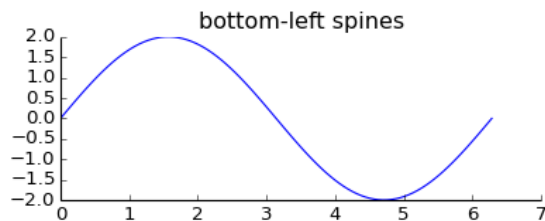# Which Chart to Use?



Chart Suggestions—A Thought-Starter

© 2006 A. Abela — a.v.abela@gmail.com

# Bar Chart

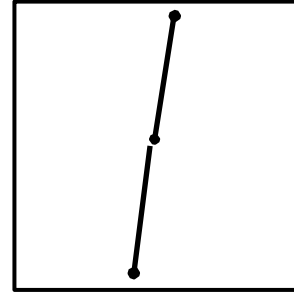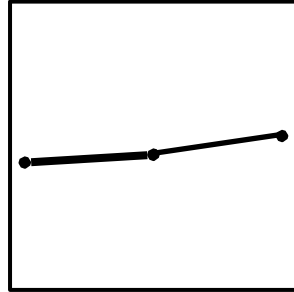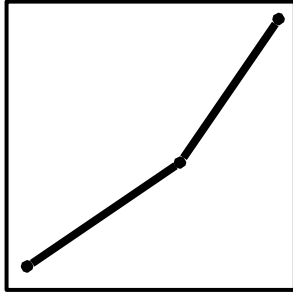# Line Charts - Trends Over Time

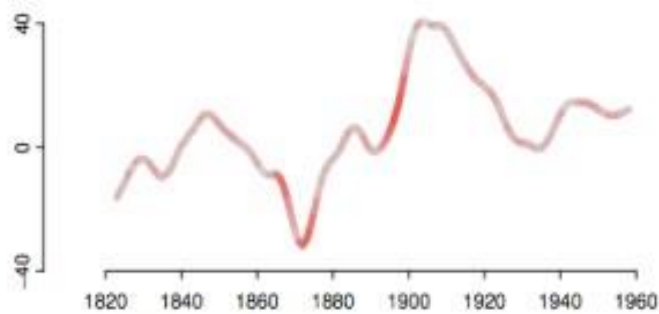

matplotlib gallery

# Aspect Ratios

# Banking to 45°

Two line segments are maximally discriminable when their average absolute angle is 45°
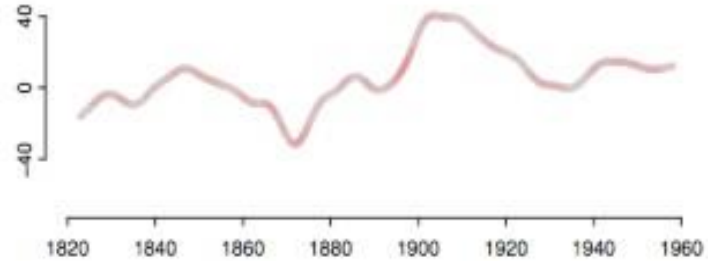
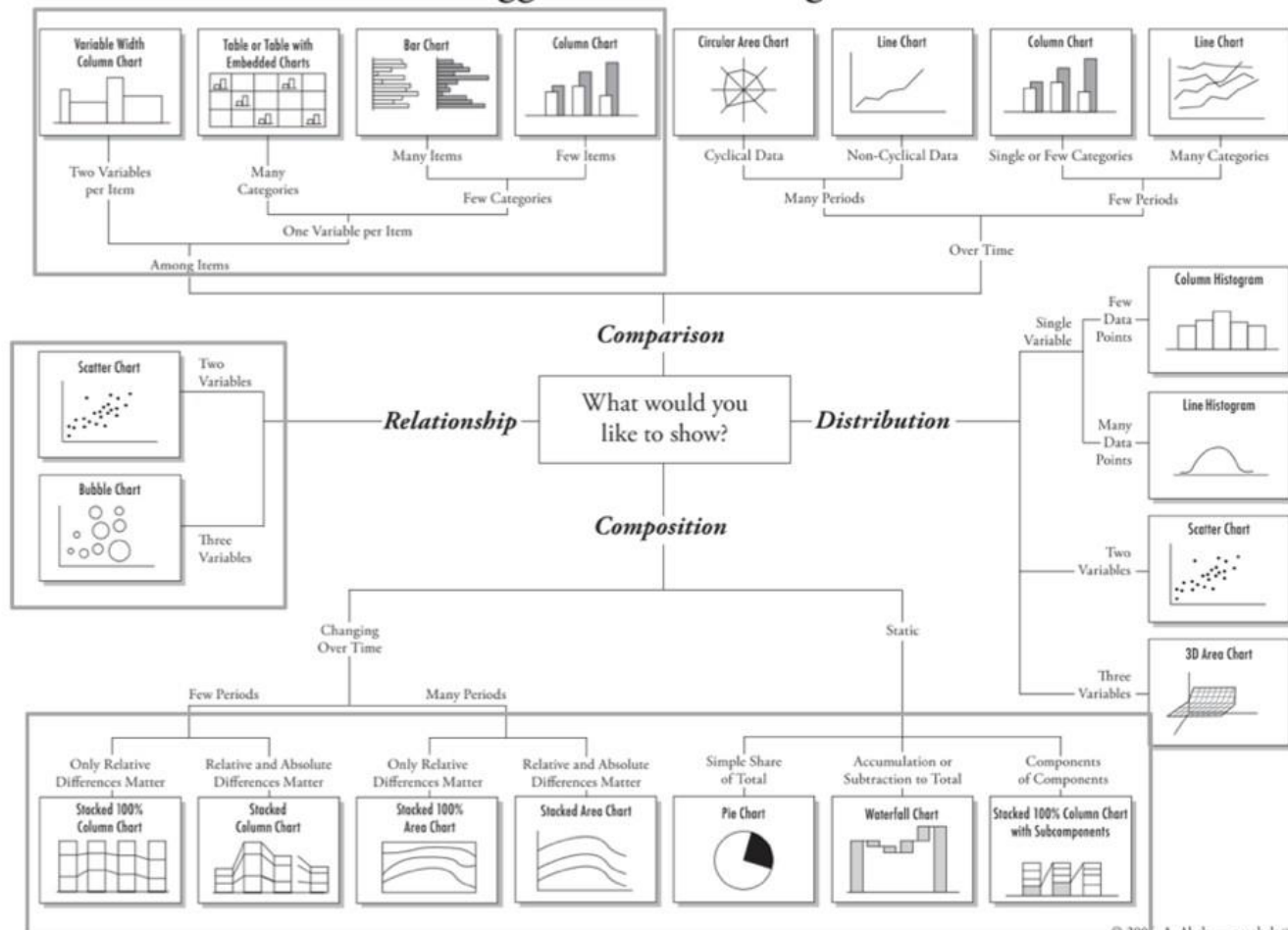W. Cleveland

