# Expanding Our Text Mining Toolkit

Sentiment Analysis and Term Selection in JMP Pro 16

Ross Metusalem, PhD

JMP Systems Engineer

jmp. STATISTICAL DISCOVERY FROM SAS
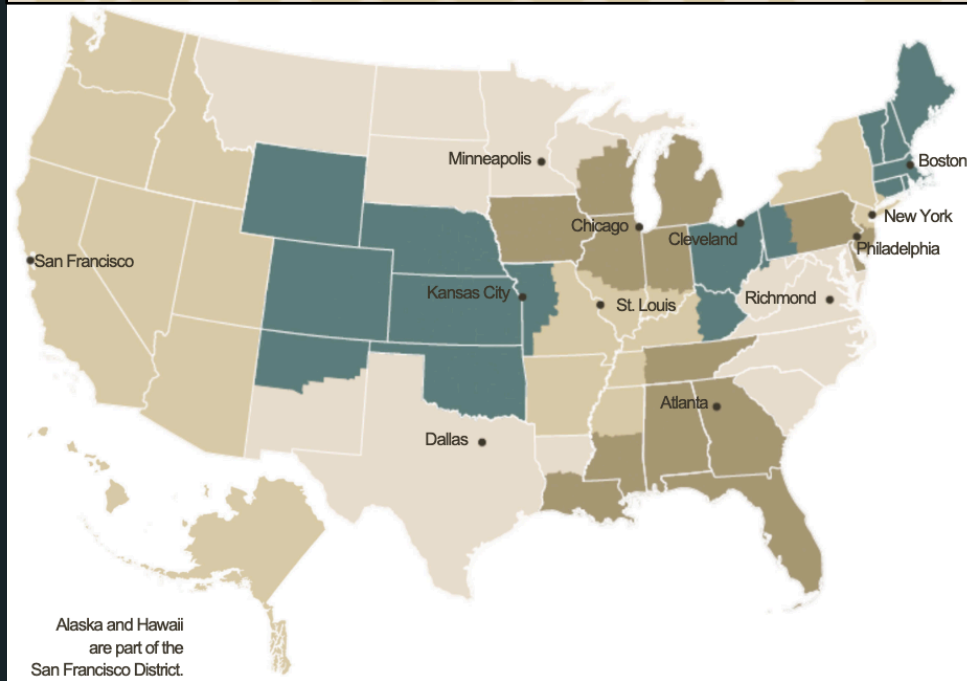
# Text Explorer

# Sentiment Analysis

How emotionally positive or negative is a text?

# The Beige Book

## Summary of Commentary on Current Economic Conditions By Federal Reserve District

Minneapolis

Boston

Chicago

Cleveland

New York

Philadelphia

San Francisco

Kansas City

St. Louis

Richmond

Atlanta

Dallas

Alaska and Hawaii
are part of the
San Francisco District.

locations by major retail chains. Computer sales, on the other hand, have been severely depressed and are not expected to recover until yearend or early 1972. A second director, who is the chairman of a large diversified manufacturing firm (major divisions include auto components, defense and space products, and electronics), expressed the view that economic activity will probably improve during 1971. He also reported that there are a number of favorable "straws in the wind" from his firm's point of view:

there is considerable agreement that (1) consumers remain pessimistic, (2) outlays for plant and equipment are being cut back or deferred, (3) labor markets are easing but wage costs are not, and (4) growth in the money supply is too fast. Reports from around the District indicate widespread pessimism on the part of consumers. Large department stores in the region report poor sales for large luxury items all the way down to small inexpensive goods. One department store executive quipped that "... only

jmp STATISTICAL DISCOVERY FROM SAS

# Sentiment Term Scores

| Term | Score |
|------|------|
| fantastic | 90 |
| favorable | 40 |
| feared | -60 |
| fortunate | 40 |
| friendly | 40 |
| functional | 20 |
| garbage | -60 |
| glad | 60 |
| good | 60 |
| gracious | 70 |
| great | 80 |
| greatest | 90 |
| grim | -75 |

✕

# Intensifiers

| Term | Multiplier |
|------|------|
| incredibly | 1.40 |
| insanely | 1.90 |
| little | 0.30 |
| mightily | 1.60 |
| more | 1.20 |
| most | 1.50 |
| much | 1.20 |
| not only | 1.00 |
| not the only | 1.00 |
| only | 0.30 |
| outrageously | 1.70 |
| over | 1.00 |
| overly | 1.00 |

✕

# Negators (-1)

| Term |
|------|
| no |
| non |
| none |
| nor |
| not |
| shouldn't |
| shouldn't |
| wasn't |
| wasn't |
| weren't |
| weren't |
| without |
| won't |

jmp STATISTICAL DISCOVERY FROM SAS

# Leading indicator of recession?



Sentiment Over Time

# The Dotcom Bust

# The Great Recession

# Some Applications of Sentiment Analysis

Consumer research

Product improvement

Customer support

Public policy

jmp. STATISTICAL DISCOVERY
FROM SAS

# Let's see it in JMP

JMP STATISTICAL DISCOVERY
FROM SAS

# Term Selection

Which words are most strongly associated with an important variable?

