

# How Predictive Modelers Should use Data to Tell Data Stories

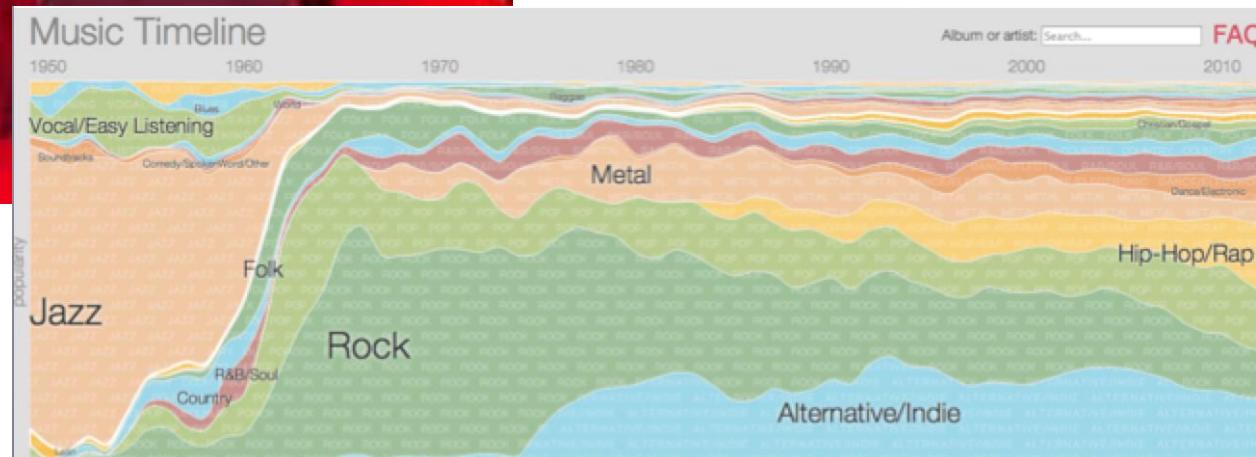
Dean Abbott

Co-Founder and Chief Data Scientist, SmarterHQ

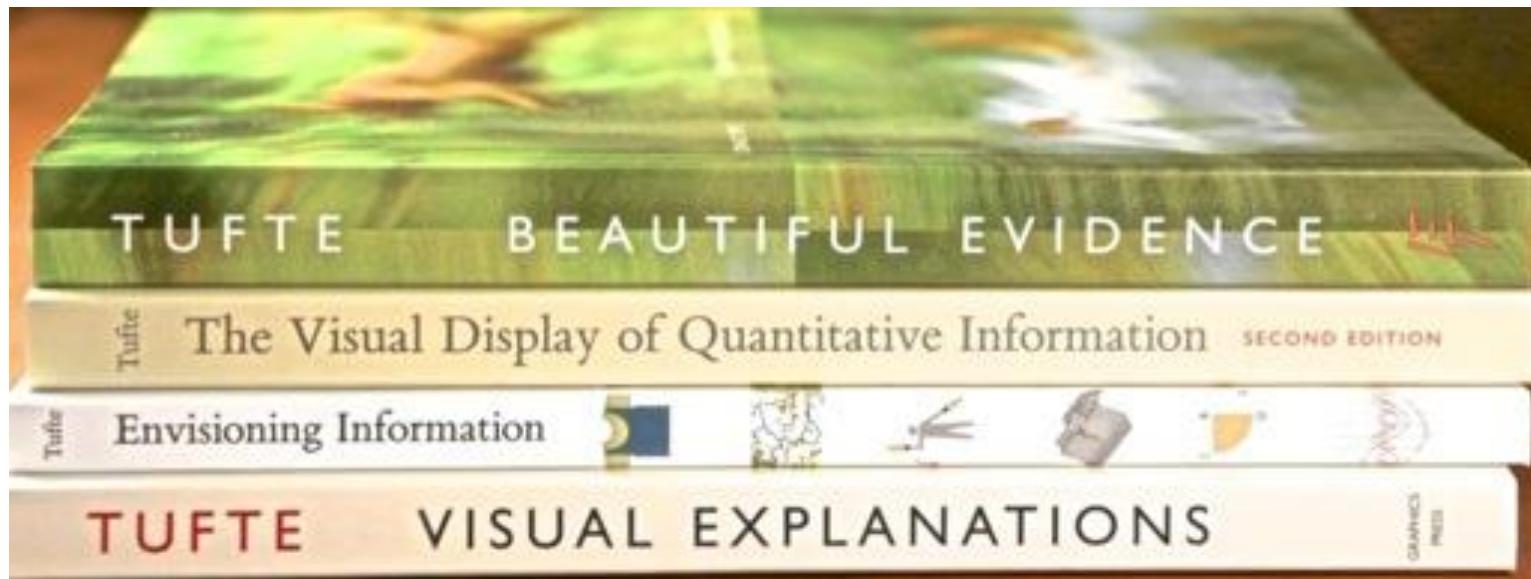
President, Abbott Analytics

Twitter: @deanabb

# Telling Data Stories is Not Just Visualization

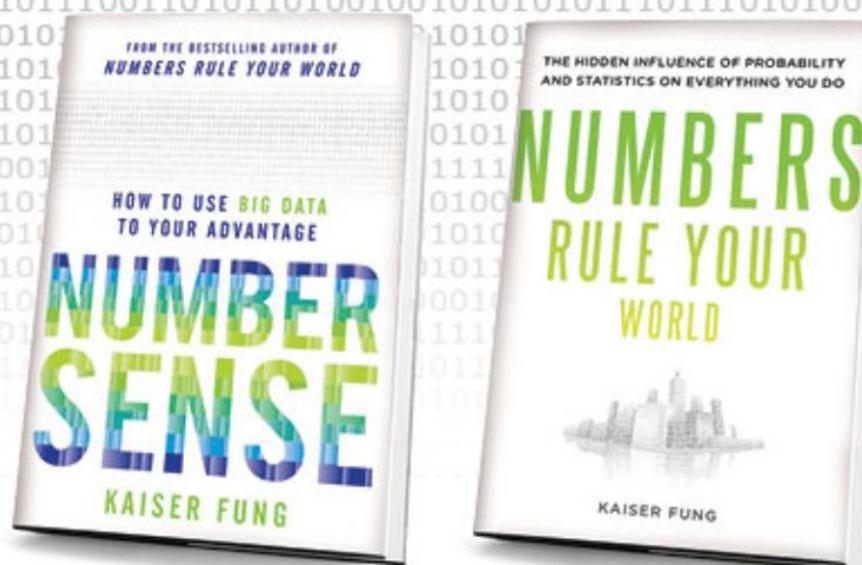


If you want in-depth  
Data Visualization Thought Leadership



<https://www.edwardtufte.com/tufte/>

If you want in-depth  
Data Visualization Thought Leadership



Kaiser Fung  
Big Data, Plainly Spoken

<http://www.kais erfung.com/blogs/>

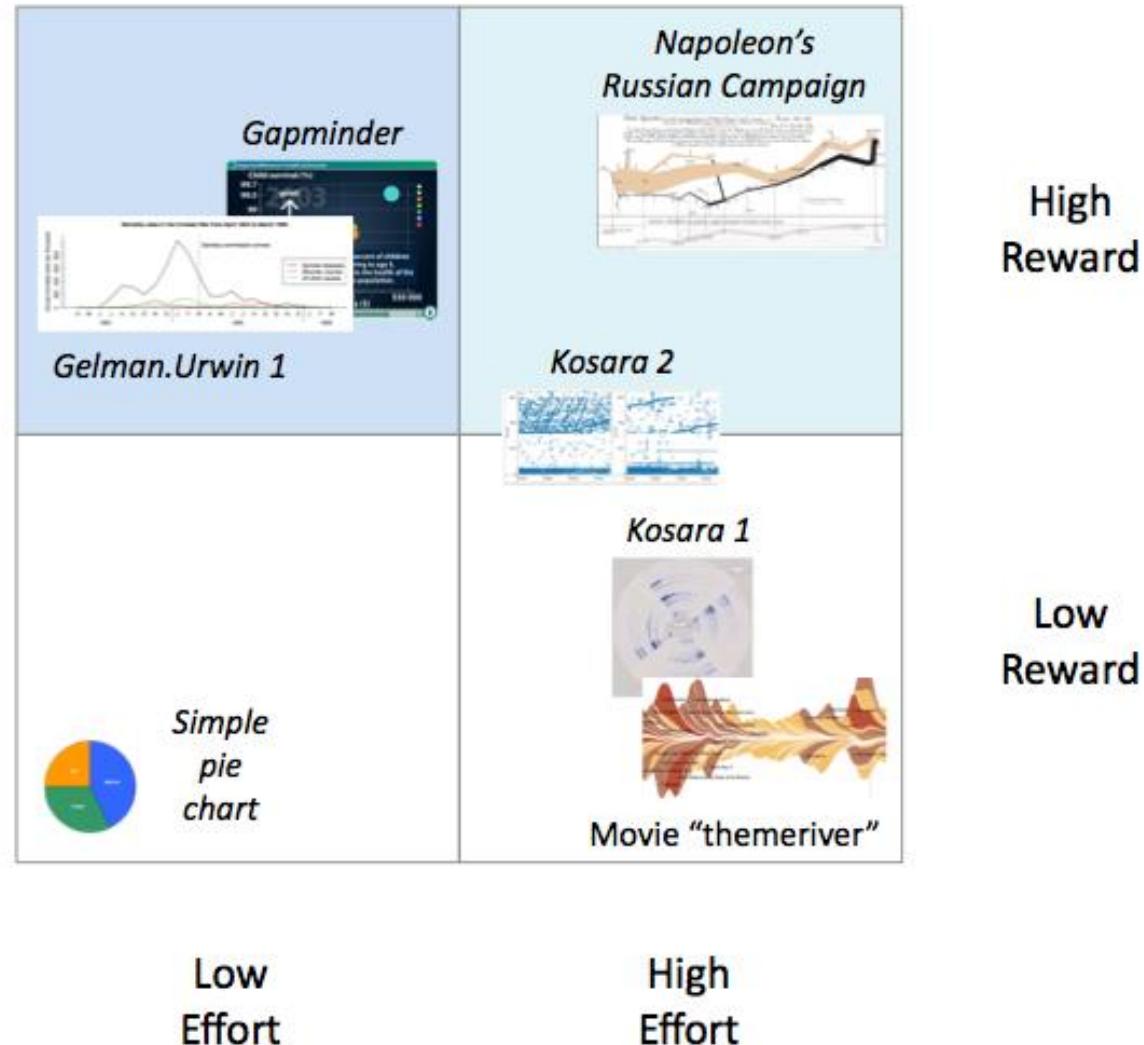
# Not All Visualization is Created Equally!

Fung: “Even acknowledged ‘perfect’ statistical graphics can require a high degree of effort.

[sometimes] the reward justifies the level of effort, and the return on effort is high.”

<https://statisticsforum.wordpress.com/2011/07/31/one-difference-between-statistical-graphics-and-infoviz-is-the-return-on-effort/>

## RETURN ON EFFORT MATRIX

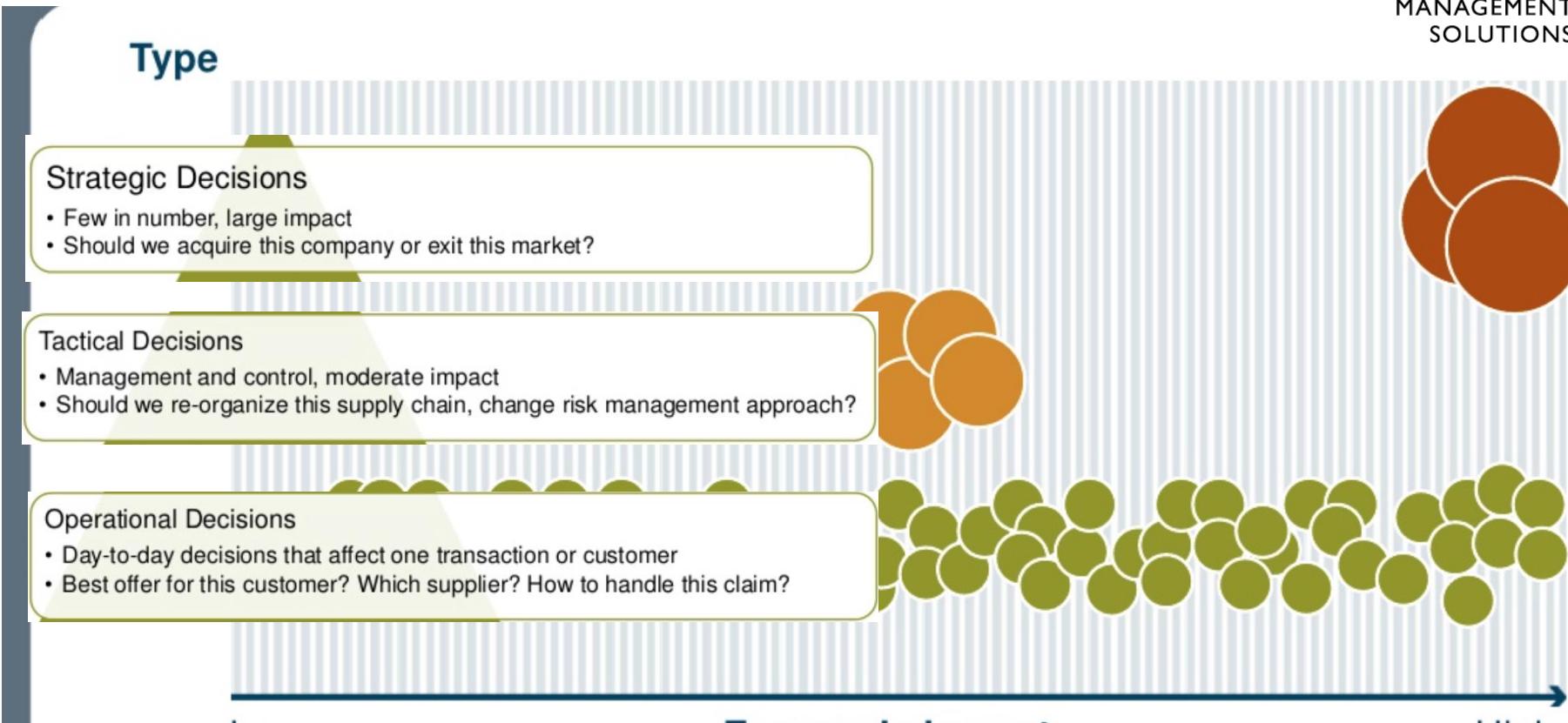


# We Don't Have Time to Create Masterpieces



<https://blog.udacity.com/2015/01/15-data-visualizations-will-blow-mind.html>

# James Taylor: Importance of (Frequent) Operational Decisions

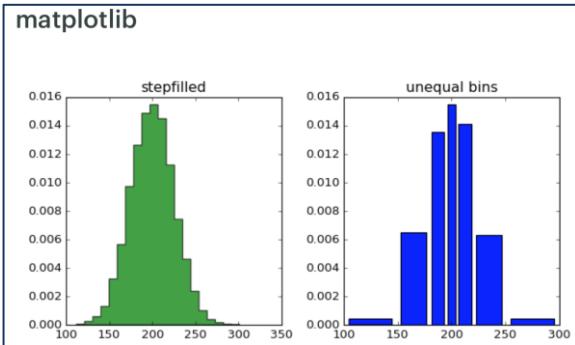


<https://www.slideshare.net/jamet123/10-best-practices-in-operational-analytics-6871966>