# Banking to 45°



Error Prone

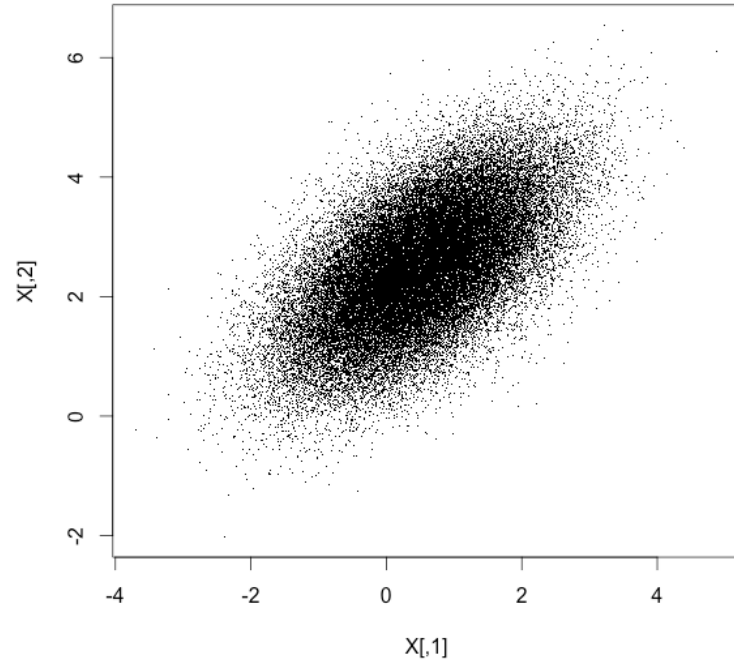Optimal Aspect Ratio

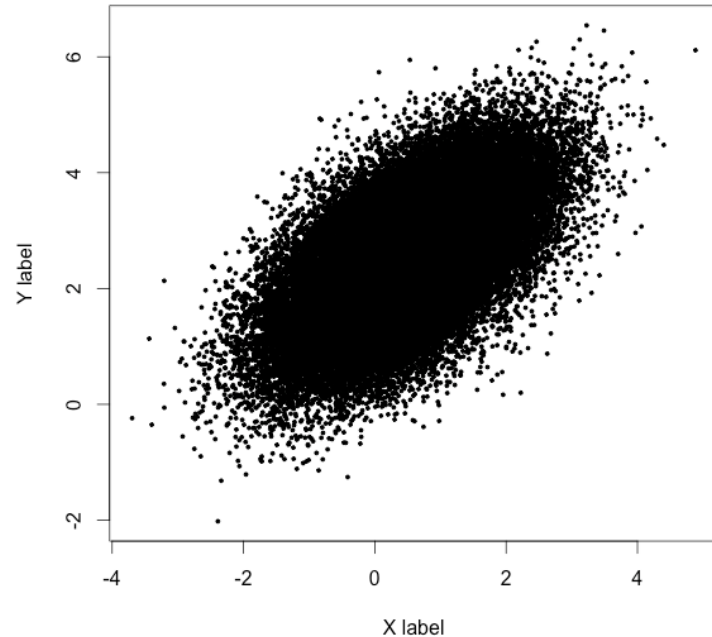# Chart Suggestions—A Thought-Starter

# **Scatter Plots / Bubble Charts**

- Scatter plots show the values of each point, and are a great way to present 2D data sets.

- For data sets with three or four variables, use bubble charts.

- Higher dimensional datasets can be projected to 2D through principle component analysis.
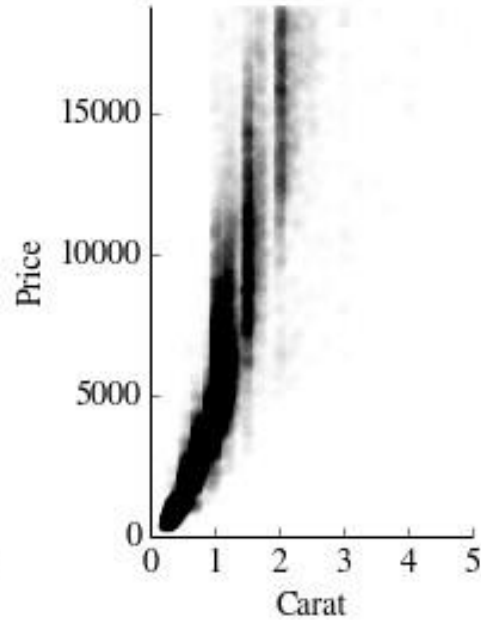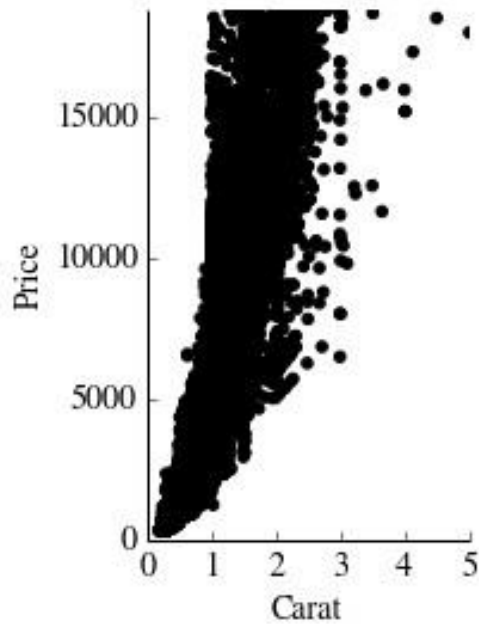
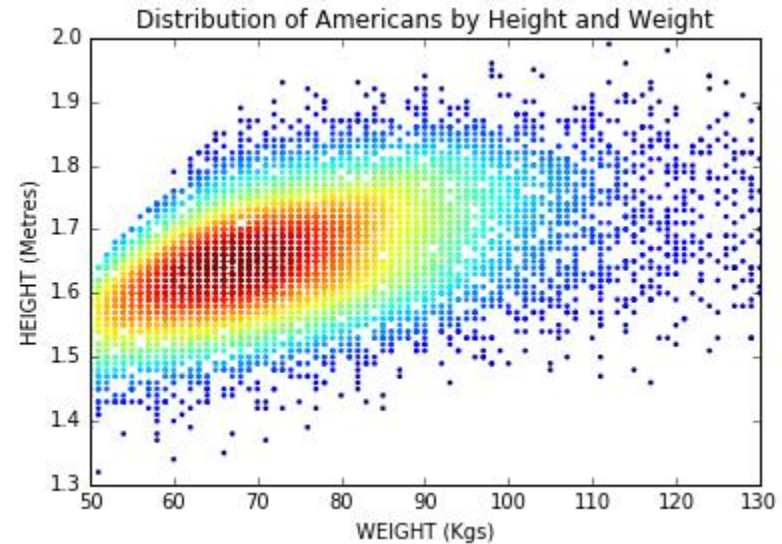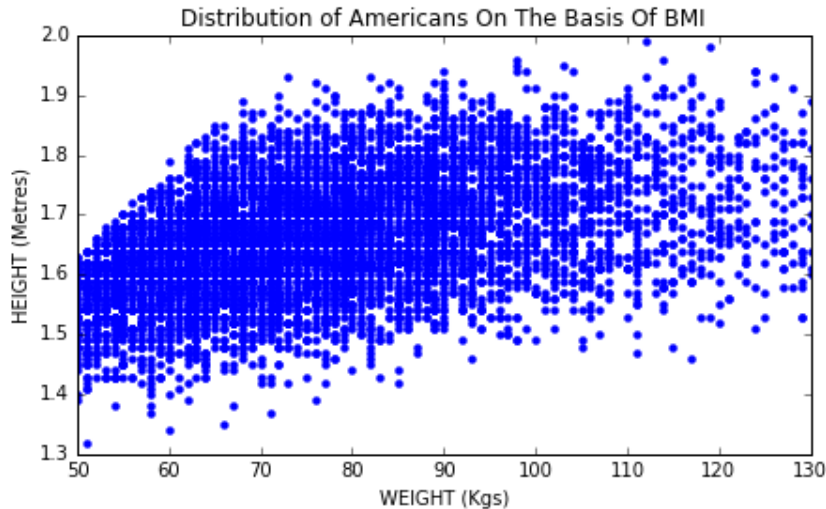# Reduce Overplotting by Small Points