# Which words are associated with recessions?

# Document-Term Matrix

# Generalized Regression

$$Log\left(\frac{P(Recession)}{P(No\ Recession)}\right) = -0.89 + (1.94 * "pandemic") - (1.02 * "gain")\ ...$$

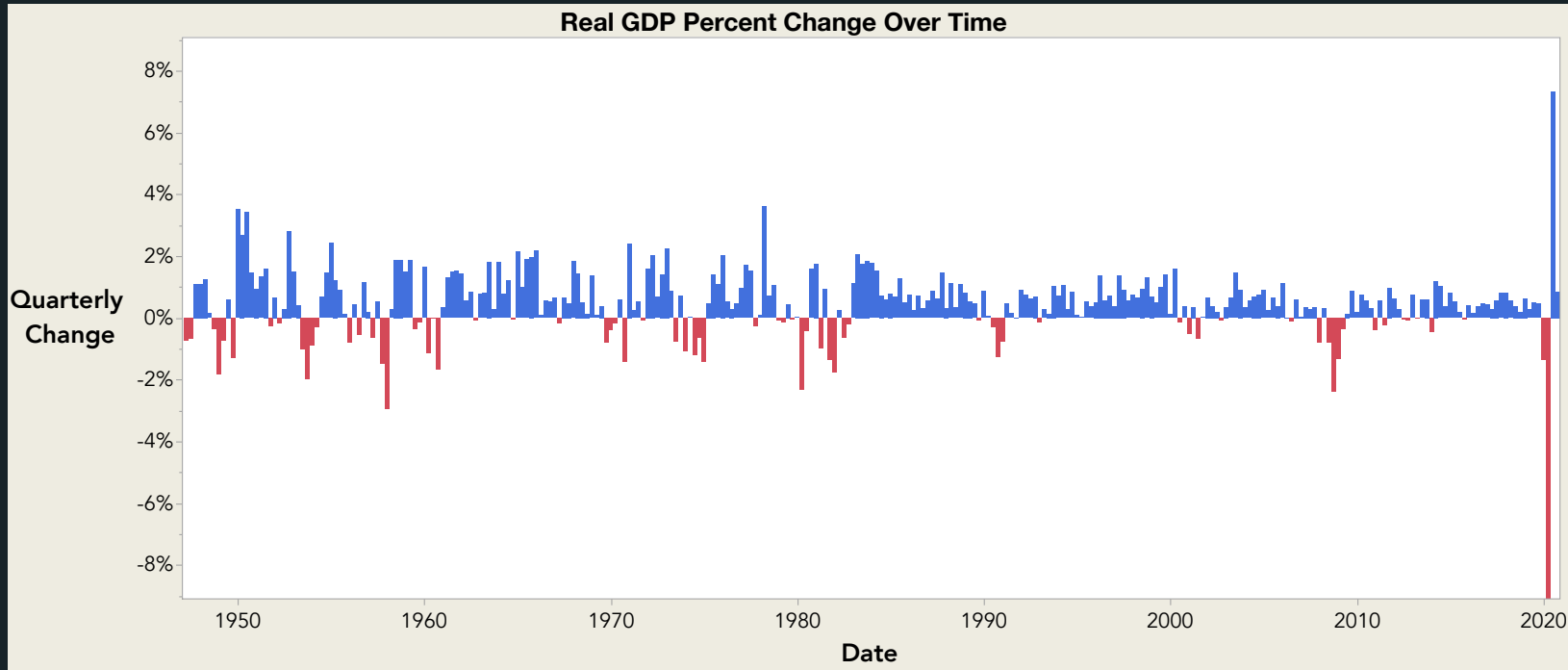## Terms associated with <u>recession</u>

| Term | Coefficient ⌄ | LogWorth | Count |
|---|---|---|---|
| pandemic | 1.942 | 1.677 | 33 |
| postpon· | 0.922 | 1.531 | 34 |
| cancel· | 0.754 | 1.321 | 60 |
| foreclosur· | 0.688 | 0.979 | 33 |
| deterior· | 0.682 | 1.974 | 187 |
| pessimist· | 0.624 | 0.853 | 45 |

## Terms associated with <u>NOT recession</u>

| Term | Coefficient ⌃ | LogWorth | Count |
|---|---|---|---|
| gain· | -1.027 | 4.770 | 657 |
| strengthen· | -0.608 | 1.952 | 416 |
| competit· | -0.536 | 1.311 | 291 |
| manufactur· | -0.207 | 0.129 | 2922 |
| acceler· | -0.194 | 0.383 | 181 |
| stronger | -0.168 | 0.261 | 331 |

jmp STATISTICAL DISCOVERY FROM SAS

Terms Assocated with Probability of Recession

# Let's see it in JMP

jmp STATISTICAL DISCOVERY
FROM SAS

# Term Selection

Identify words associated with an important variable

## Terms associated with <u>recession</u>

| Term | Coefficient | LogWorth | Count |
|---|---|---|---|
| pandemic | 1.942 | 1.677 | 33 |
| postpon· | 0.922 | 1.531 | 34 |
| cancel· | 0.754 | 1.321 | 60 |
| foreclosur· | 0.688 | 0.979 | 33 |
| deterior· | 0.682 | 1.974 | 187 |
| pessimist· | 0.624 | 0.853 | 45 |

## Terms associated with <u>NOT recession</u>

| Term | Coefficient | LogWorth | Count |
|---|---|---|---|
| gain· | -1.027 | 4.770 | 657 |
| strengthen· | -0.608 | 1.952 | 416 |
| competit· | -0.536 | 1.311 | 291 |
| manufactur· | -0.207 | 0.129 | 2922 |
| acceler· | -0.194 | 0.383 | 181 |
| stronger | -0.168 | 0.261 | 331 |

## Sentiment Analysis

Quantify positive-negative emotion in texts

### Sentiment Summary

| | N | Mean Score |
|---|---|---|
| All Scored Documents | 5267 | 4.1 |
| Net Positive Documents | 3836 | 22.8 |
| Net Negative Documents | 1409 | -15.4 |
| No Sentiment Documents | 1 | 0.0 |



Scores of Documents with Net Sentiment

jmp STATISTICAL DISCOVERY FROM SAS