# Principle 1: Know Your Building Blocks (And Use the Good Ones!)



# Most of us Create These...



```
import numpy as np
import matplotlib.pyplot as plt
np.random.seed(0)

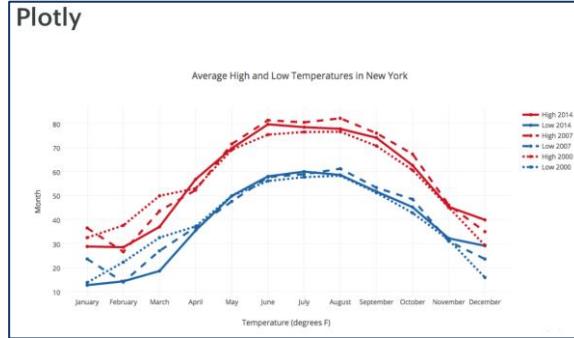
mu = 200
sigma = 25
x = np.random.normal(mu, sigma, size=100)

fig, (ax0, ax1) = plt.subplots(ncols=2, figsize=(8, 4))

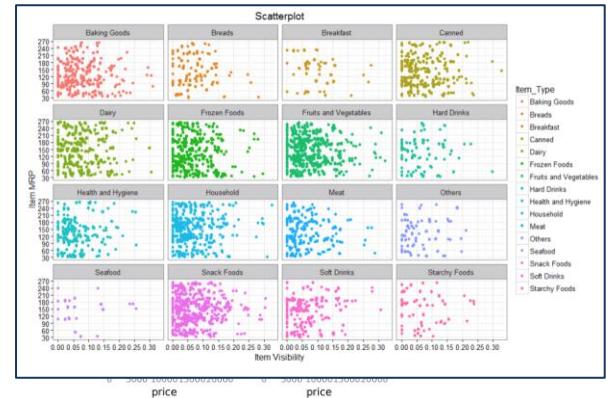
ax0.hist(x, 20, normed=1, histtype='stepfilled', facecolor='g', alpha=0.75)
ax0.set_title('stepfilled')

# Create a histogram by providing the bin edges (unequally spaced).
bins = [100, 150, 180, 195, 205, 220, 250, 300]
ax1.hist(x, bins, normed=1, histtype='bar', rwidth=0.8)
ax1.set_title('unequal bins')

fig.tight_layout()
plt.show()
```

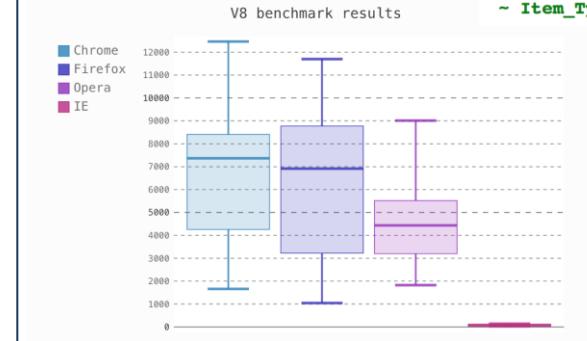


**ggplot**



```
ggplot(train, aes(Item_Visibility, Item_MRP)) +
  geom_point(aes(color = Item_Type)) +
  scale_x_continuous("Item Visibility", breaks =
seq(0,0.35,0.05))+ 
  scale_y_continuous("Item MRP", breaks = seq(0,270,by
= 30))+ 
  theme_bw() + labs(title="Scatterplot") + facet_wrap(~ Item_Type)
```

**pygal**

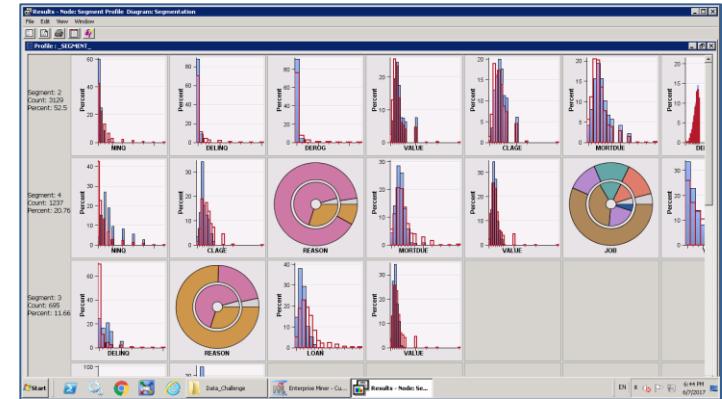
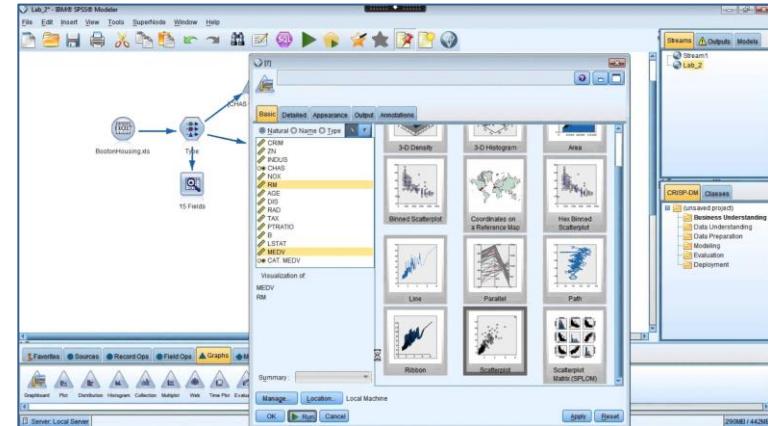
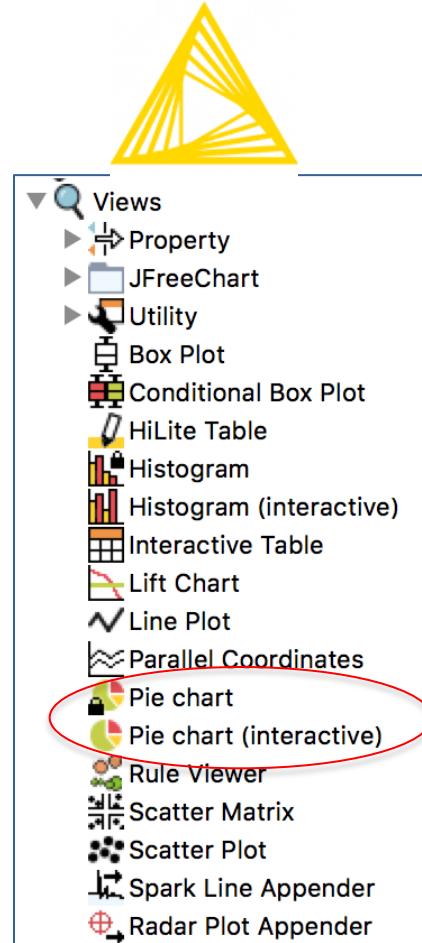


# Or Use These...

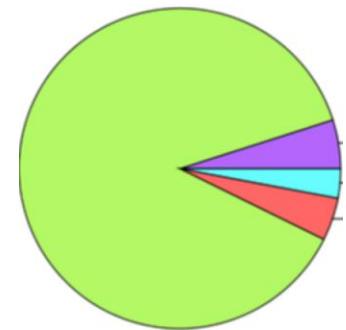


**Visualize**

Tree Viewer	Box Plot	Distribut...	Scatter Plot
Sieve Diagram	Mosaic Display	FreeViz	Linear Projection
Heat Map	Venn Diagram	Silhouette Plot	Pythagor...
Pythagor... Forest	CN2 Rule Viewer	Nomogr...	Geo Map



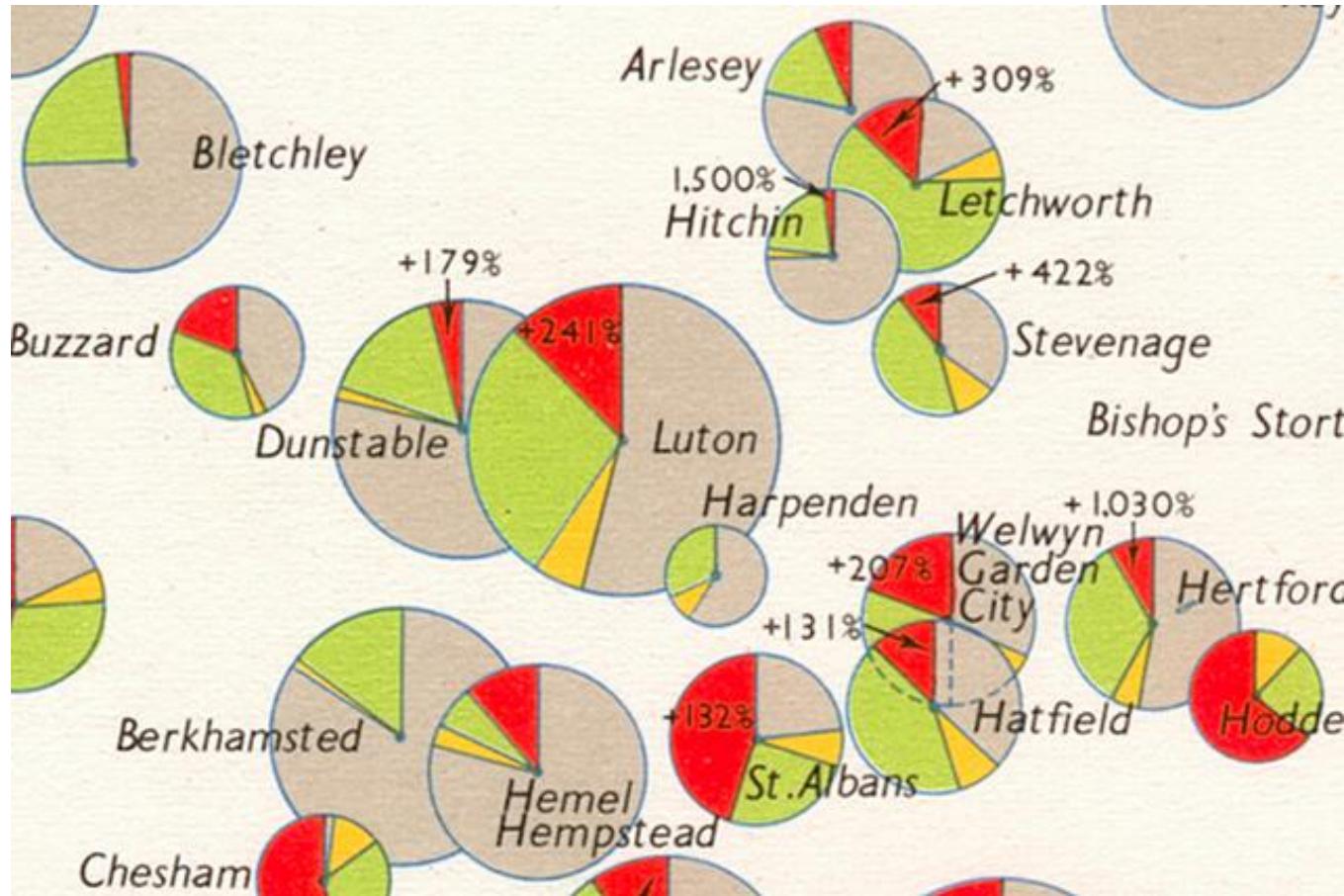
## Pie Charts



- Opinions of experts
  - Edward Tufte: “One of the prevailing orthodoxies of this forum - one to which I whole-heartedly subscribe - is that **pie charts are bad and that the only thing worse than one pie chart is lots of them.**”

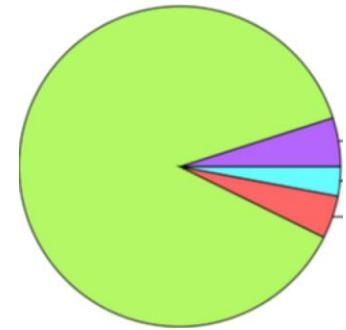
[https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg\\_id=00018S](https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg_id=00018S)

# Tufte Example



[https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg\\_id=00018S](https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg_id=00018S)

