

TopCATI Notes Sentiment Analysis

Proof of Concept and Preliminary Results

Matthew Burgess

Economist

CES Data Collection Branch

Bureau of Labor Statistics

FedCASIC 2015

3/5/2015



Overview

- CES Overview - TopCATI Data Collection and Notes
- Sentiment Analysis Proof of Concept (R and SAS)
- Preliminary Results
- Conclusions/Going Forward

Background on the Current Employment Statistics Survey

- The BLS Current Employment Statistics (CES) survey is also known as the payroll survey or the establishment survey
- The CES Survey is the Largest multi-modal survey
 - Survey of about 143,000 businesses and government agencies, representing approximately 588,000 individual worksites
- CES publishes employment, hours, and earnings data for the nation, states, and metropolitan areas at total nonfarm and detailed industry levels

CES Data Collection

- Monthly survey of establishments to collect **employment, hours, and earnings** data
- Collect employment data for the pay period that includes the 12th of the month

Computer Assisted Telephone Interviewing

- 26% of CES reports are collected by CATI. We also use other methods like web and fax
- Interviewers at Data Collection Centers (Atlanta, Dallas, Kansas City, and Fort Walton Beach) call respondents and collect data using TopCATI software
- TopCATI software allows interviewers to take notes, schedule call times, and review reported data

Project Background

- CES is required to permanently save TopCATI interviewer notes
- TopCATI interviewer notes contain valuable, qualitative information about businesses participating in the CES survey
- Data mining techniques such as sentiment analysis may be used to quantify information contained in the notes

Sentiment Analysis Explained

- Estimating sentiment
 - ▶ Many complex algorithms have been developed – all are limited by the computer's ability to interpret language
 - ▶ Pros: facilitates analysis of millions of text notes in a short period of time
 - ▶ Cons: computers will incorrectly interpret nuanced phrases, sarcasm, etc.

