

# Designing Exploratory Visualization Tools

---

# Basic charts for EDA

---

Table

Bar graph, dot plot

Histogram, frequency polygon, strip plot

Line graph

Scatter plot, bubble (size for third variable)

Box plot

Heatmap matrix

Bar/line or bar/dot combination

What makes for a good visual EDA tool?

# Seamless data interaction

**Dimensions**

- Customer
  - Customer Name
  - Segment
- Order
  - Order Date
  - Order ID
  - Ship Date
  - Ship Mode
- Location
  - Country
  - State
  - City
  - Postal Code
- Product
  - Category
  - Sub-Category
  - Manufacturer
  - Product Name

**Measures**


- Discount
- Profit
- Profit Ratio
- Quantity
- Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records

**Sets**

- Top Customers by Profit

**Parameters**

- Profit Bin Size
- Top Customers

 Columns	$\text{SUM}([\text{Profit}]) / \text{SUM}([\text{Sales}])$
 Rows	Profit (bin)

Calculated Field

[Ship Date]

The calculation is valid.

Apply

OK

	user_id	text	diag_date	from_diag	created_date	target
158437	2906418756	RT @karimaro11: 先月から家に居るコマちゃん。30日間を30秒間にまとめてみま...	2014-06-11	582.0	2016-01-14	1
24971	86740068	Home made fried oreos. Yum! <a href="http://yfrog.com/k...">http://yfrog.com/k...</a>	None	NaN	2011-08-07	0
122398	1389531	Marti Belle and Tina San Antonio back when the...	2015-11-10	60.0	2016-01-09	1
84403	248388481	HOLY FUCKING SHIT!!!!!!!!!!!!!! #HannibalFinale @...	None	NaN	2014-05-24	0
47950	45013513	MobaXTerm is pretty much the best thing ever.	None	NaN	2015-04-26	0

Data transformation operations (e.g.,  
`df.groupbyKey( )`).