# Reduce Overplotting by Opacity
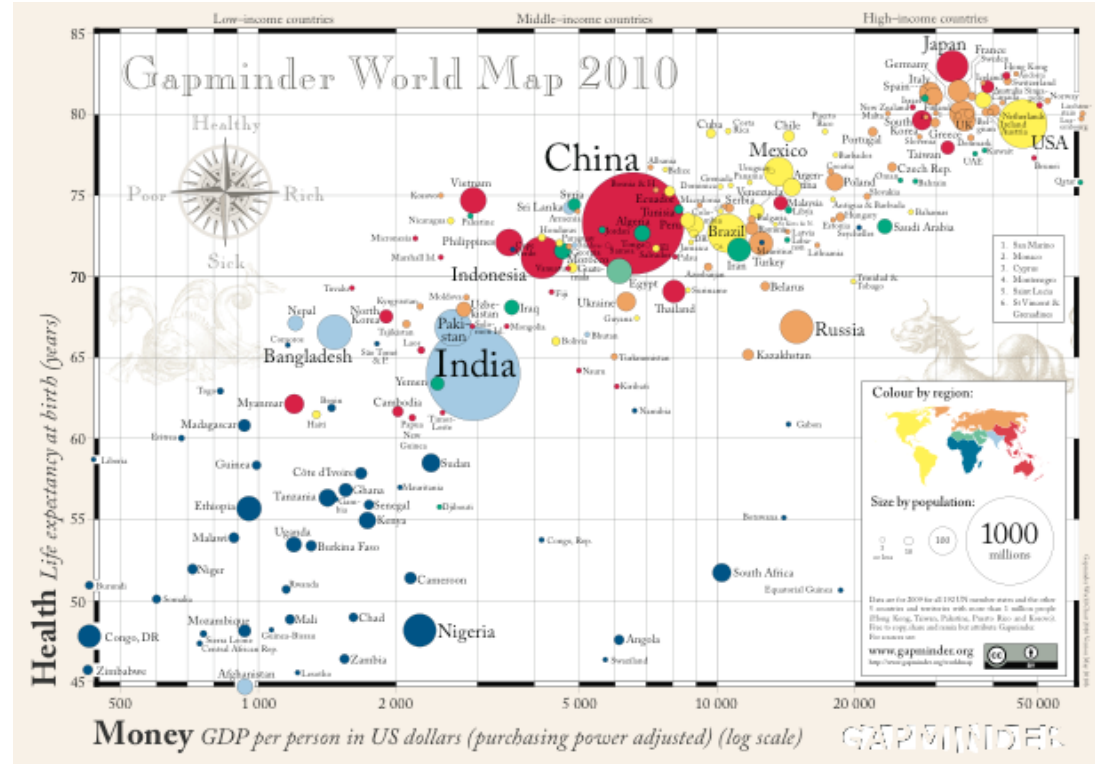


alpha = 1/100

# Heatmaps Reveal Finer Structure

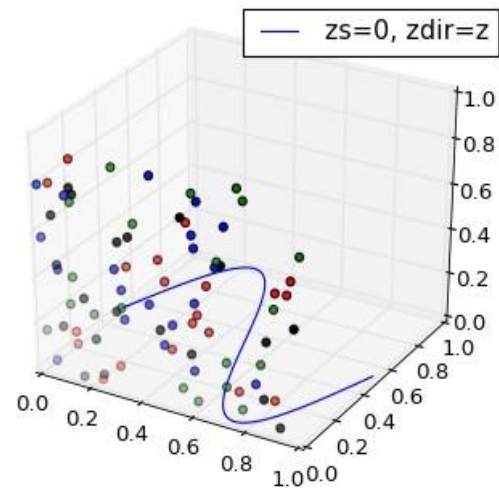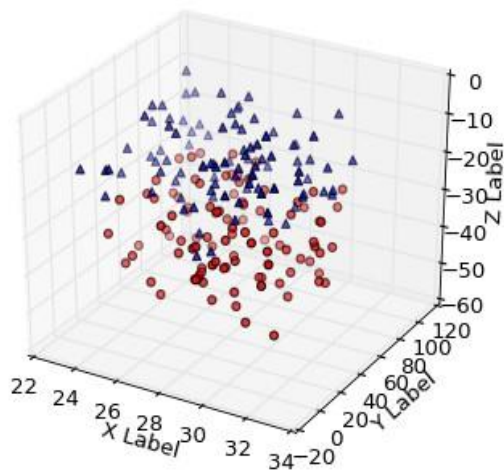Color points on the basis of frequency