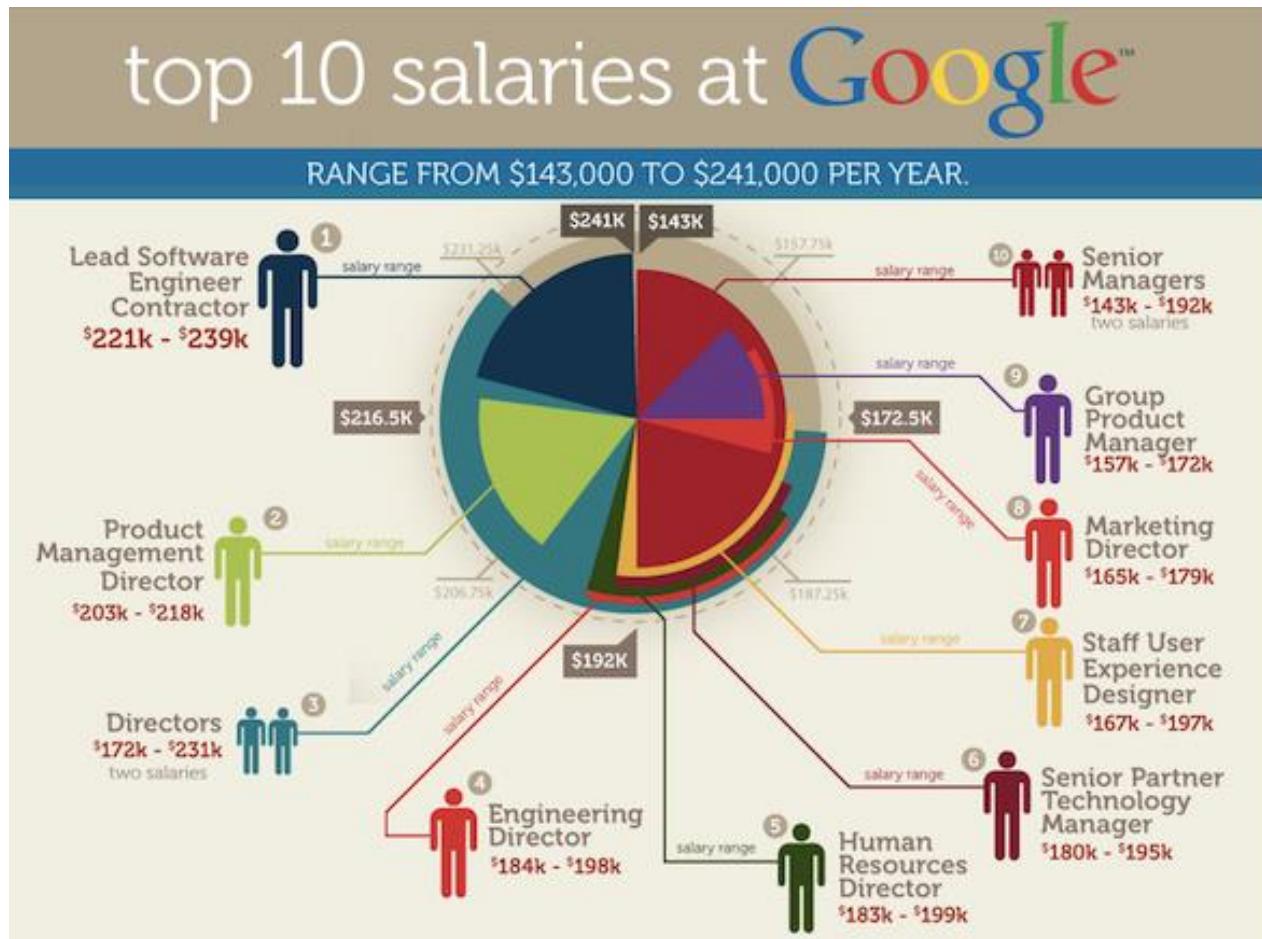
## Pie Charts



- Opinions of experts
  - Kaiser Fung: “The world would be a better place if **pie charts were banned**”

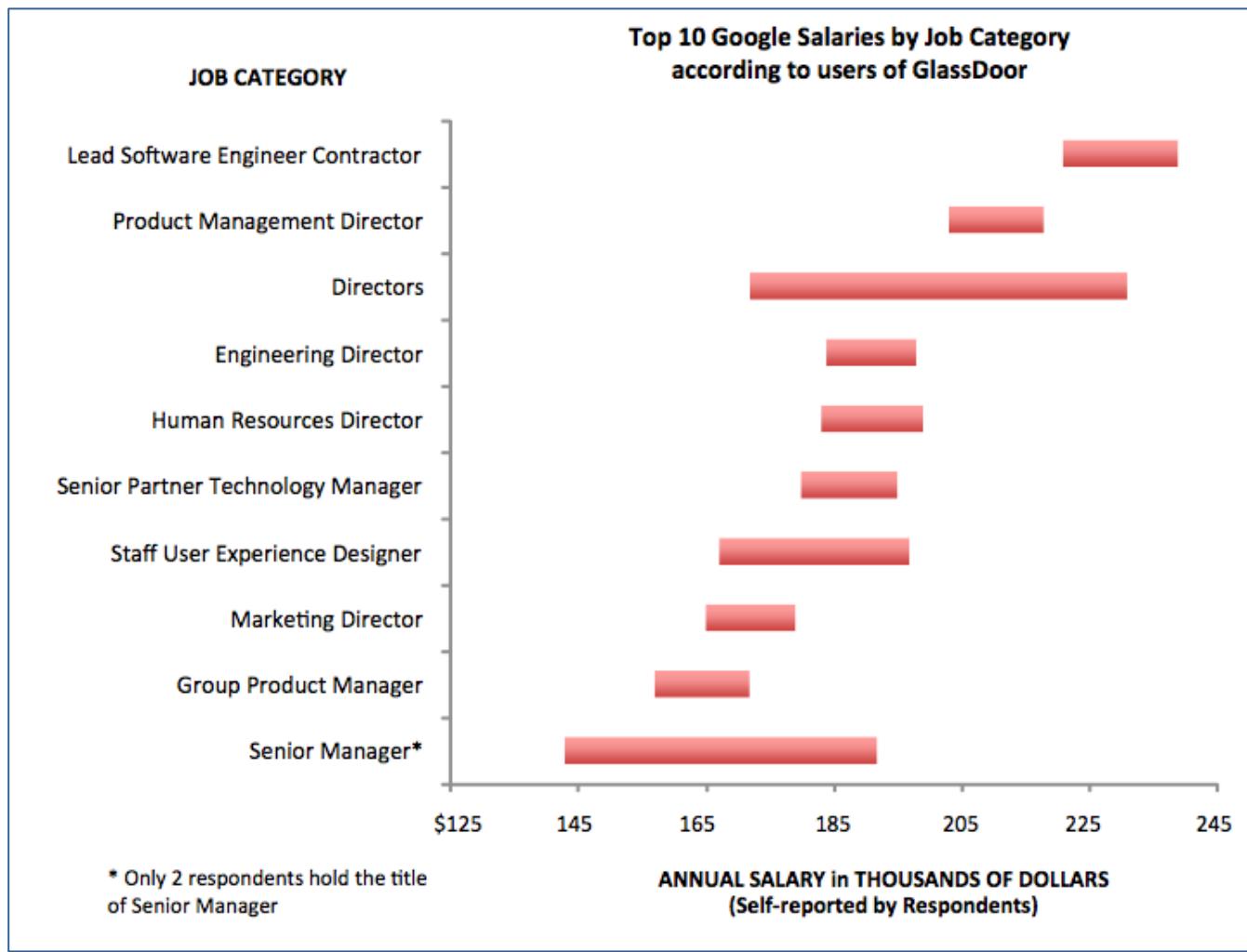
[http://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](http://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)

# Sample Pie chart story from Fung



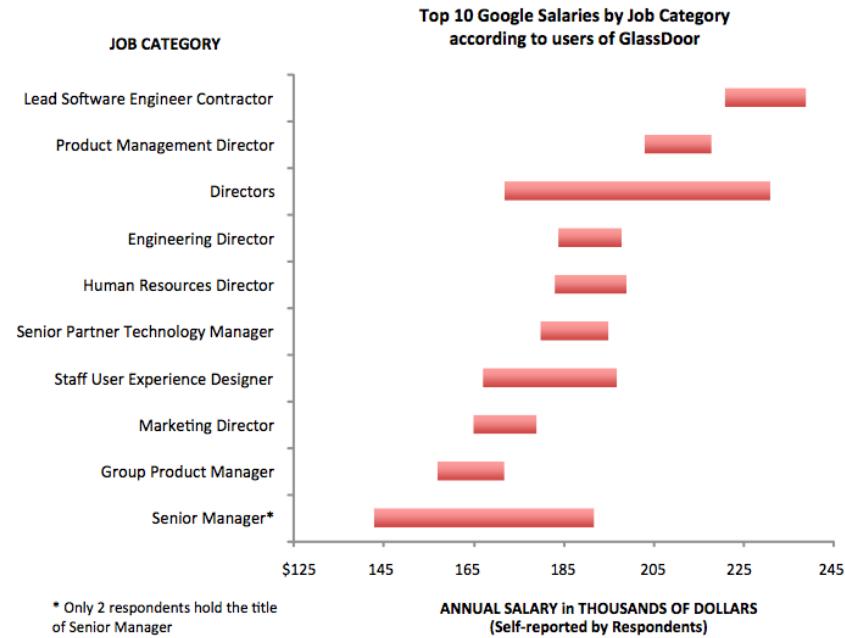
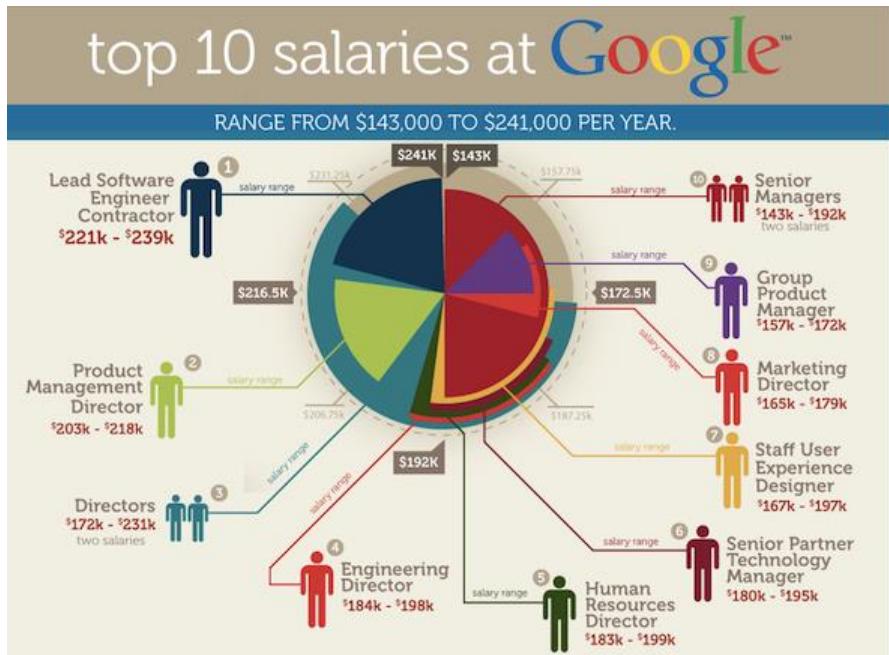
[http://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](http://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)

# Information From Last Pie Chart Clearer as a Bar Chart



[http://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](http://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)

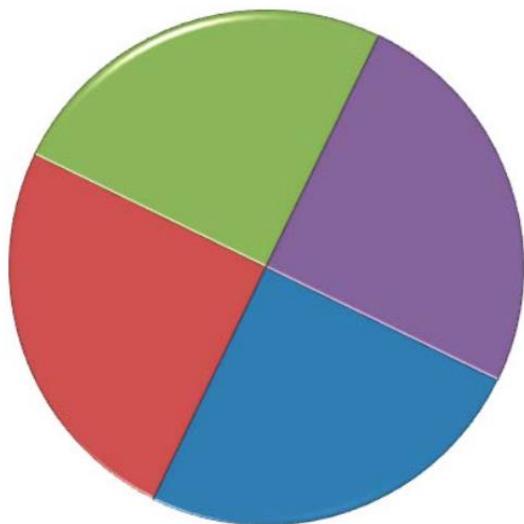
# Bringing A Tufte Concept into this Deck: Comparisons are Best Made on the Same Page



[http://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](http://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)

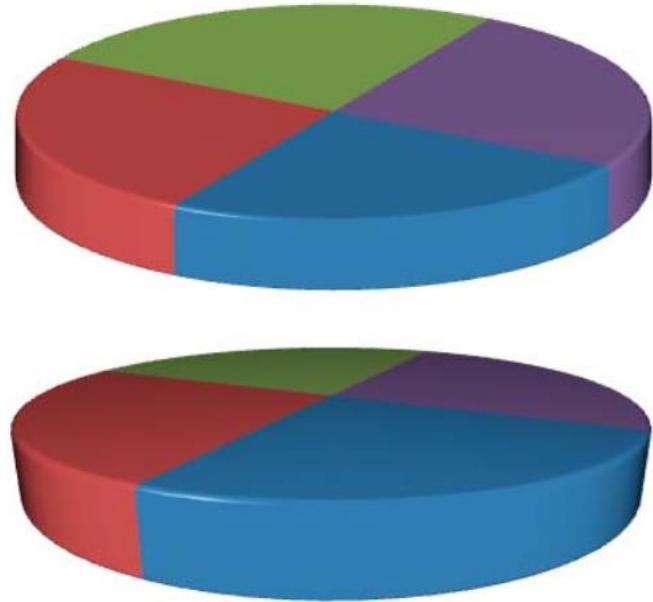
# A Bad Idea: 3-D Pie Charts

**“the only thing worse than a pie chart is a 3-D pie chart” – Dean Abbott**



Color	Size
Green	25%
Blue	25%
Red	25%
Purple	25%

What is the relative size of the green wedge?  
What is the relative size of the blue wedge?



# Pie Charts to Show Relative Proportion (that sum to 1)

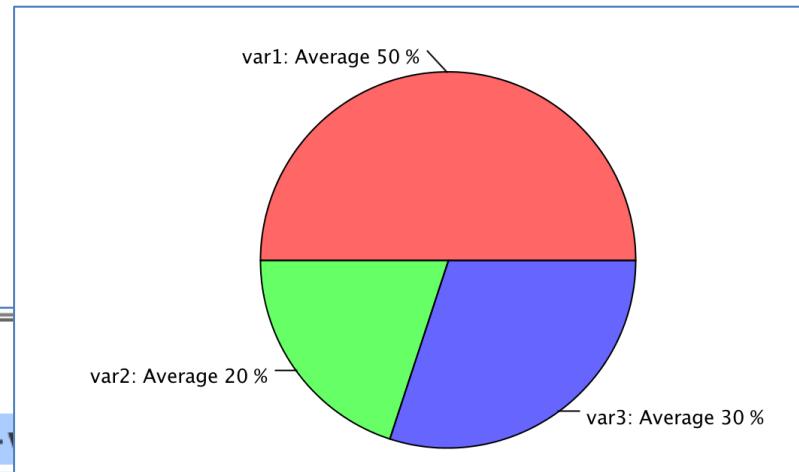
**The message:** var1 has the most influence in the regression model predictions

Statistics on Linear Regression

Variable	Coeff.	Std. Err.	t-Value	p-Value
var1	0.5	1.33E-16	3.75E15	0.0
var2	0.2	1.34E-16	1.49E15	0.0
var3	0.3	1.34E-16	2.24E15	0.0
Intercept	-1.09E-13	4.68E-15	-23.3392	0.0