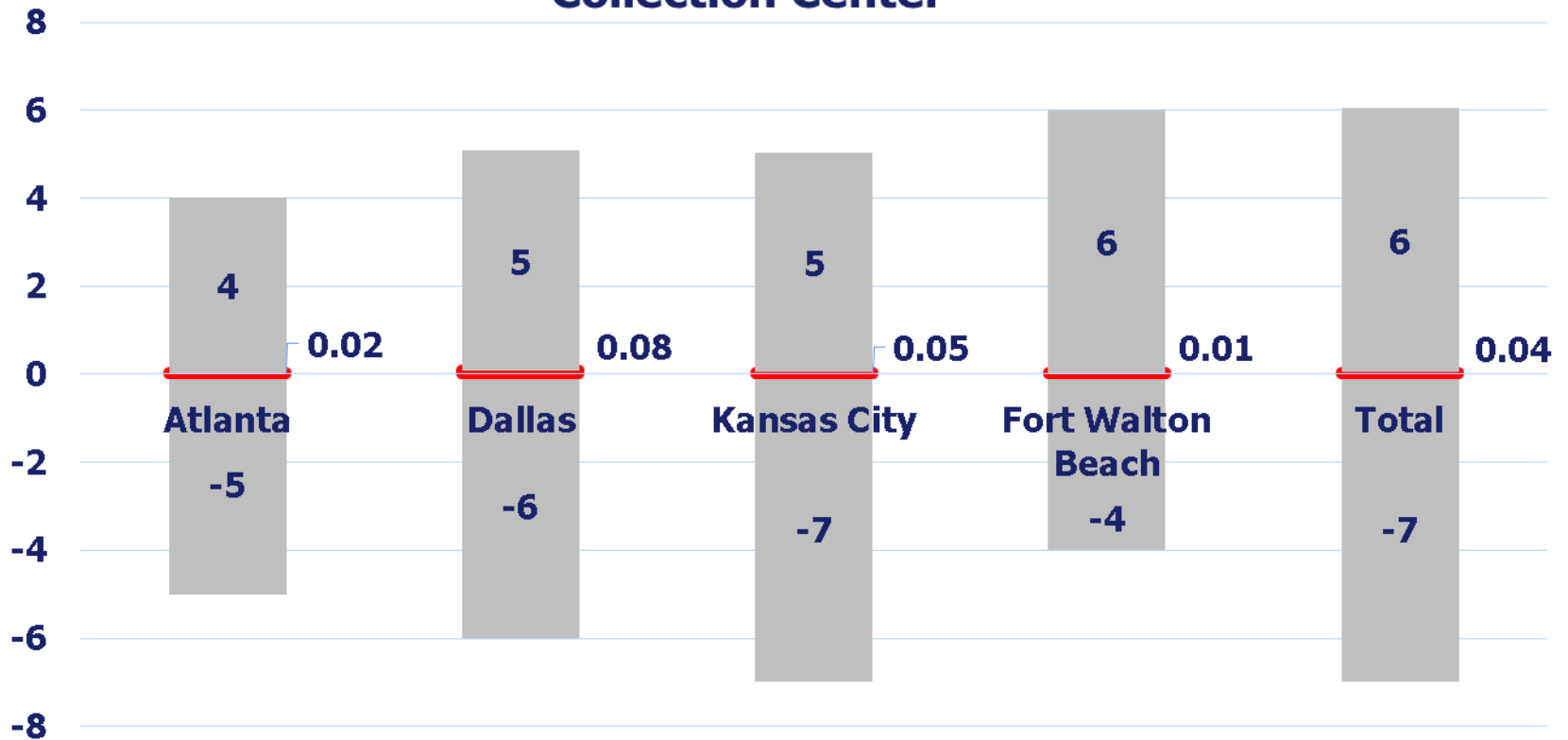
Project Methods

■ Project methods

- ▶ This project used an algorithm that counts the number of “positive” and “negative” words and computes an overall sentiment score for each note
- ▶ This analysis used a list of positive and negative words categorized by researchers Hu and Liu in their “opinion lexicon” of about 6,800 words
- ▶ The project was a proof of concept exercise and only analyzed notes from Wisconsin respondents. Sentiment Scoring code was developed by CES staff in both SAS and R; both software systems produced matching results

Preliminary Results

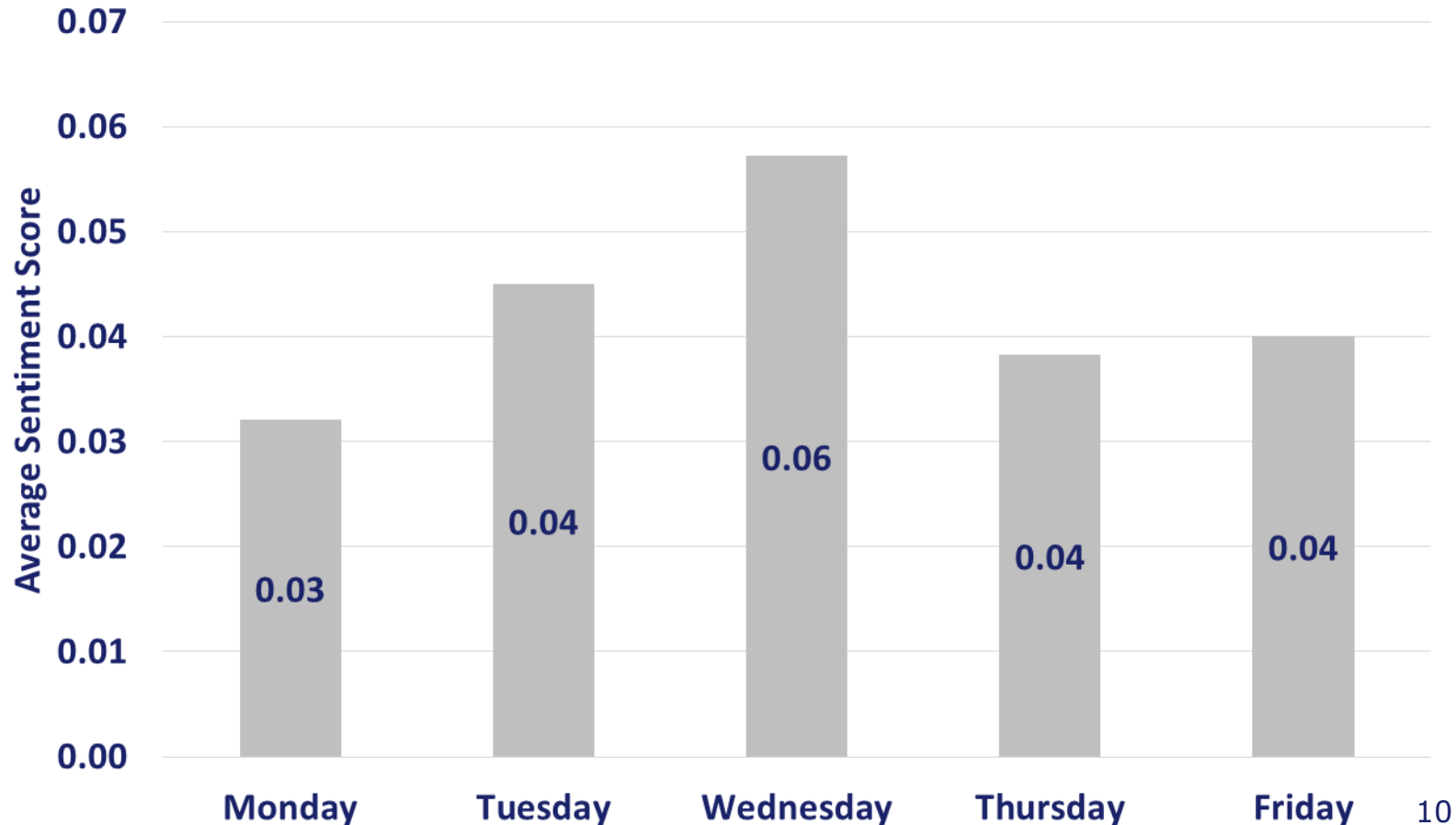
Min, Max, and Average Sentiment Score by Data Collection Center