# Bubble Charts for Extra Dimensions

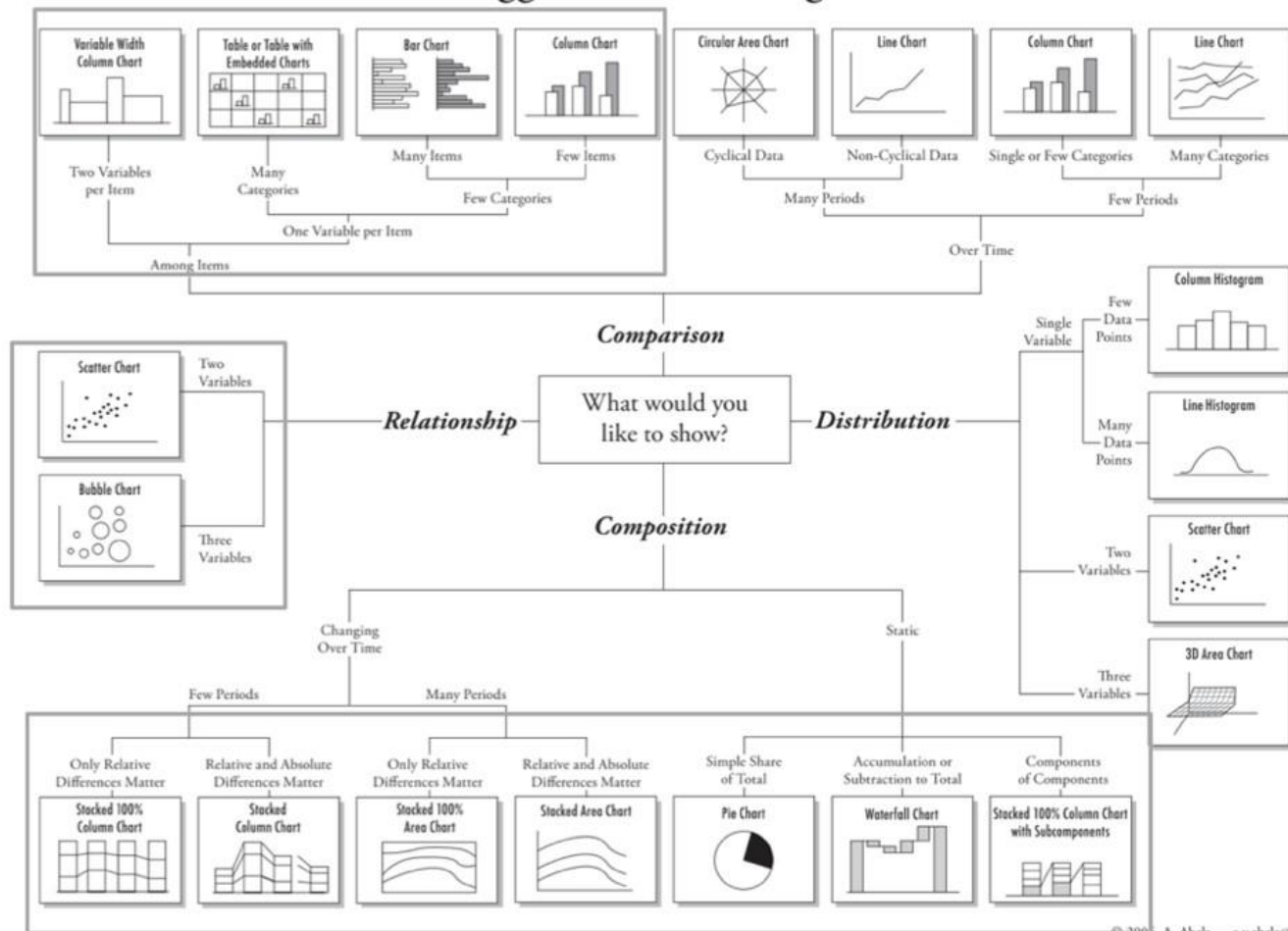Using color, shape, and size of "dots" enables dot plots to represent additional dimensions.

http://www.gapminder.org/videos/200-years-that-changed-the-world-bbc/

# Don't

# Chart Suggestions—A Thought-Starter

**Variable Width Column Chart** — Two Variables per Item

**Table or Table with Embedded Charts** — Many Categories

**Bar Chart** — Many Items

**Column Chart** — Few Items

Few Categories

One Variable per Item

Among Items

**Circular Area Chart** — Cyclical Data

**Line Chart** — Non-Cyclical Data

Many Periods

**Column Chart** — Single or Few Categories

**Line Chart** — Many Categories

Few Periods

Over Time

**Comparison**

**Scatter Chart** — Two Variables

**Bubble Chart** — Three Variables

**Relationship**

What would you like to show?

**Distribution**

Single Variable — Few Data Points — **Column Histogram**

Many Data Points — **Line Histogram**

Two Variables — **Scatter Chart**

Three Variables — **3D Area Chart**

**Composition**

Changing Over Time

Static

Few Periods:
Only Relative Differences Matter — **Stacked 100% Column Chart**
Relative and Absolute Differences Matter — **Stacked Column Chart**

Many Periods:
Only Relative Differences Matter — **Stacked 100% Area Chart**
Relative and Absolute Differences Matter — **Stacked Area Chart**