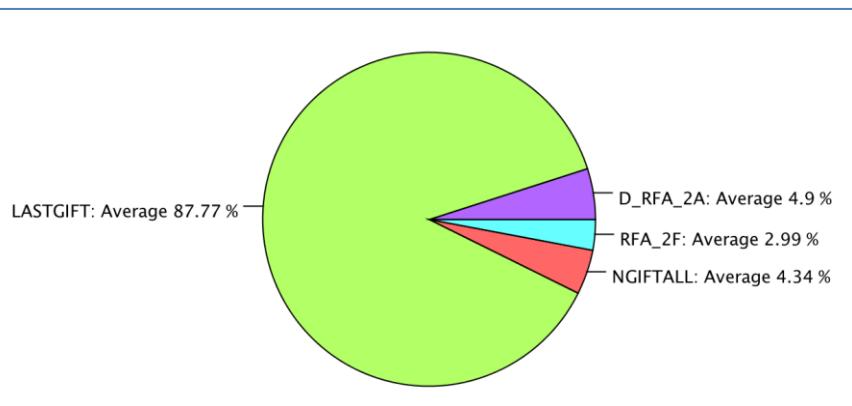
Multiple R-Squared: 1

Adjusted R-Squared: 1



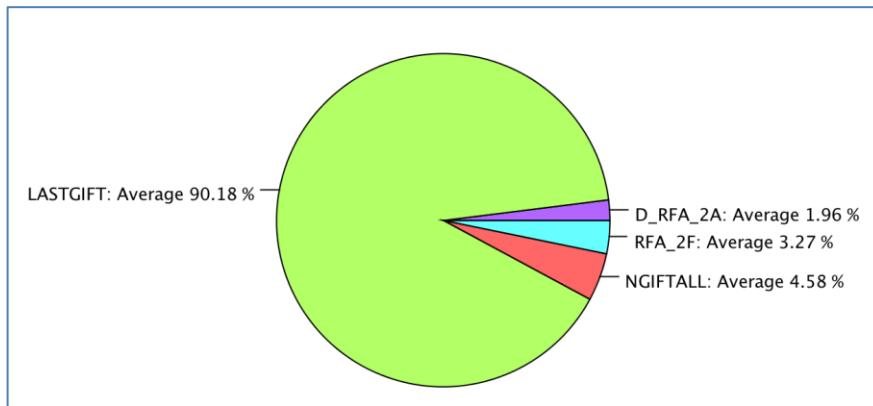
# Pie Chart to Compare Multiple Measures (But One Wedge)

**Coefficient**

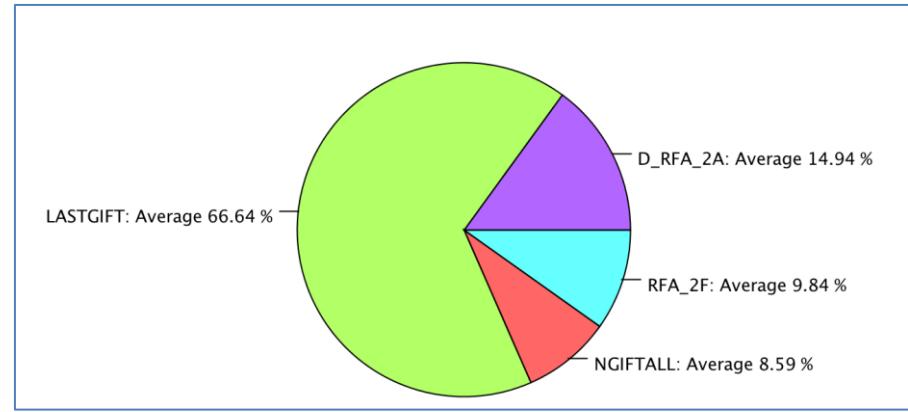


Use `abs()` for calculations

**Direct Proportion**

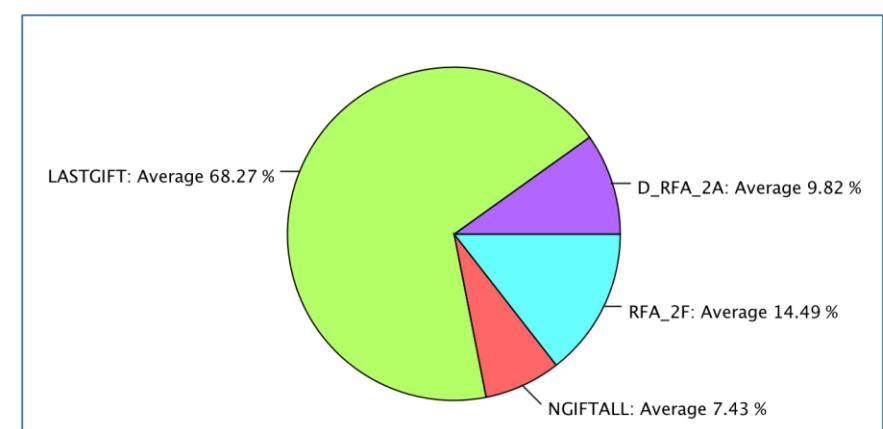


**t-Proportion**

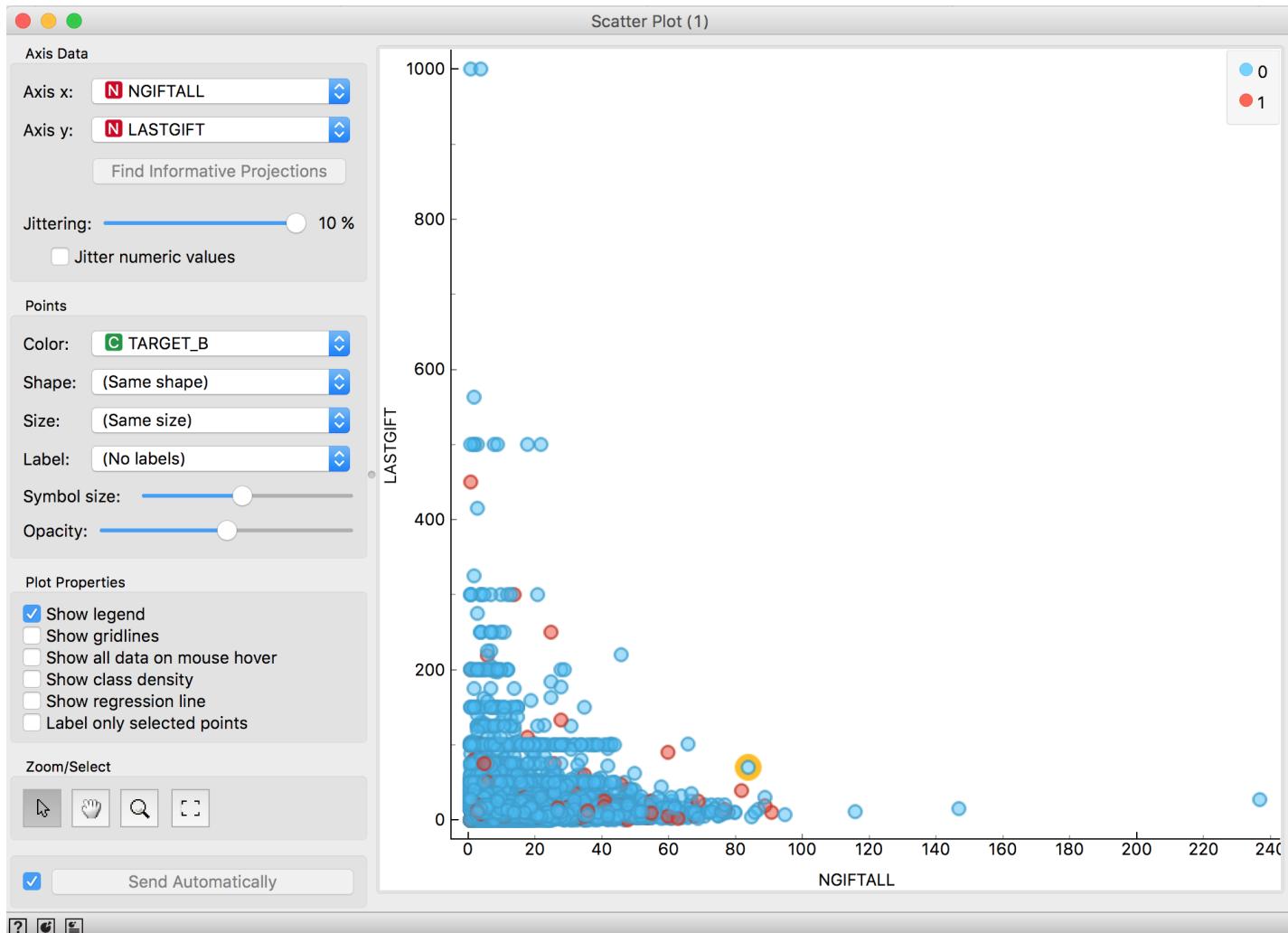


Use `abs()` for t-proportion calculations

**Input Shuffling Proportion**



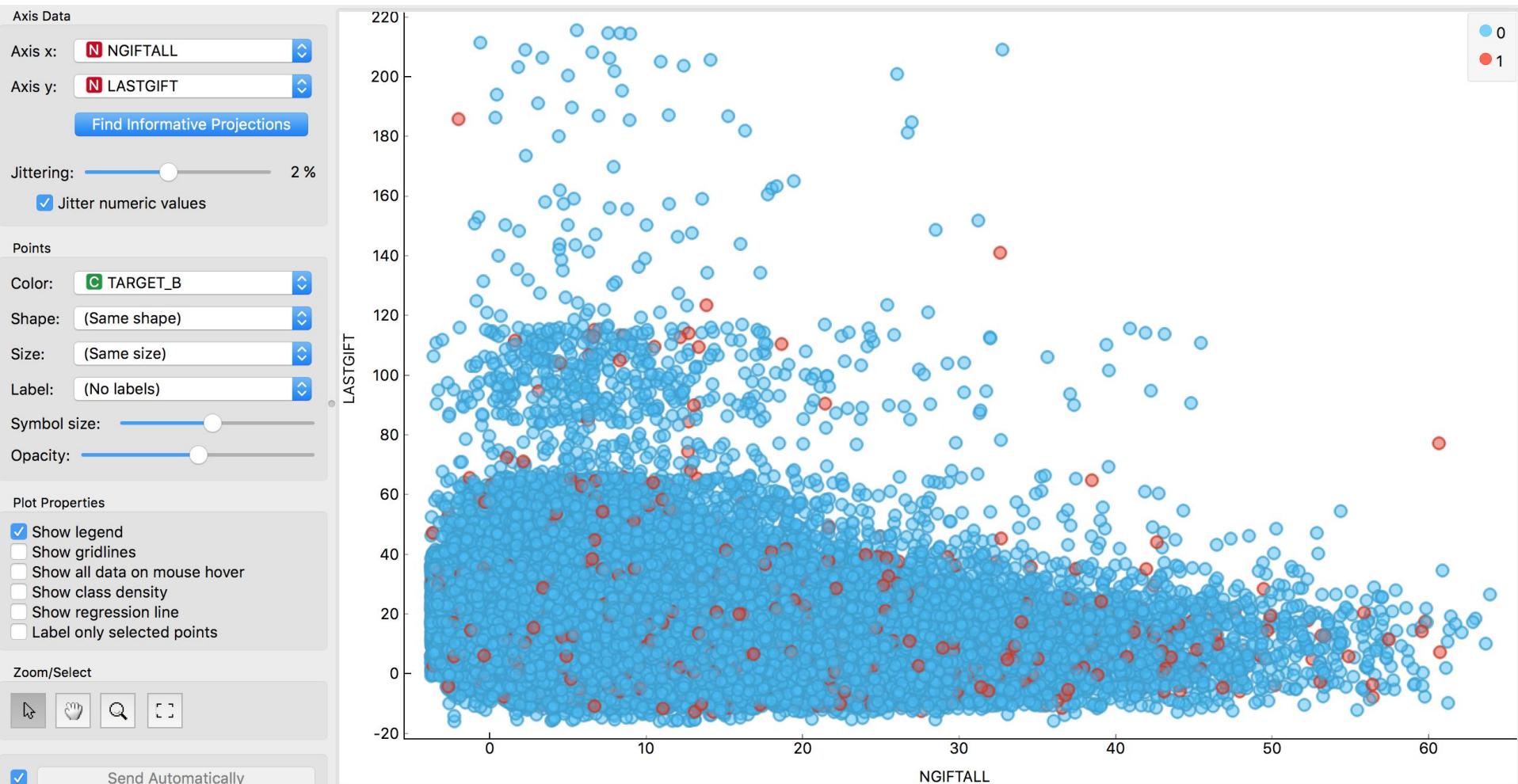
# Many Scatterplots end up looking like this: Skewed Distributions Create Lots of Empty Space



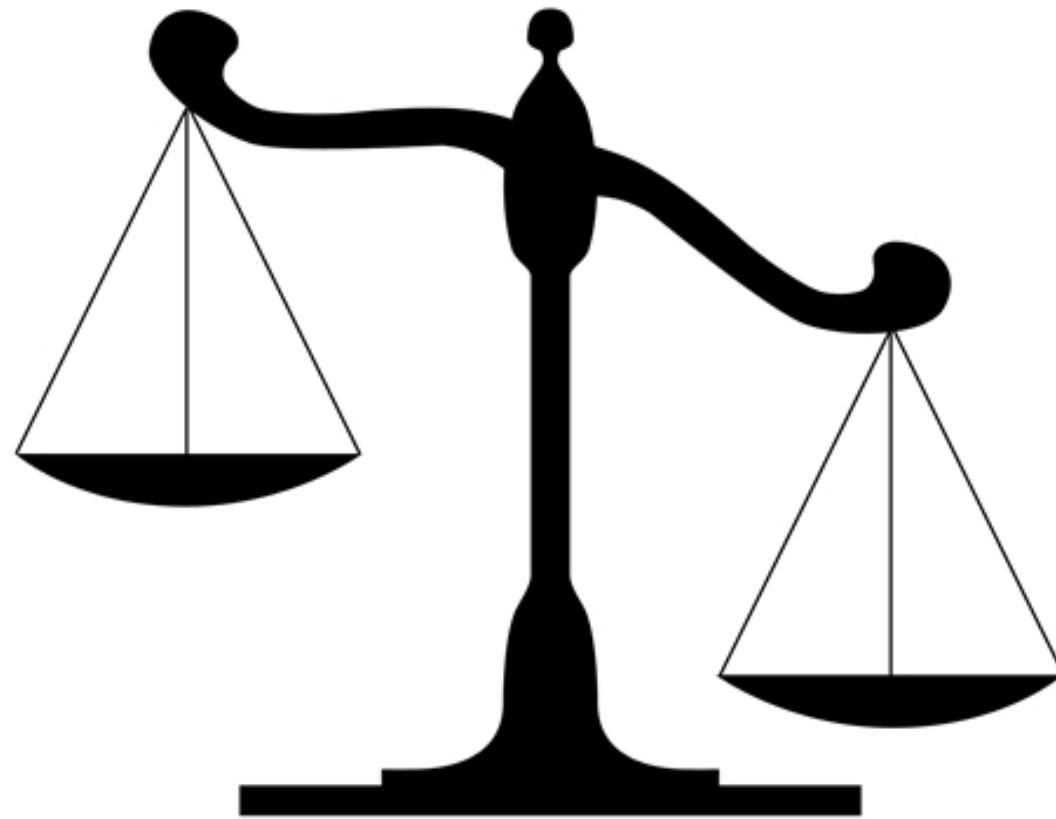
# Normalized Axes Make the Space More Dense, but at the Expense of Interpretability