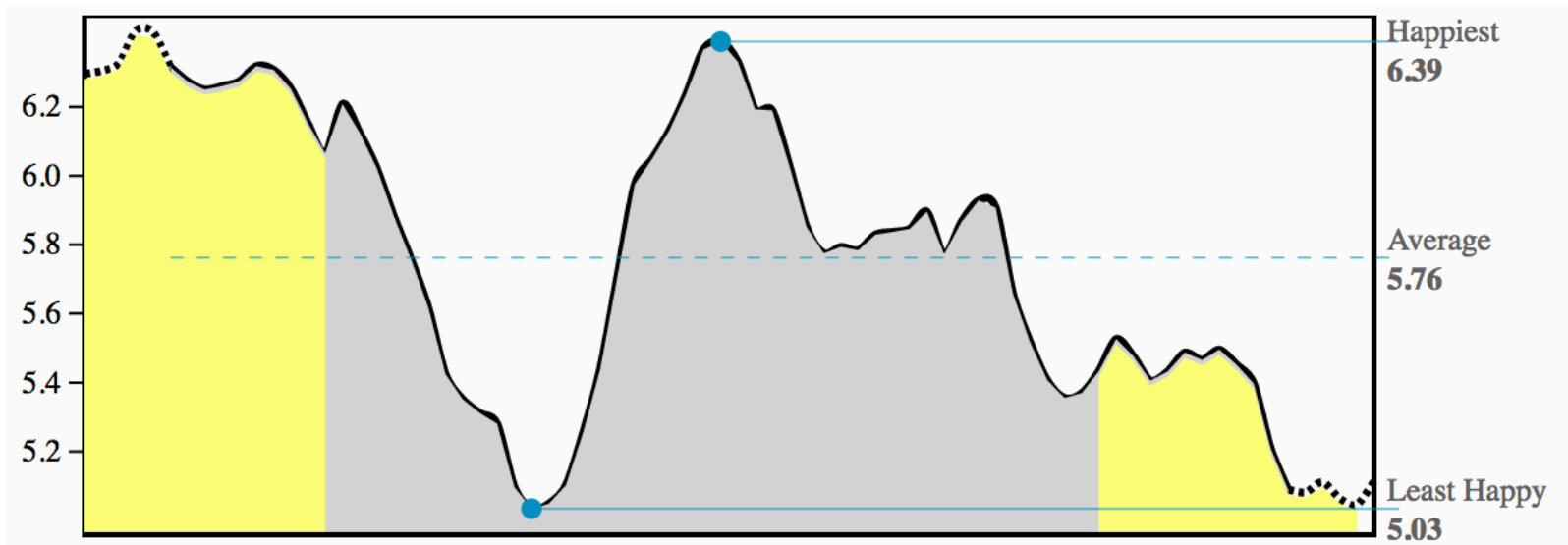


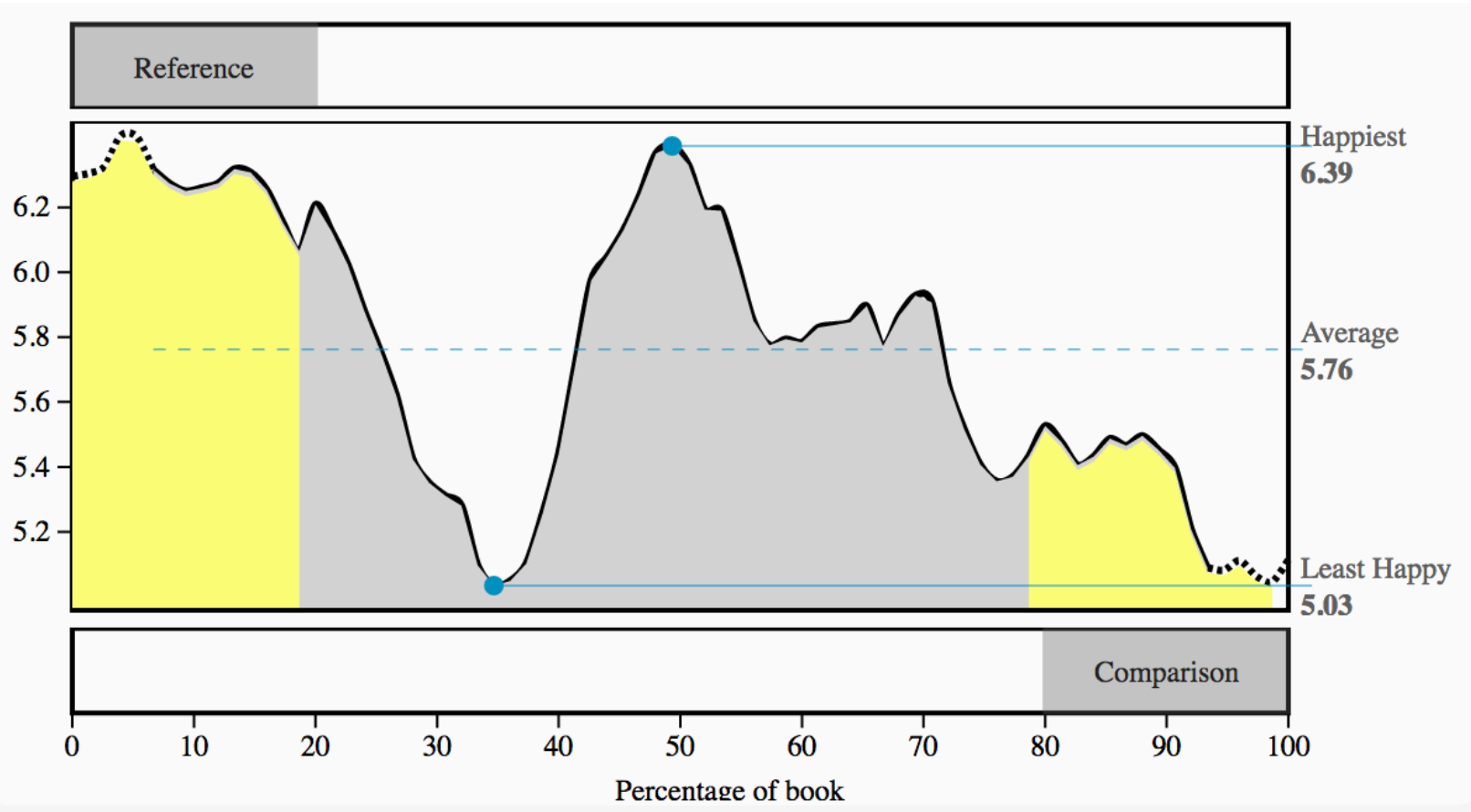
```

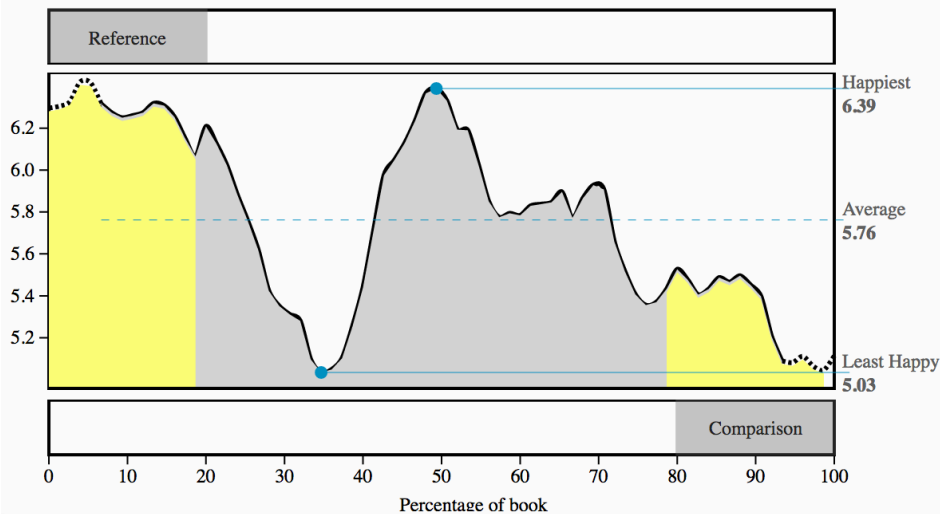
collapse_count = 0
# first thing, need to collapse everything onto user-days
# i'm going to hack the dataframe object to do this
target_users = dict()
for df_id, r in tqdm(q.iterrows()):
    if r.user_id in target_users:
        if r.from_diag in target_users[r.user_id]:
            target_users[r.user_id][r.from_diag] += " "+r.text
            collapse_count += 1
        else:
            target_users[r.user_id][r.from_diag] = r.text
    else:
        target_users[r.user_id] = dict()
        target_users[r.user_id][r.from_diag] = r.text
empty_count = 0
# turn each observation into a sparse, normalized word vector
for user in tqdm(target_users):
    for date in target_users[user]:
        target_users[user][date] = my_LabMT.wordVecify(dictify(listify(target_users[user][date])))
        # print(target_users[user][date].sum())
        # print(target_users[user][date])
        # print(target_users[user][date]/target_users[user][date].sum())
        if target_users[user][date].sum() != 0:
            target_users[user][date] = target_users[user][date]*10/(target_users[user][date].sum())
        else:
            empty_count+=1
    target_users[user][date] = csr_matrix(target_users[user][date])
control_users = dict()
collapse_count_control = 0
for df_id, r in tqdm(p.iterrows()):
    if r.user_id in control_users:
        if r.created_date in control_users[r.user_id]:
            control_users[r.user_id][r.created_date] += " "+r.text
            collapse_count_control += 1
        else:
            control_users[r.user_id][r.created_date] = r.text
    else:
        control_users[r.user_id] = dict()
        control_users[r.user_id][r.created_date] = r.text
for user in tqdm(control_users):
    for date in control_users[user]:
        control_users[user][date] = my_LabMT.wordVecify(dictify(listify(control_users[user][date])))
        if control_users[user][date].sum() != 0:
            control_users[user][date] = control_users[user][date]*10/control_users[user][date].sum()
        control_users[user][date] = csr_matrix(control_users[user][date])