Simple Share of Total — **Pie Chart**

Accumulation or Subtraction to Total — **Waterfall Chart**

Components of Components — **Stacked 100% Column Chart with Subcomponents**

© 2006 A. Abela — a.v.abela@gmail.com

http://extremepresentation.typepad.com/blog/files/choosing_a_good_chart.pdf

# Pie vs. Bar Charts



65% of the market is controlled by companies B and C
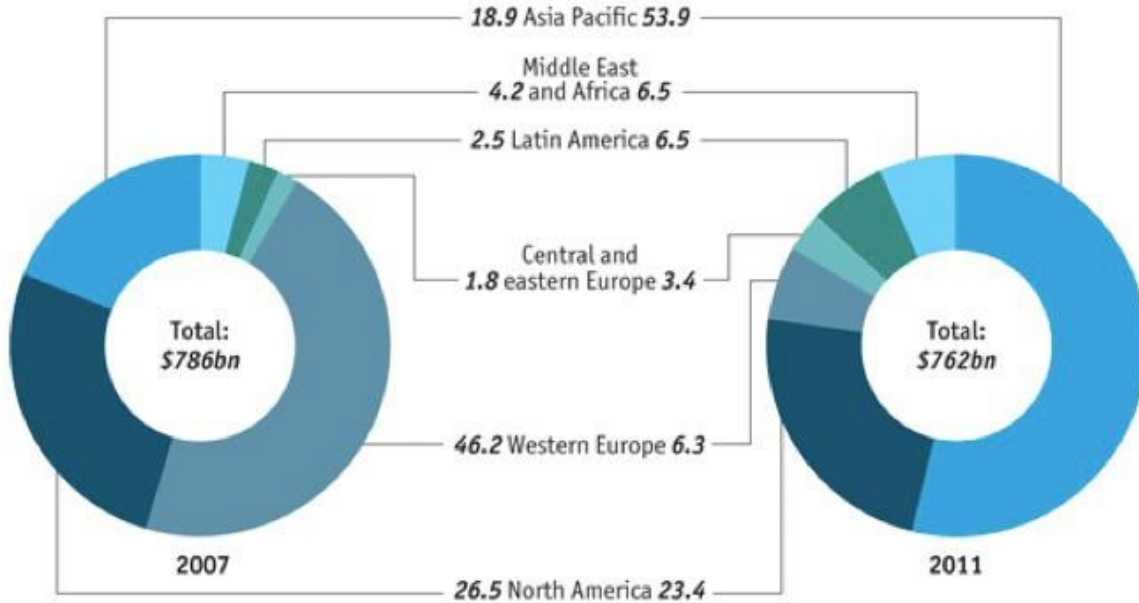
# Donut Chart



**Pre-tax profits of the 1,000 largest banks**
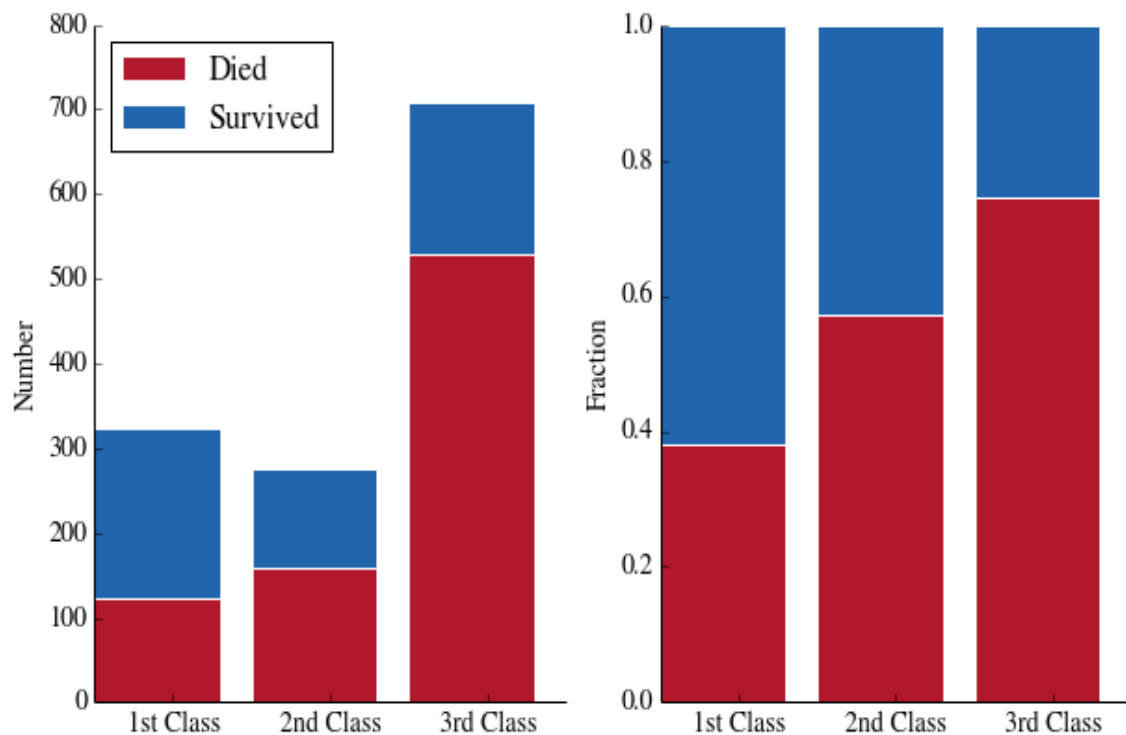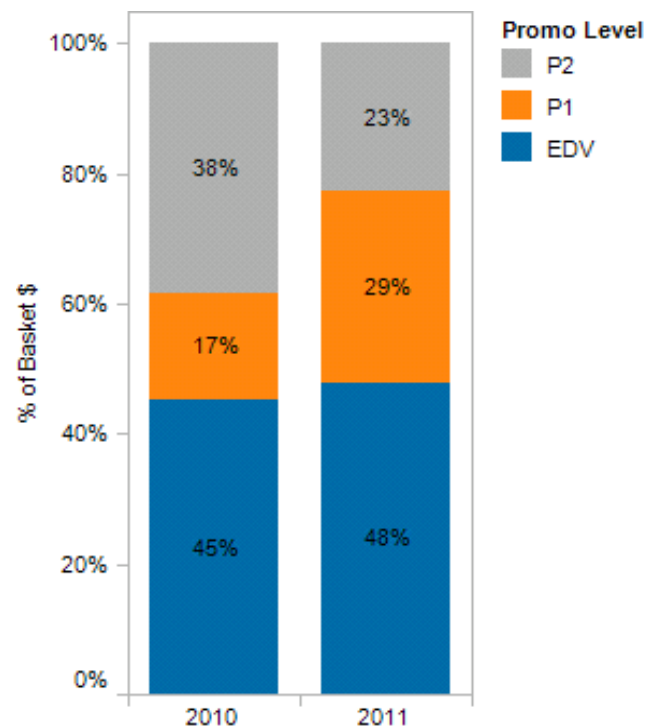By tier-one capital and domicile, % of total