# Truncated Axes Make the Space Fuller and Intuitive, but Could Miss Interesting Patterns at the Edges

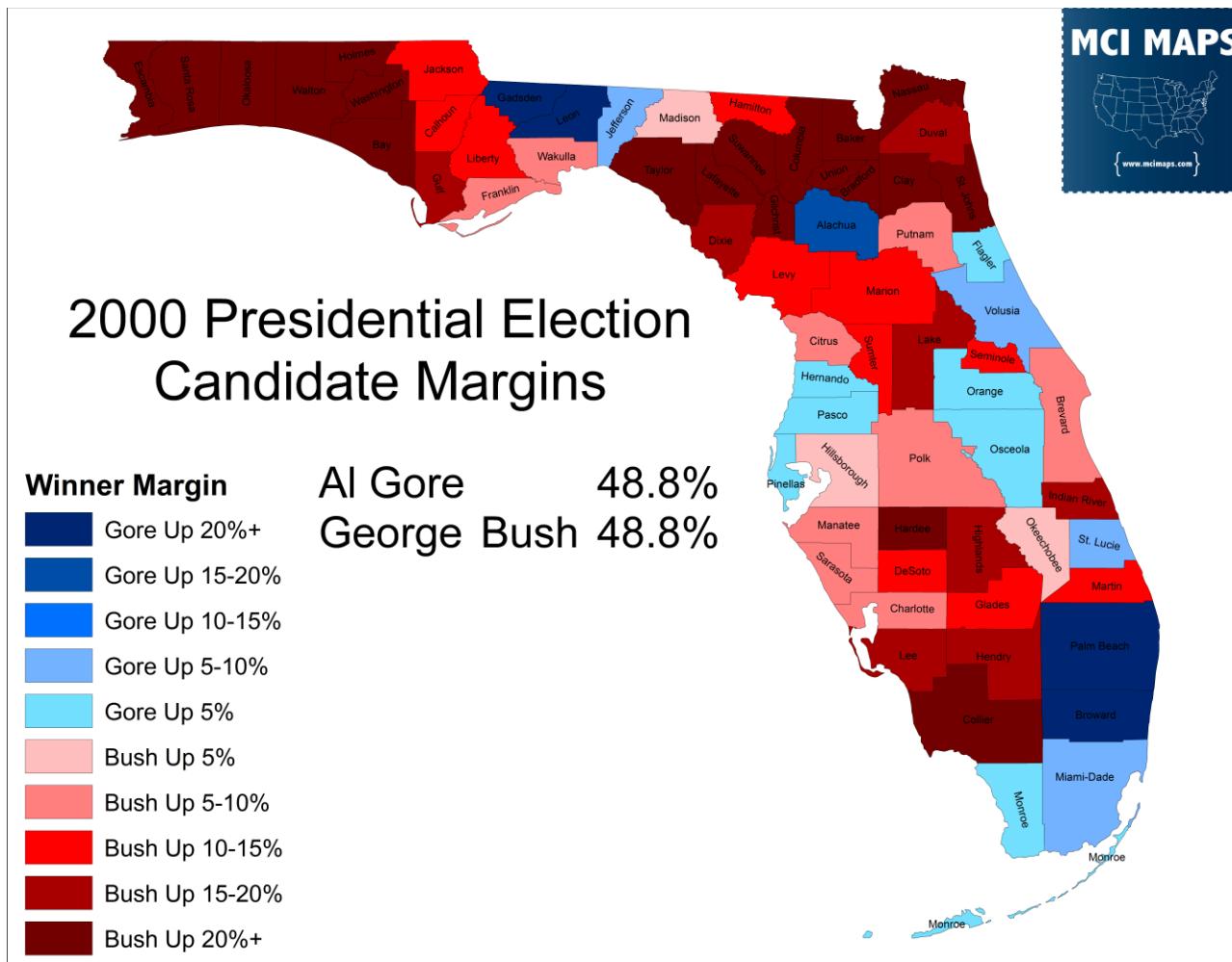


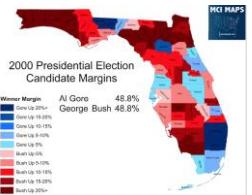
## Principle #2: Tell the Story You Need to Tell: But Be Transparent about the Bias



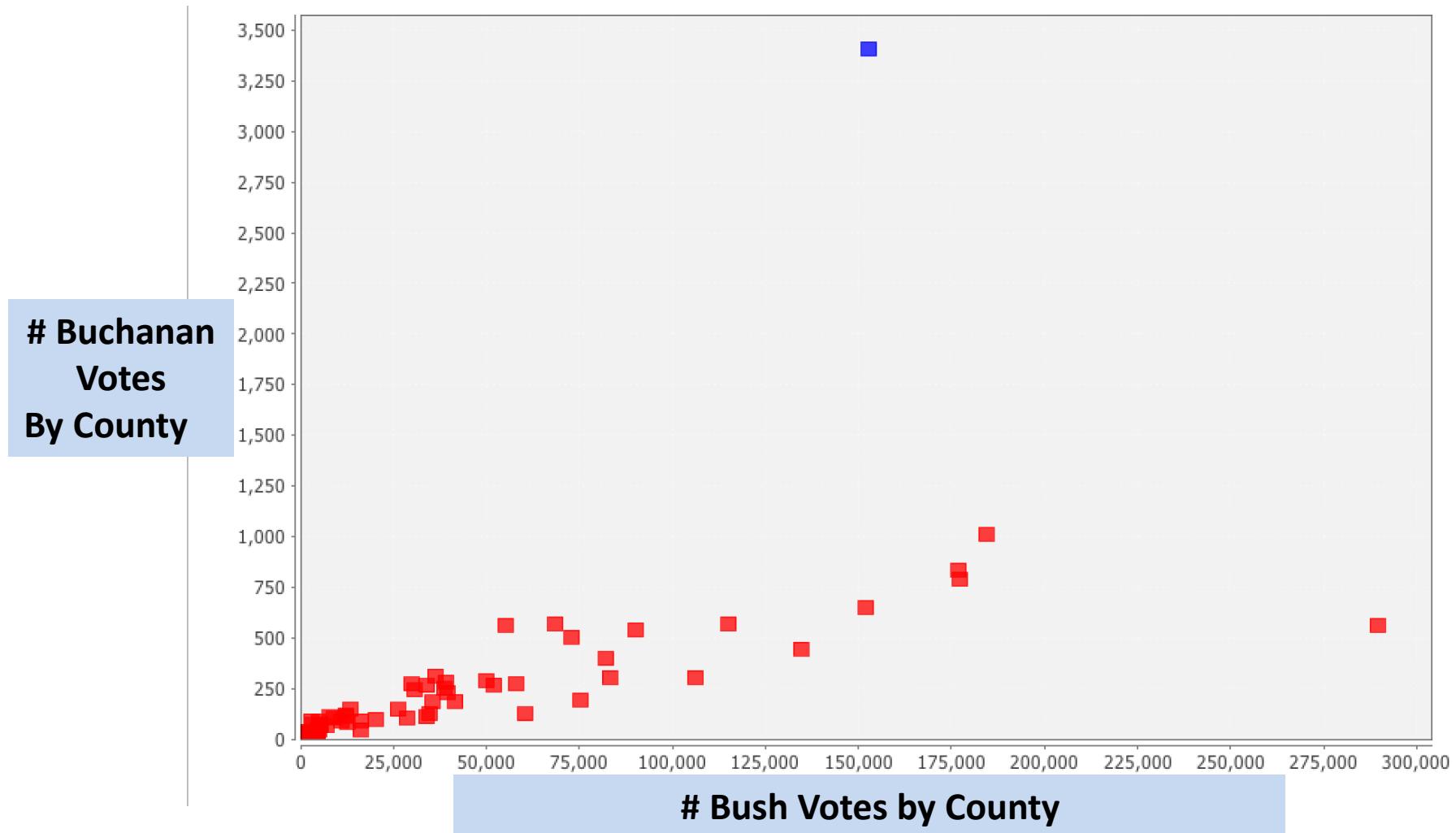
Think about the data, not the politics  
What is the point of the visualization?

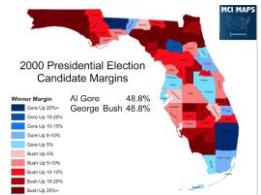
# Do We Really Want to Go Back to the 2000 Election?





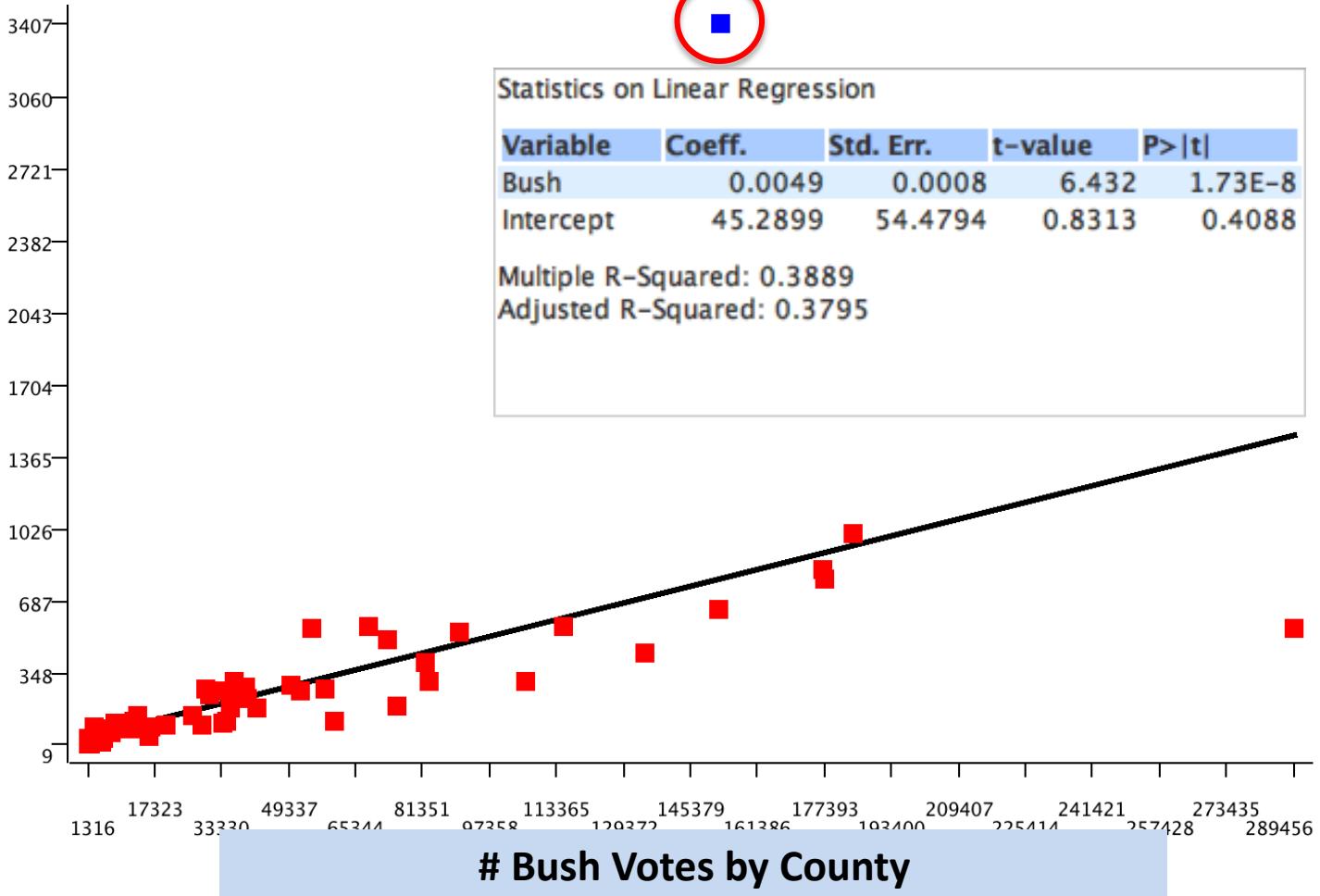
# The Classic Scatterplot of FL Results: Can you Guess Which Point is Palm Beach Cty?

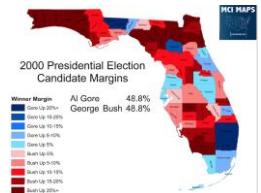




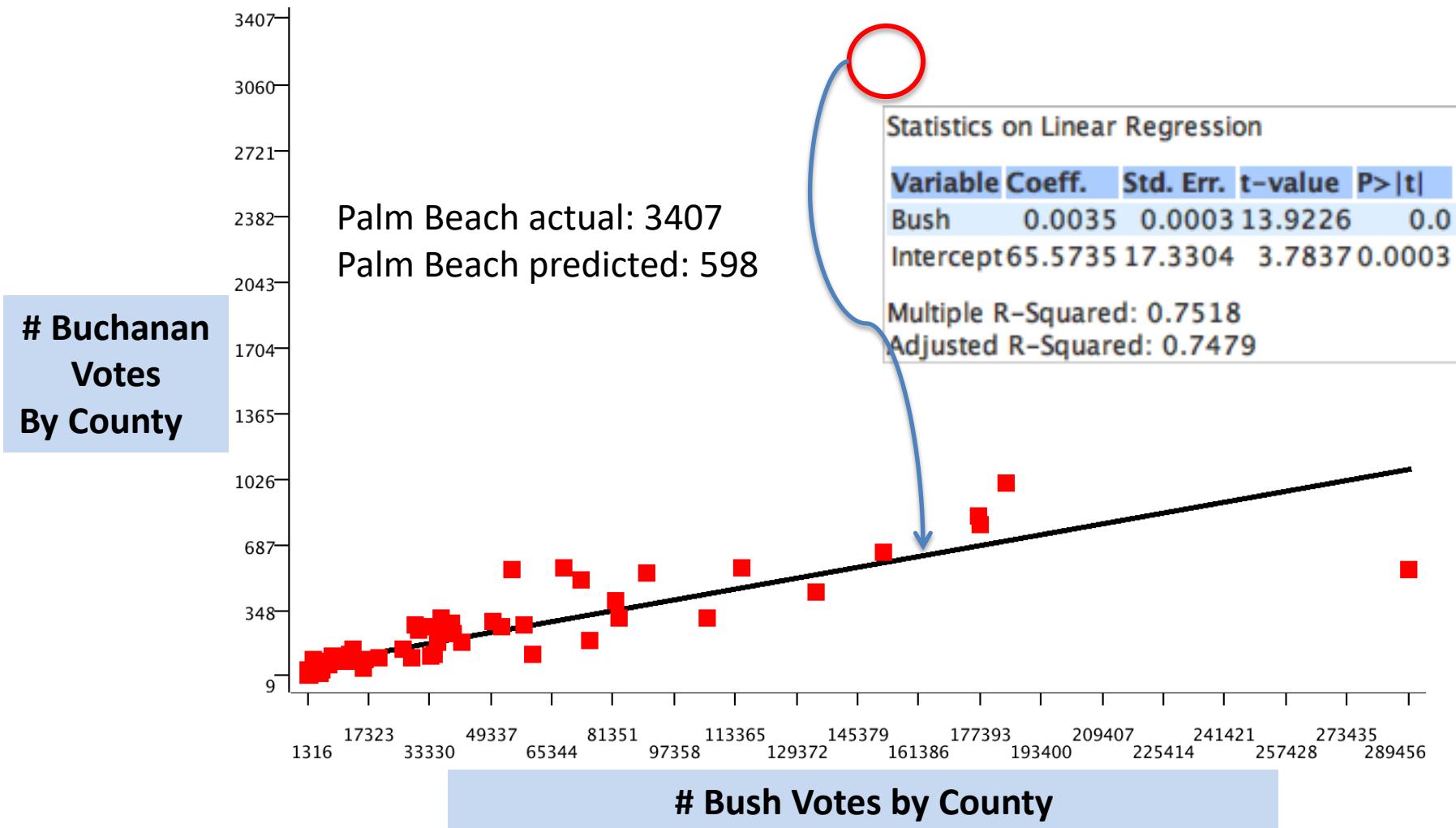
# How Bad is this Outlier?