```

# Rich data comparison





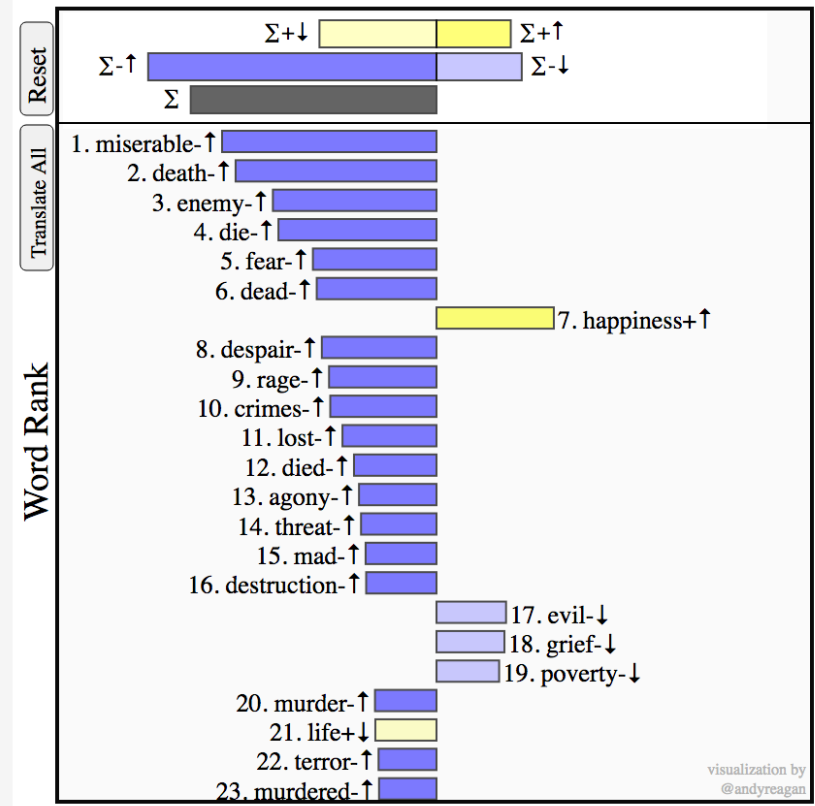


## Word Shift:

**Why comparison section is less happy than reference section:**

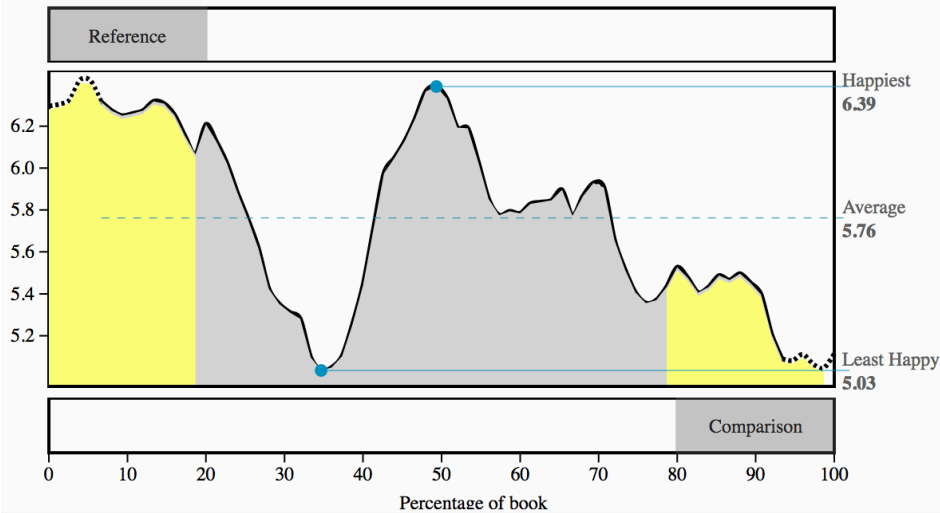
Reference section's happiness: 6.31

Comparison section's happiness: 5.34



Per word average happiness shift

visualization by  
@andyreagan