18.9 Asia Pacific 53.9

Middle East
4.2 and Africa 6.5

2.5 Latin America 6.5

Central and
1.8 eastern Europe 3.4

Total: $786bn

Total: $762bn

46.2 Western Europe 6.3

2007

26.5 North America 23.4

2011

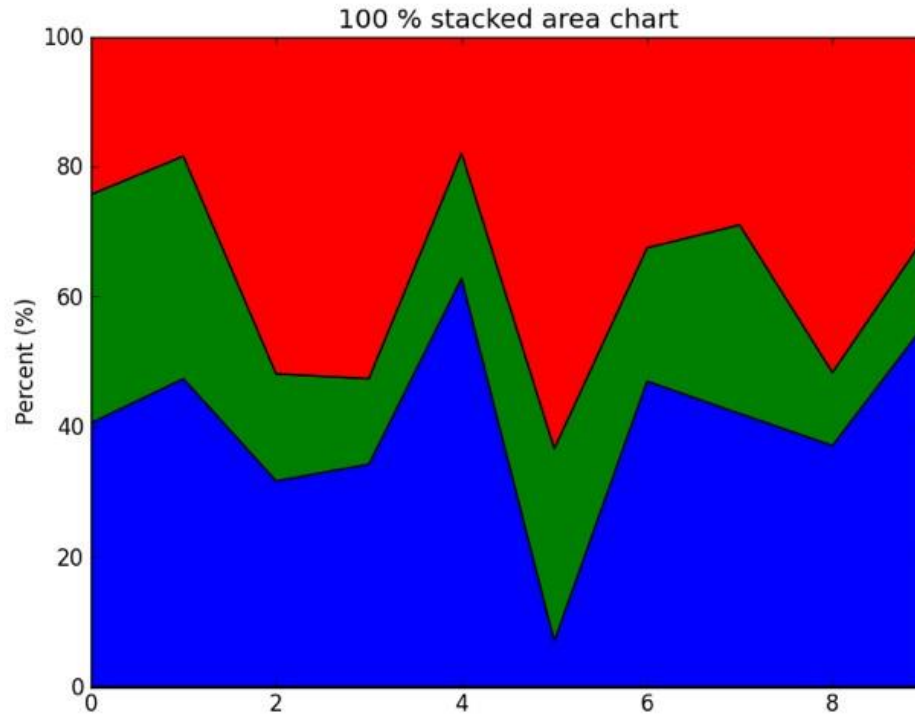Source: *The Banker Top 1000*

# Stacked Bar Chart

# Stacked Bar Chart



v.s.

# Stacked Area Chart
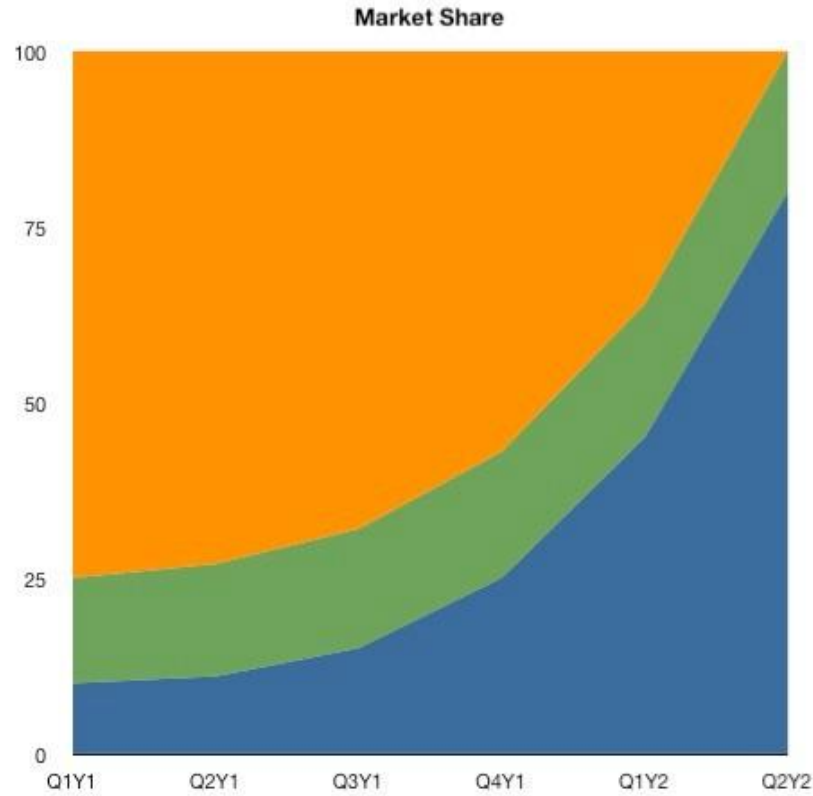
# 100% Stacked Area Chart
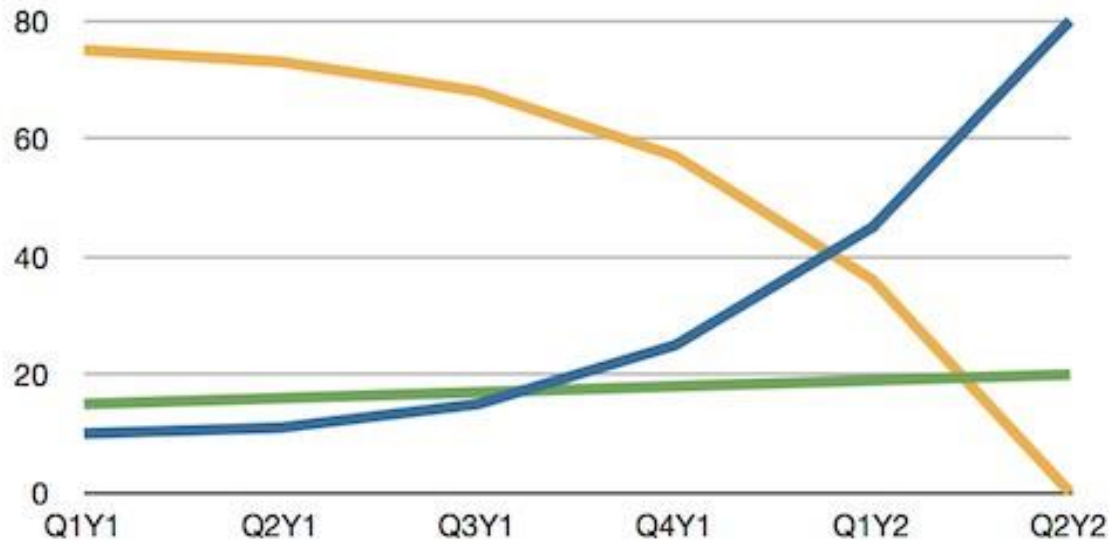