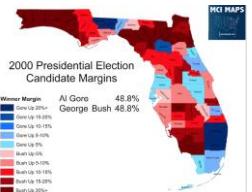
# Buchanan  
Votes  
By County



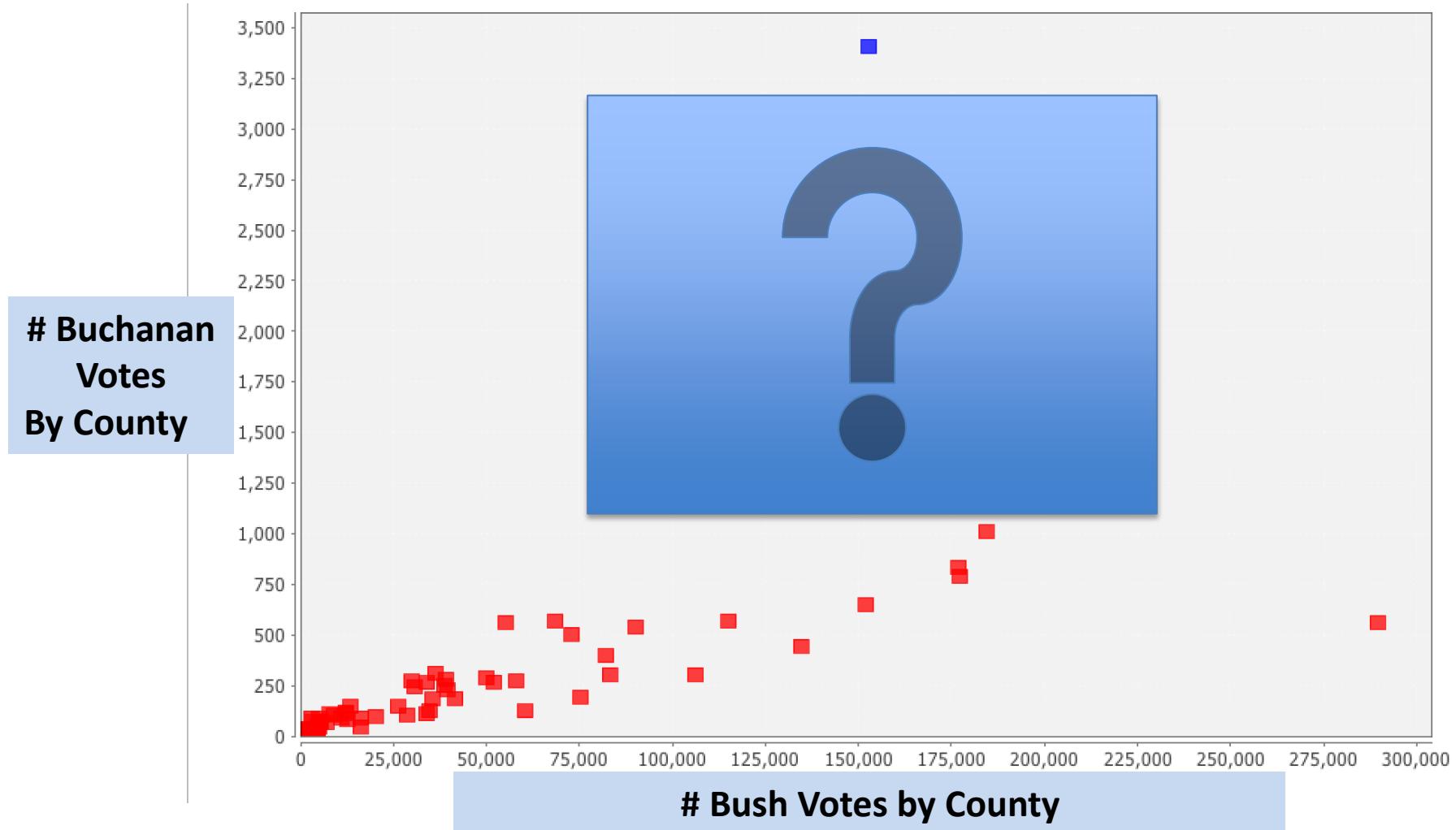


# Remove the Outlier, and Refit



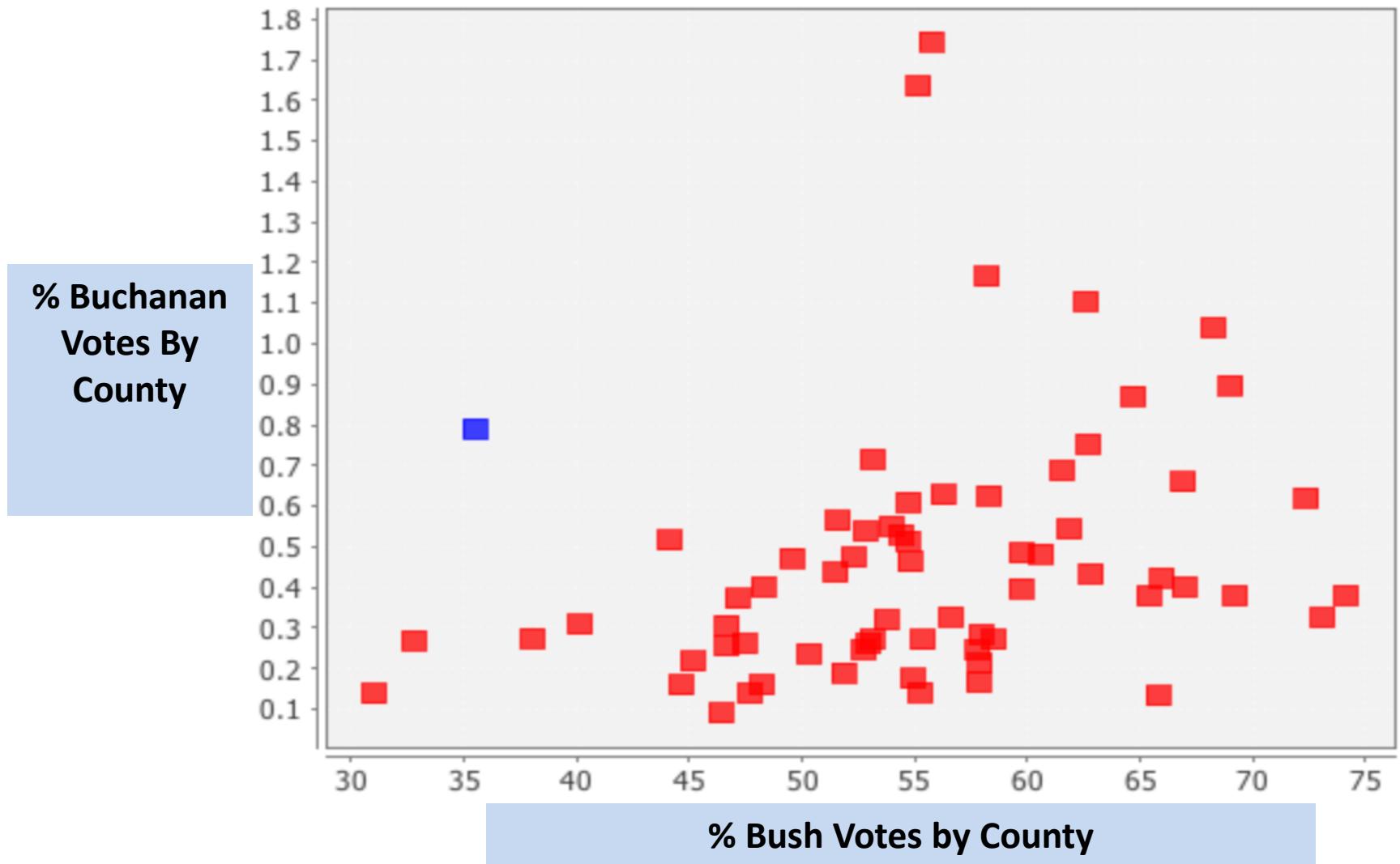


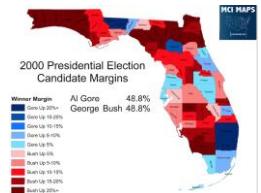
# What Biases Could Exist in This Data? What Assumptions are We Making?





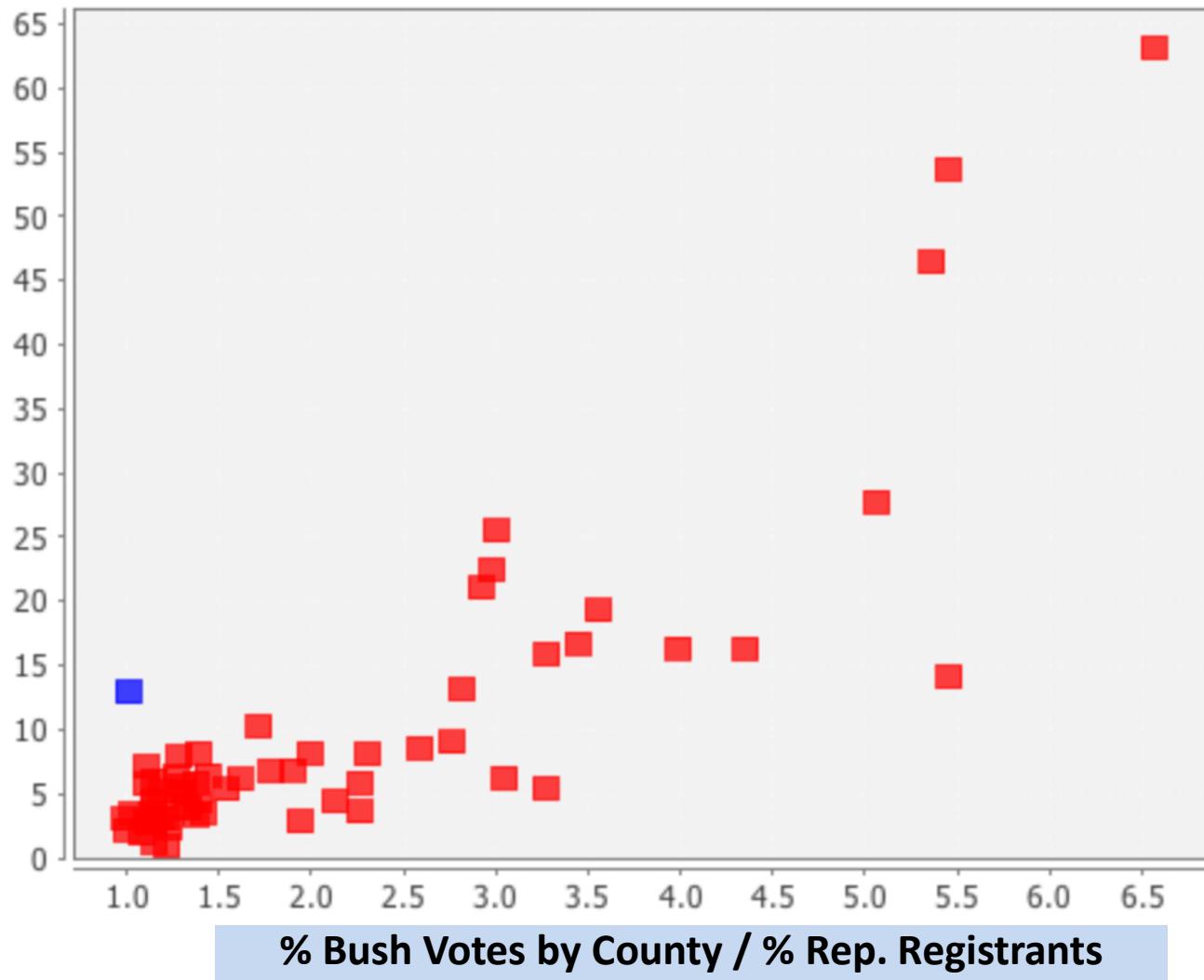
# Normalizing Scales: Remove Magnitude Bias