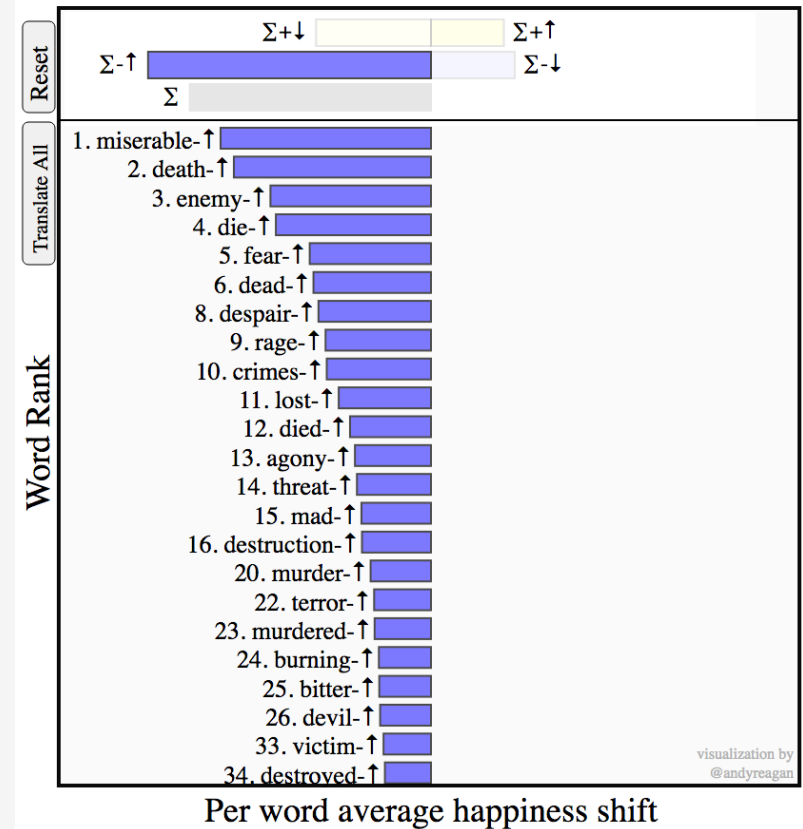


### Word Shift:

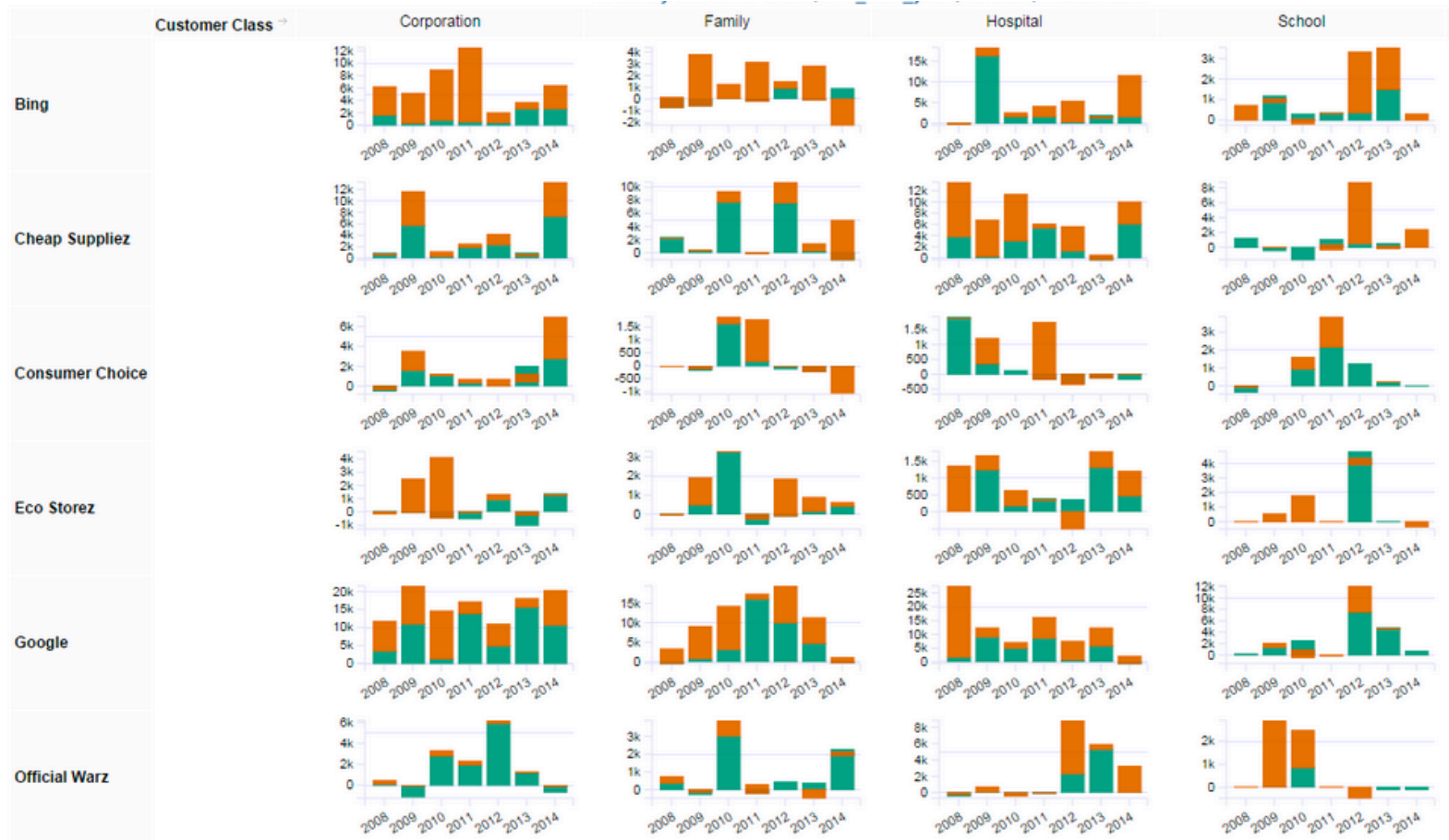
**Why comparison section is less happy than reference section:**

Reference section's happiness: 6.31

Comparison section's happiness: 5.34



# Multifaceted views





# Integrated statistical calculations

# Data access and integration

Berkeley SCHOOL OF  
INFORMATION