100 % stacked area chart

# Stacked Area vs. Line Graphs



Market Share
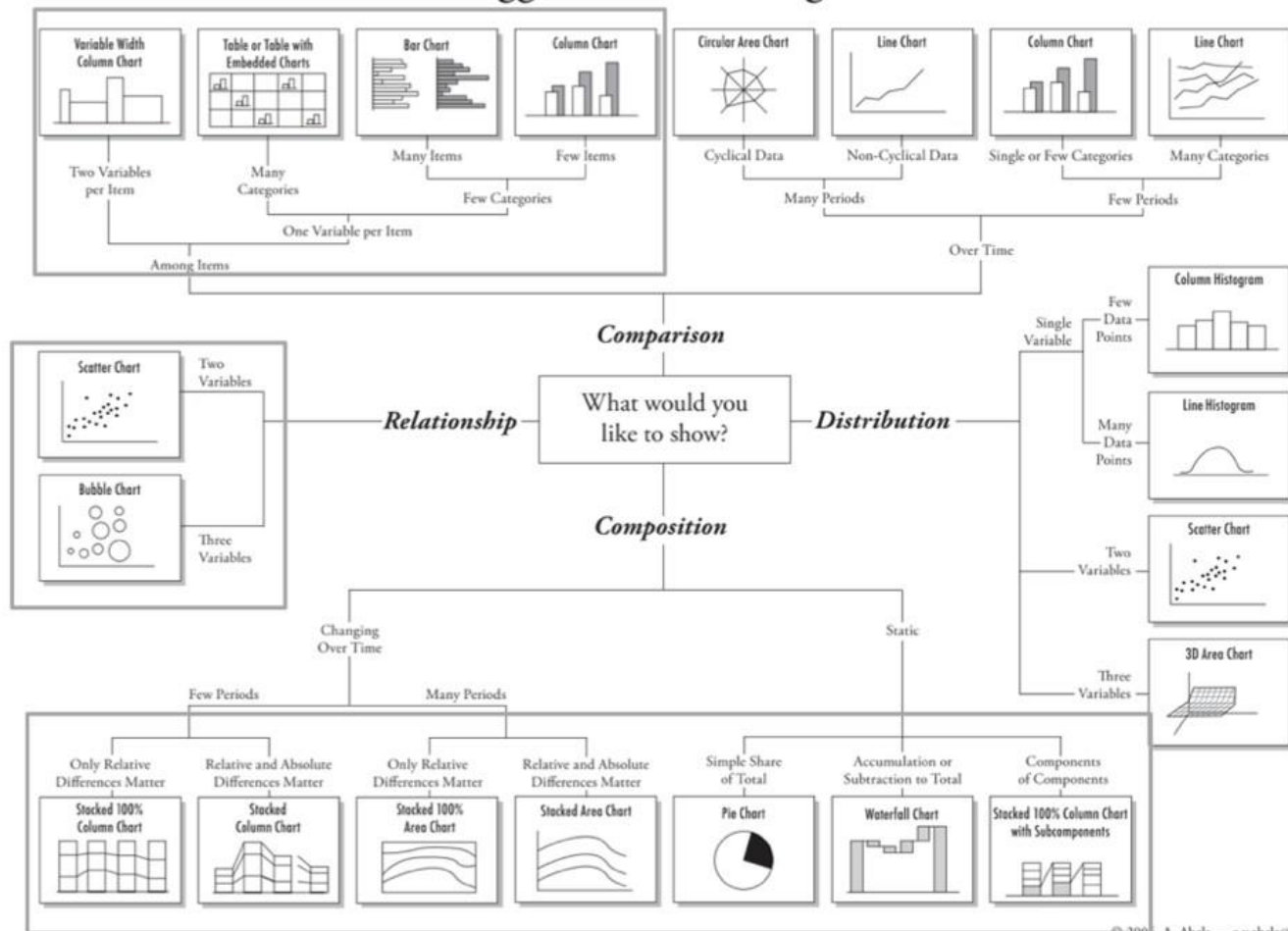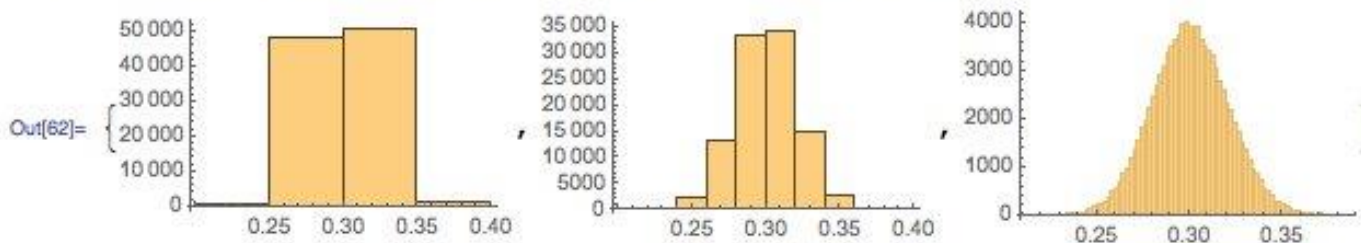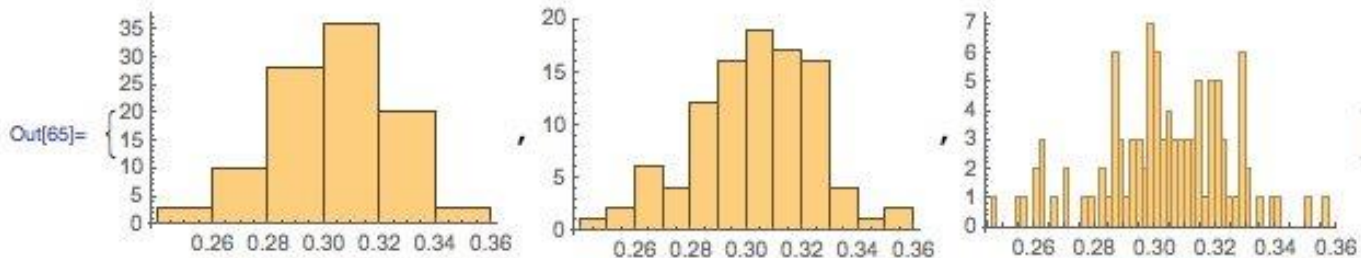
# Stacked Area vs. Line Graphs