■ Average Score

Preliminary Results

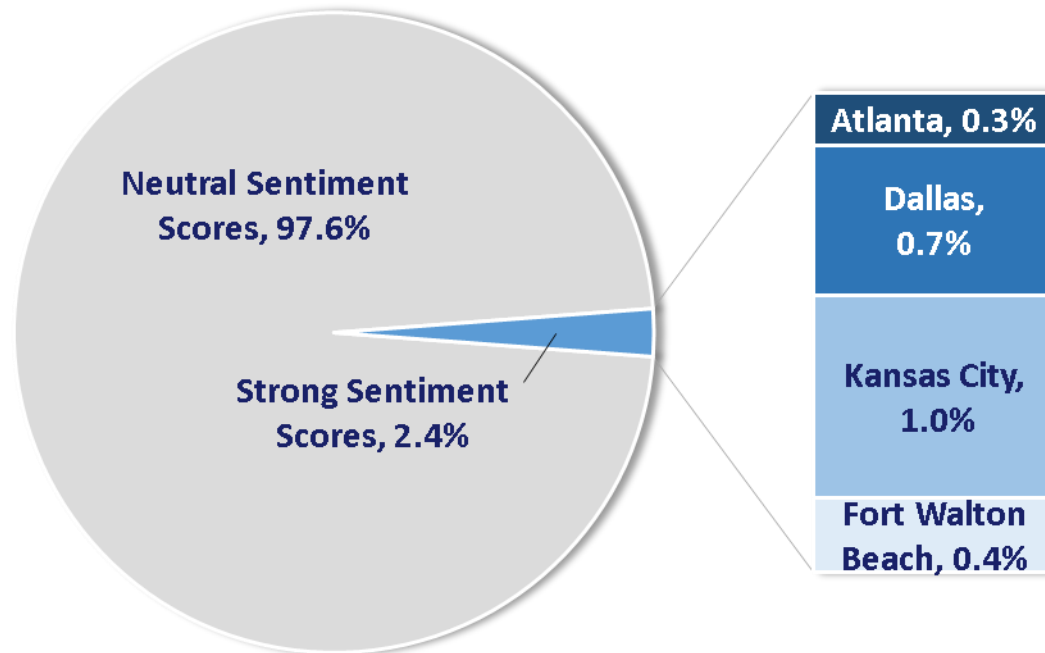
Average Sentiment Score by Day of the Week