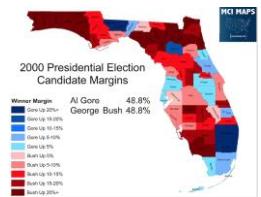


# Adding “Expected” Results: Priors

**% Buchanan  
Votes By  
County  
/**  
**% Reform  
Party  
Registrants**

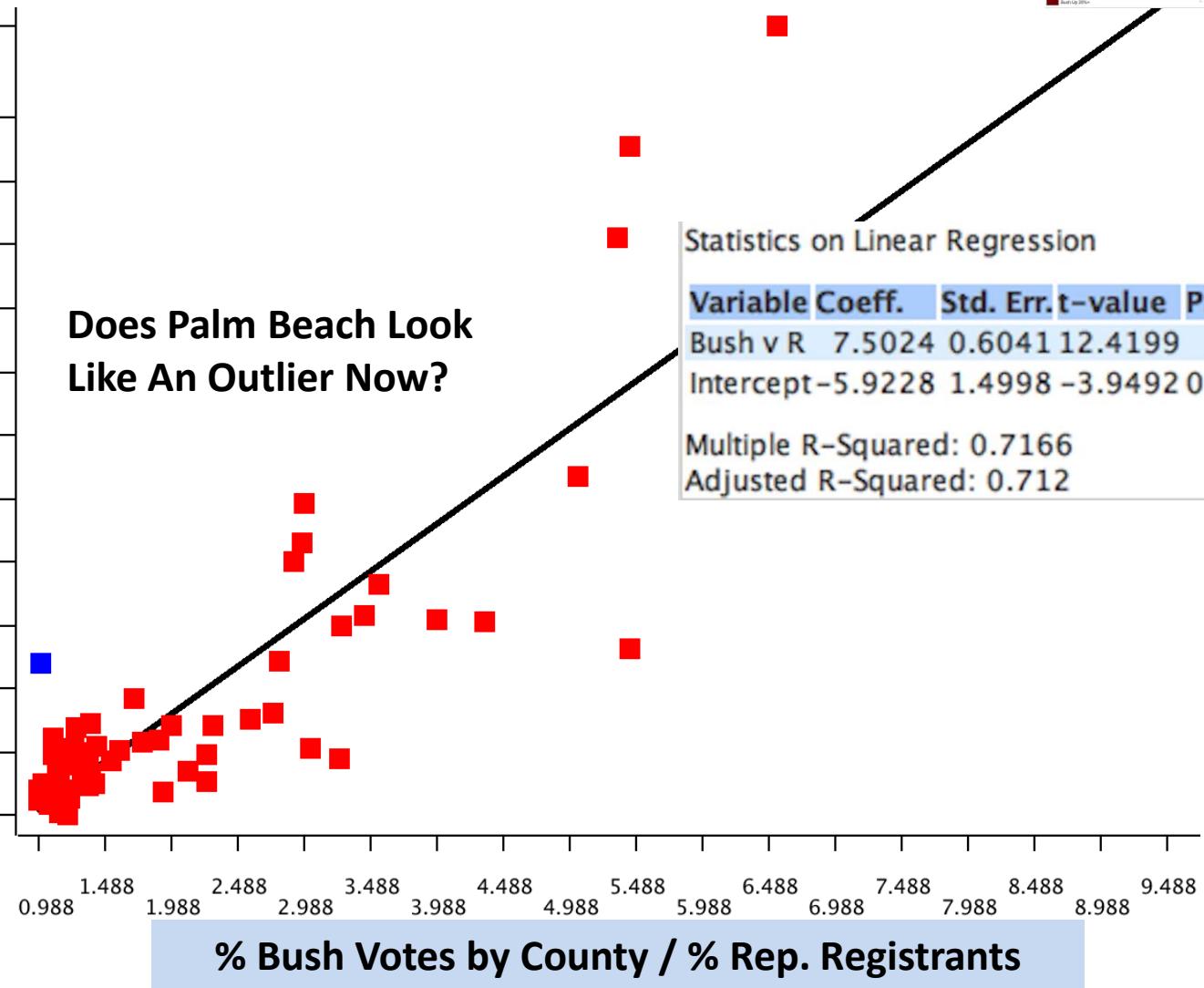


# Refit



% Buchanan  
Votes By  
County  
/  
% Reform  
Party  
Registrants

Does Palm Beach Look  
Like An Outlier Now?



## Principle 3: Tell the Data Story so the Hearer Can Understand



# Tell the Data Story so Your Boss Can Understand

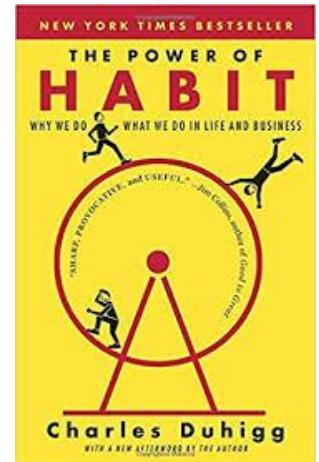


# YMCA Survey Data

- 32,811 surveys, 57 questions
  - All survey questions in analysis had 4, 5, or 6 responses,
  - coded 1 for top/highest answer
- **Q1 - Satisfaction = 1:** 31%
- Q48 - Recommend to Friend = 1: 54%

<http://www.abbottanalytics.com/data-mining-case-study-2.php>

<http://media.salford-systems.com/presentation/A-More-Transparent-Interpretation-of-Health-Club-Surveys-YMCA.pdf>



# Regression Analysis Revealed Seven Question That Were the Most Important

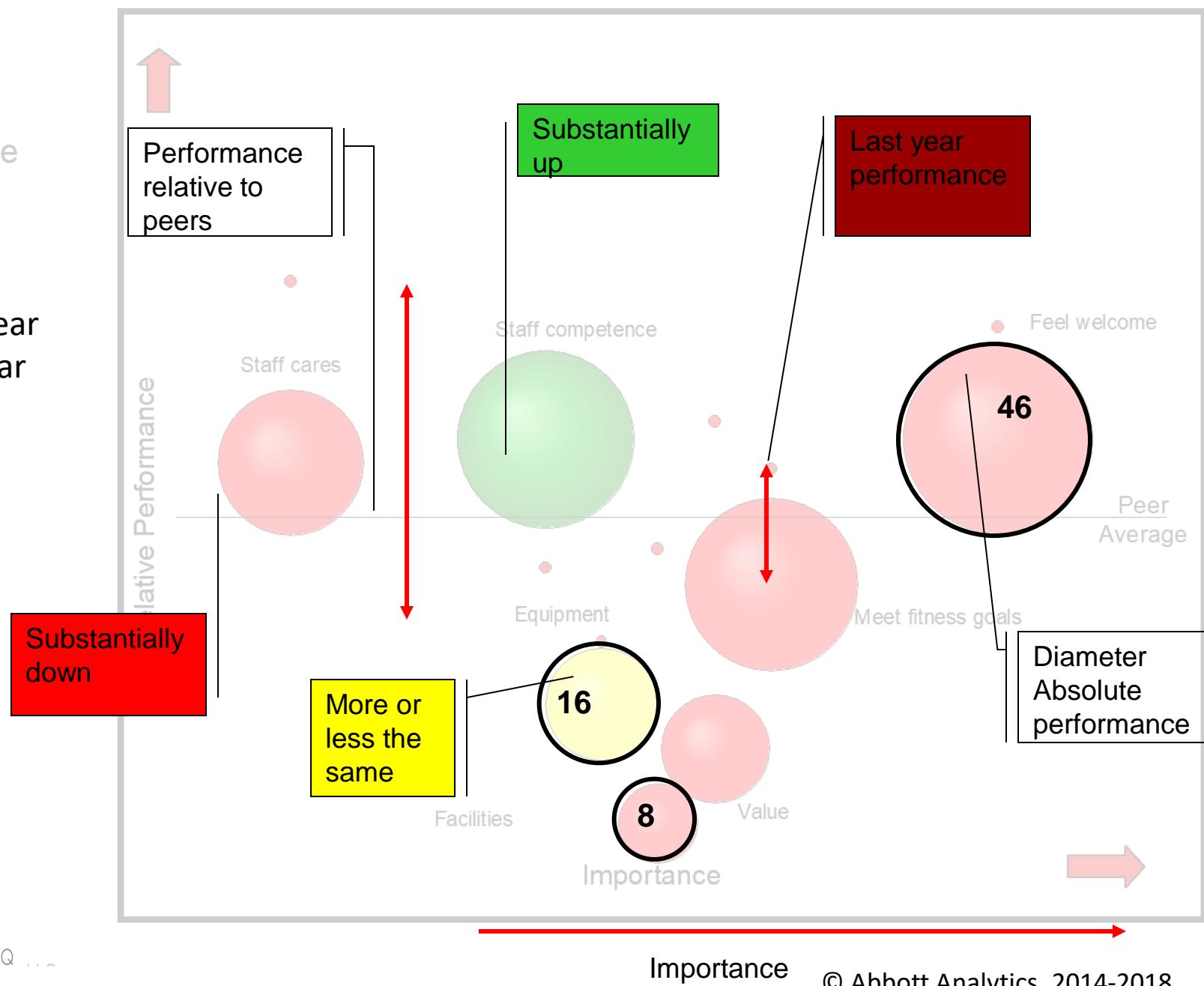
- Empirical comparison: Factors as inputs vs. Top-loading question in factor as input
  - Top-loading or most interesting question on factor as representative of that factor produced slightly better models
  - Use of top-loading question makes final model more easily understood
  - This flies in the face of traditional theory, but worked better operationally
- Final regression model contained these fields:

Database_Question_Number	Factor	Description
Q13	facilities 1	Facilities clean
Q18	equipment 1	Equipment maintained
Q22	value 1	Value for the Money
Q25	relationships 1	Feel Welcome
Q44	goals 1	Y Helps meet fitness goals
Q6	staff 1	Competent Staff
Q9	staff 2	Staff Cares about Well-being

drivers of  
excellence

Current year  
vs. last year

- Better
- Same
- Worse



# I Know....Now

Kaiser Fung: “The **bubble chart is one of the most hopeless data graphics ever invented.**

It is sometimes useful for conceptual charts but trying to express data with it is a lost cause.”

[http://junkcharts.typepad.com/junk\\_charts/2013/03/blowing-the-whistle-at-bubble-charts.html](http://junkcharts.typepad.com/junk_charts/2013/03/blowing-the-whistle-at-bubble-charts.html)