# Chart Suggestions—A Thought-Starter

# Histograms: Bin Size / Count Matters



In[62]:= {Histogram[d, 5], Histogram[d, 10], Histogram[d, 100]}

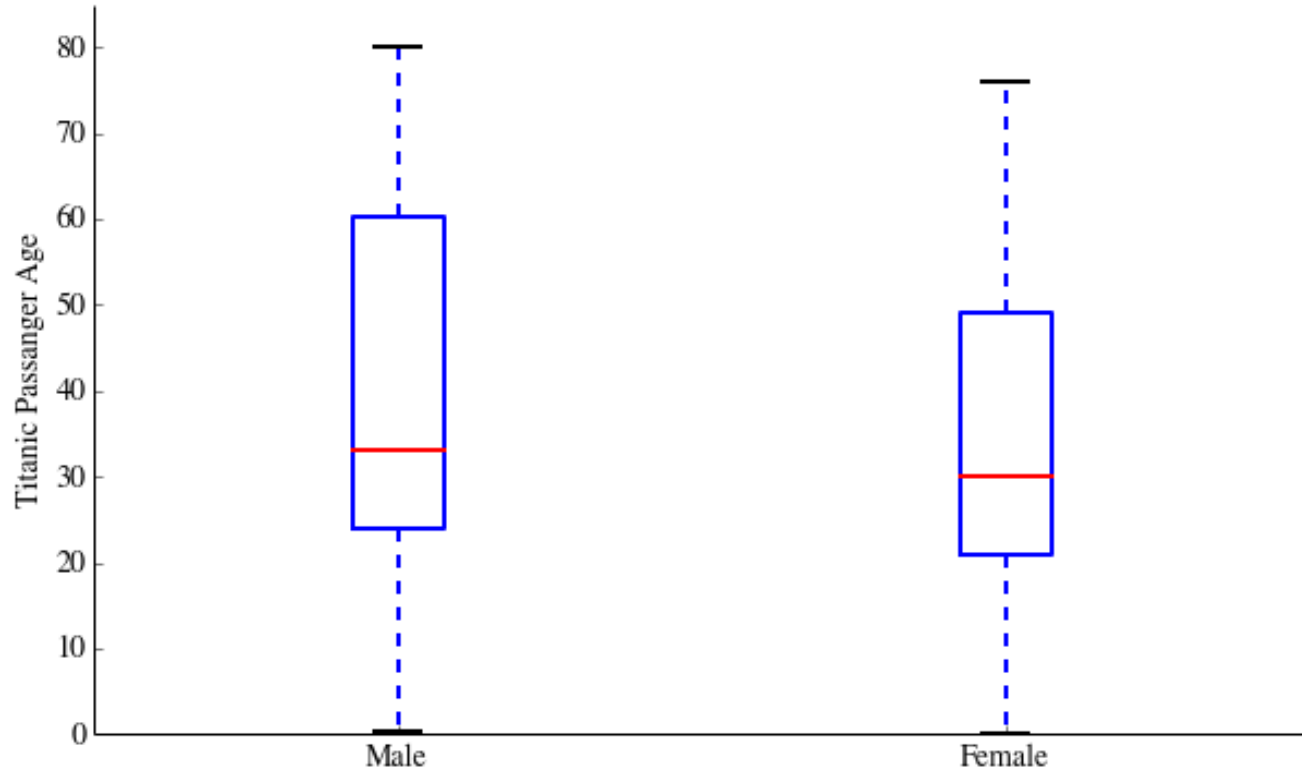d100 = Take[d, 100];

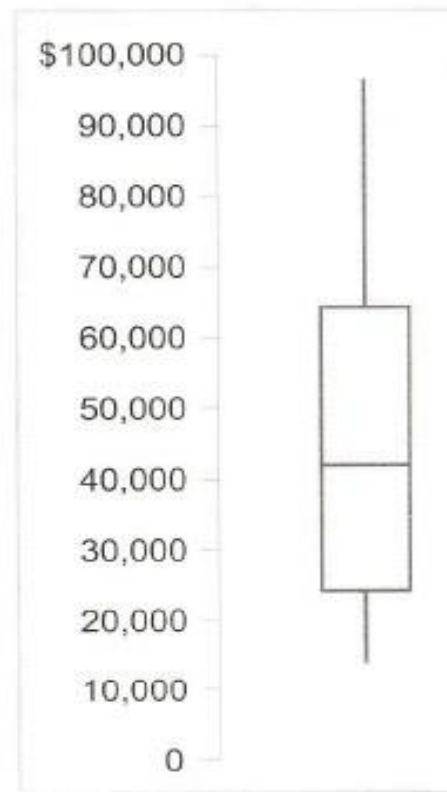In[65]:= {Histogram[d100, 5], Histogram[d100, 10], Histogram[d100, 100]}
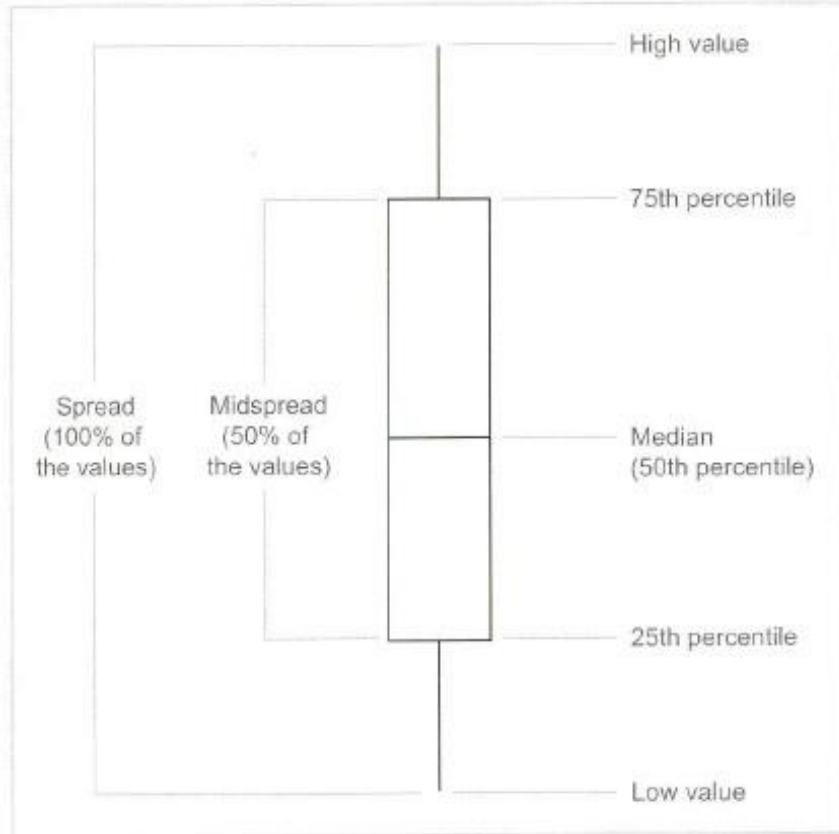
# Box & Whisker Plots

# Box & Whisker Plots

# Keep a Critical Eye

Remember Tufte's principles whenever  designing or interpreting data visualizations:

- Maximize data-ink ratio
- Minimize lie factor
- Minimize chartjunk
- Use proper scales and clear labeling

Beautiful data deserves beautiful visualization.