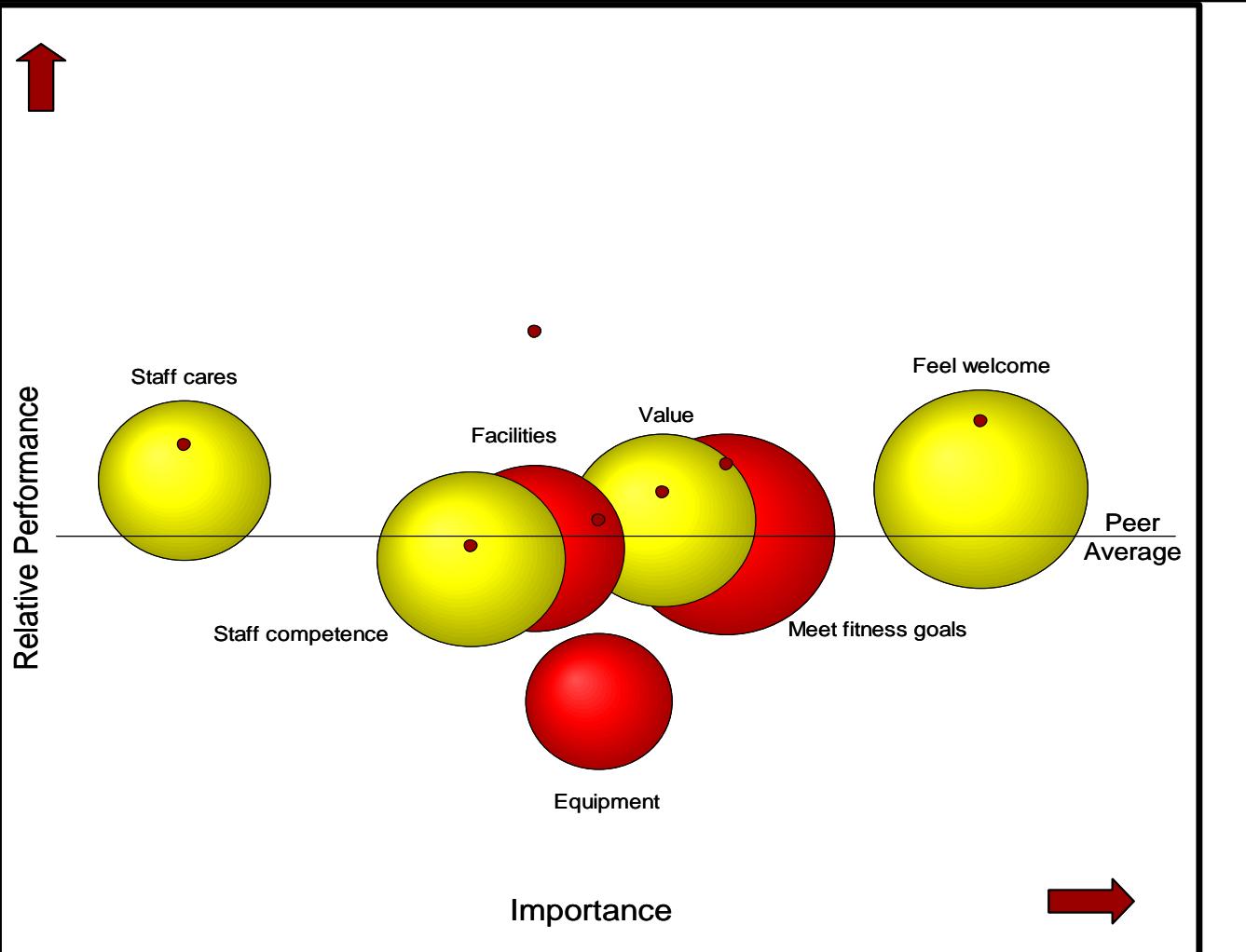


Photo © LoloStock/Shutterstock

drivers of  
excellence

Current year  
vs. last year

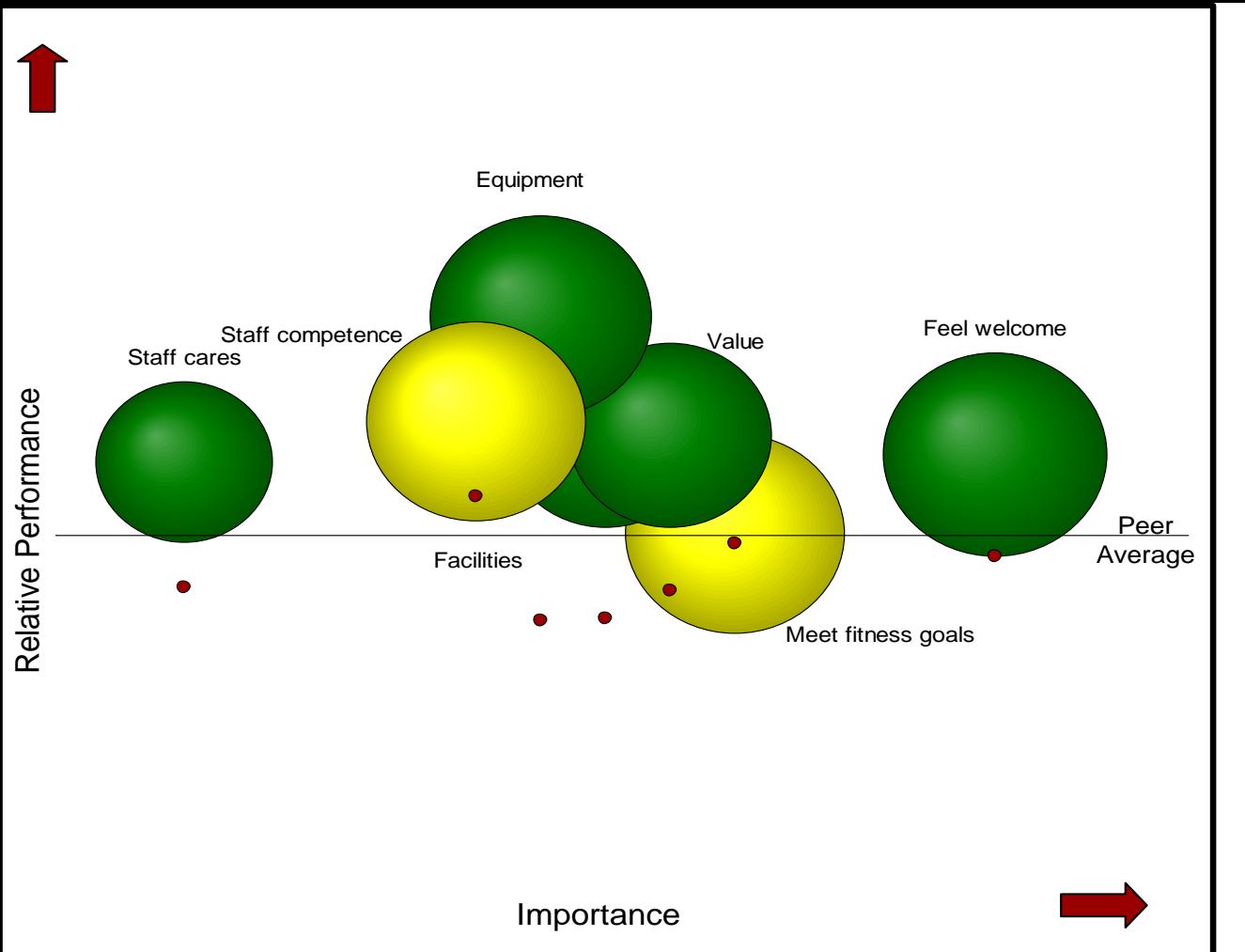
- Better
- Same
- Worse



drivers of  
excellence

Current year  
vs. last year

- Better
- Same
- Worse



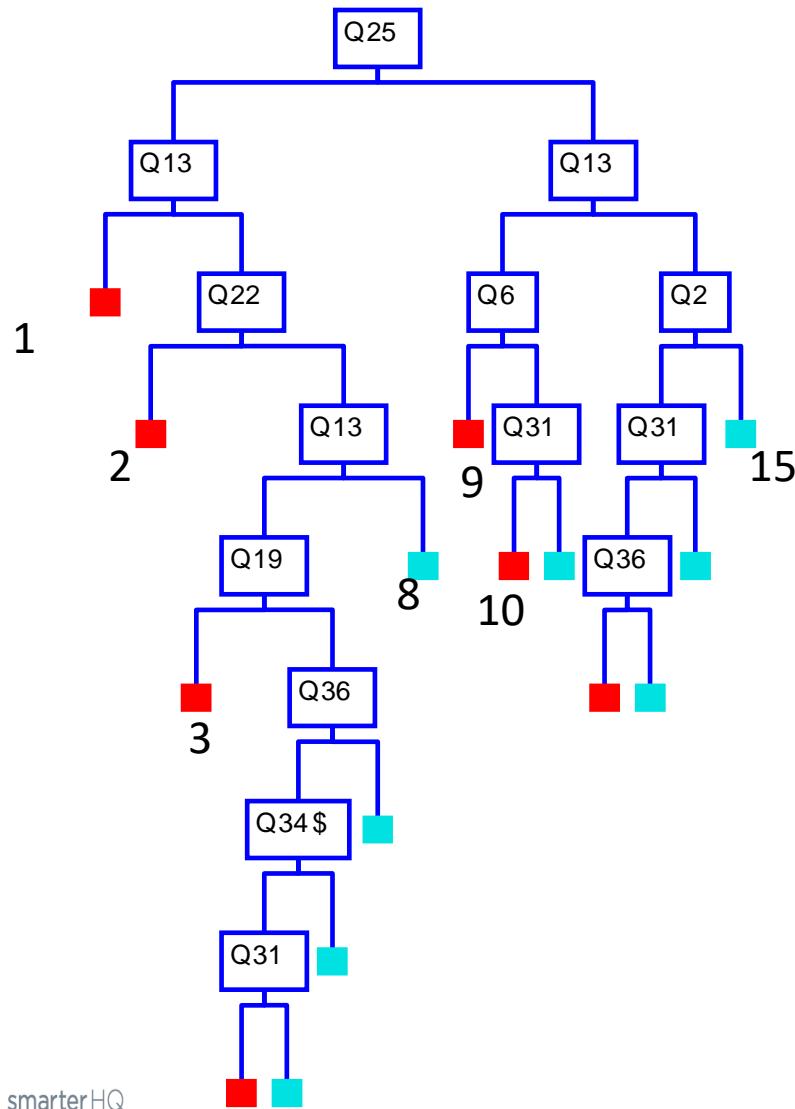
# NO LOVE

(EMINEM)



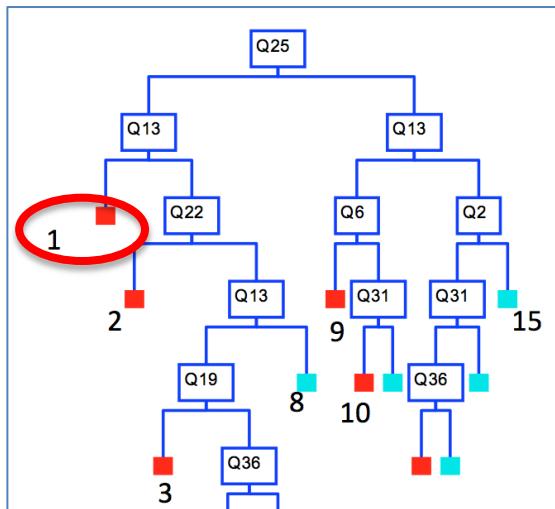
FEAT.  
LIL WAYNE

# Another Approach to Telling the Data Story



- Q25: Feel Welcome
  - Q13: Facilities are clean
- Q22: Value for Money
- Q6: Staff Competent
- Q2: Staff is Efficient
- Q31: You are Loyal to the YMCA
- Q36: Available Cardio Equipment

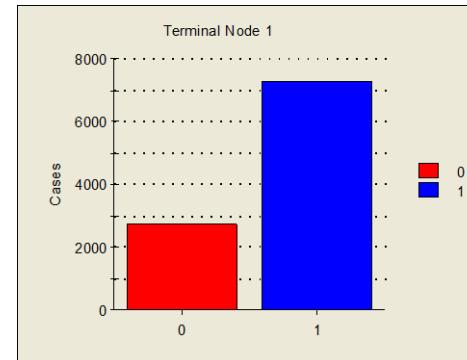
# Member Satisfaction Model: Key Rules



If strongly agree that **facilities are clean** and strongly agree that member **feels welcome**, then highly satisfied

**~3/4 of members who match this profile are Highly Satisfied**

**~1/2 of ALL Highly Satisfied Members Fit this Profile**



/\*Rules for terminal node 1\*/

Matches

- 10,014 surveys (20.8%),
- 7,289 highly satisfied (72.8%),
- 49% of all highly satisfied

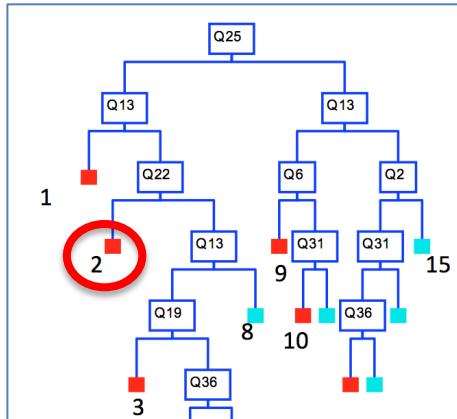
RULE:

If ( Q13 = 1 and Q25 = 1)

Then Satisfaction = 1

P(1) = 0.727881; Lift = 2.4

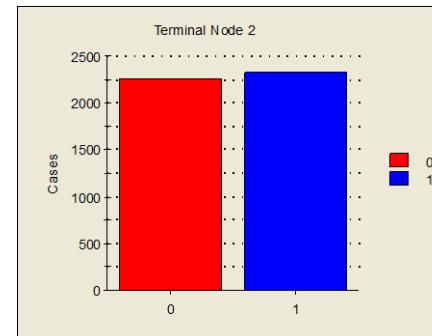
# Member Satisfaction Model: Key Rules



Even if don't strongly agree **facilities are clean**,  
If strongly agree that **feel welcome** and  
strongly agree Y is **value for money**,  
then  
highly satisfied

**~1/2 of members who match this profile are Highly Satisfied**

**~1/6 of ALL Highly Satisfied Members Fit this Profile**



/\*Rules for terminal node 2\*/

Matches

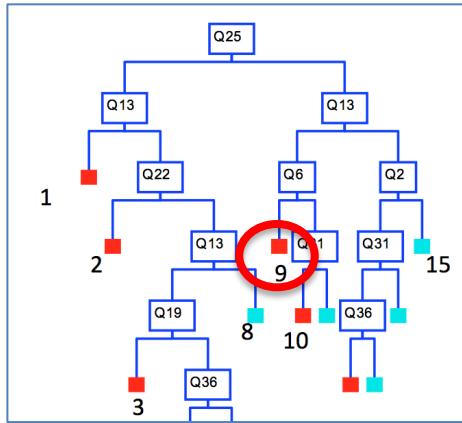
- 4,578 surveys (9.5%),
- 2,317 highly satisfied (50.6%),
- 15.6% of all highly satisfied

RULE:

If ( Q13 <> 1 and Q22 = 1 and Q25 = 1  
Then Satisfaction = 1

$$P(1) = 0.506116; \text{ Lift} = 1.6$$

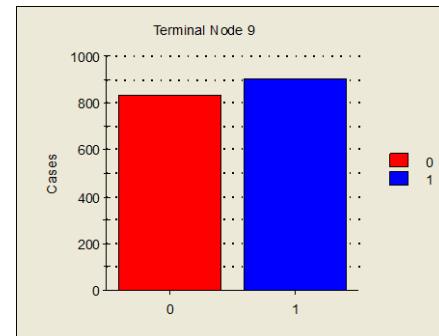
# Member Satisfaction Model: Key Rules



Even if don't strongly agree **feel welcome**,  
If strongly agree that **facilities are clean**, and  
strongly agree that **staff is competent**,  
then  
highly satisfied

**~1/2 of members who match this profile are Highly Satisfied**

**~1 in 16 of ALL Highly Satisfied Members Fit this Profile**



/\*Rules for terminal node 9\*/

Matches

- 1,739 surveys (3.6%),
- 904 highly satisfied (52.0%),
- 6.1% of all highly satisfied

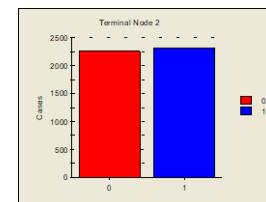
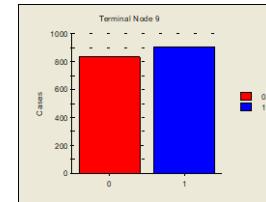
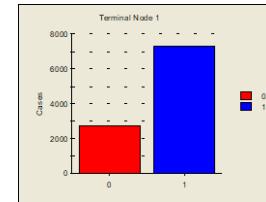
RULE

IF ( Q6 = 1 and Q13 = 1 and Q25 <> 1

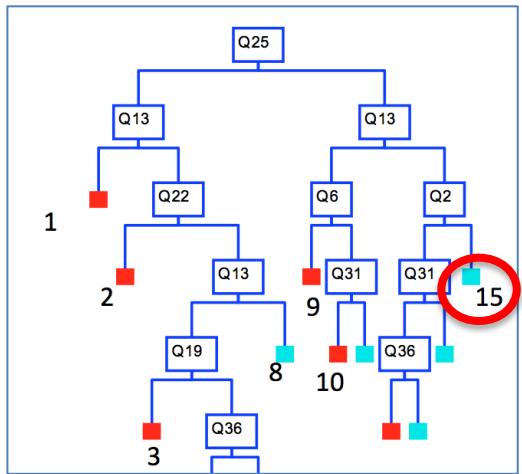
Then Satisfaction = 1  
 $P(1) = 0.52$ ; Lift = 1.7

# Member Satisfaction Model: Key Rules

- facilities are clean and member feels welcome
  - + even if don't feel welcome, if facilities are clean and staff is competent
    - + even don't agree facilities are clean, if feel welcome and Y is value for money
- Top 3 rules comprise ~ 3/4 of all highly satisfied members



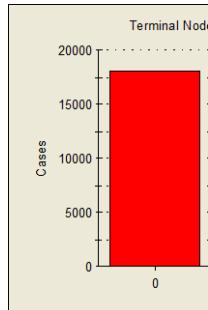
# Member Satisfaction Model: Key Negative Rules



If don't strongly agree that **staff is efficient** and don't strongly agree that **feel welcome**, and don't strongly agree that the **facilities are clean**, then member isn't highly satisfied

**~11 of 12 of members who match this profile are NOT Highly Satisfied**

**~6 in 10 of ALL NOT Highly Satisfied Members Fit this Profile**



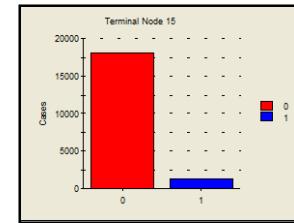
- ```
/*Rules for  
Matches
```
- 19,323 su
  - 1,231 hig
  - 8.3% of 1
  - 58.2% of 1

RULE:  
If ( Q2 <> 1  
Q25 <> 1  
Then Not  
P(1) = 0.00

# Member Satisfaction Model: Not Highly Satisfied

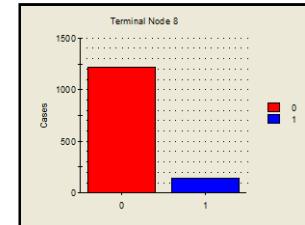
- If staff not efficient and don't feel welcome and facilities aren't clean

+



- Even if feel welcome, if facilities are not clean and Y is not a good value for the money

+



- Top 2 rules comprise ~ 2/3 of all NOT highly satisfied members

# Conclusions

- Know your software visualization building blocks
  - You will have to build most visualizations quickly, so become efficient with what you have

# Conclusions

- Know your software visualization building blocks
  - You will have to build most visualizations quickly, so become efficient with what you have
- Every visualization is biased
  - Bias isn't bad; be transparent about the bias
  - Different biases reveal complementary aspects of the data story

# Conclusions

- Know your software visualization building blocks
  - You will have to build most visualizations quickly, so become efficient with what you have
- Every visualization is biased
  - Bias isn't bad; be transparent about the bias
  - Different biases reveal complementary aspects of the data story
- Know your Audience
  - Tell a data story that the hearers understand