Preliminary Results

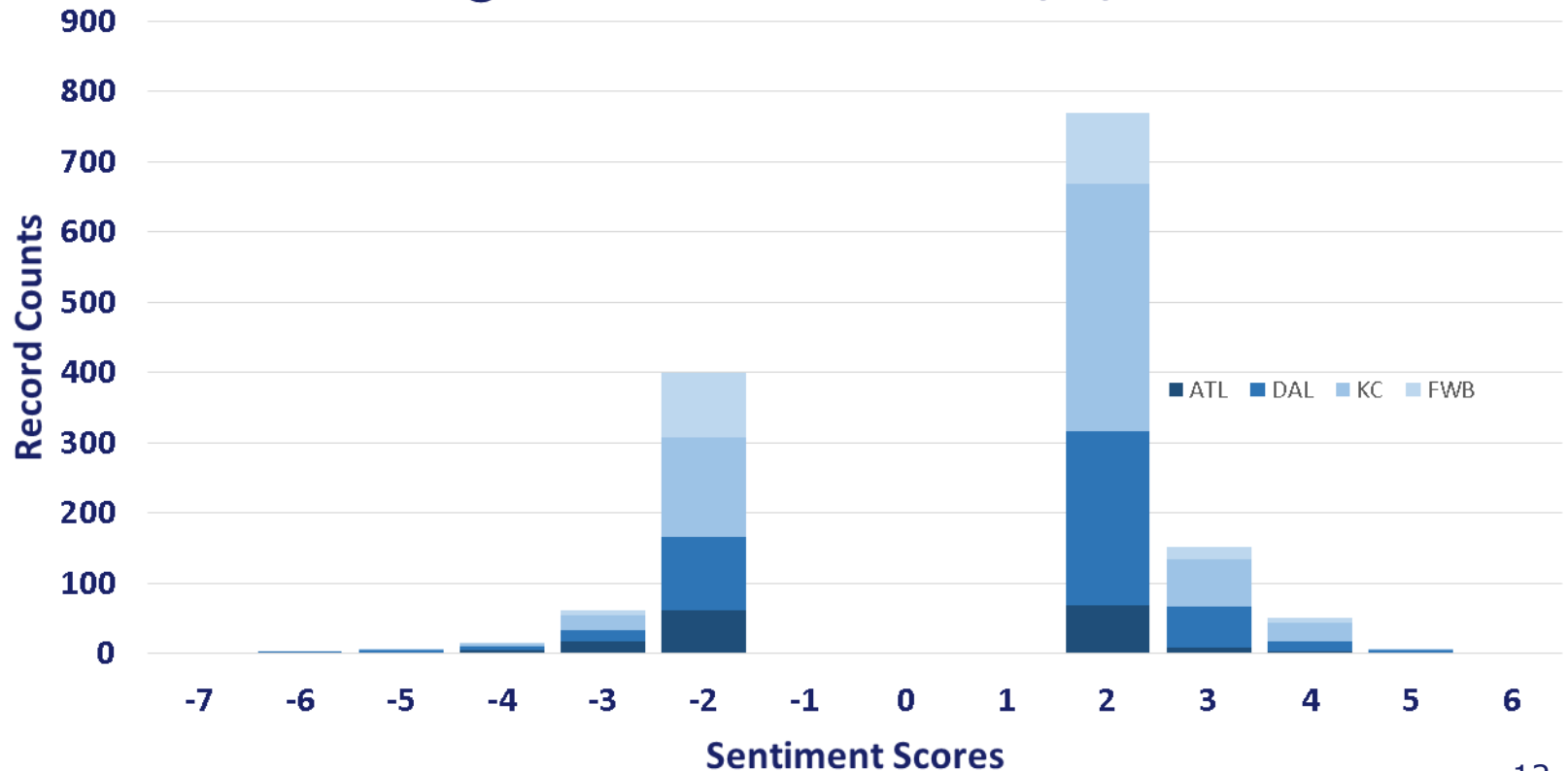
- Sentiment scores are mostly “neutral” (zero, 1, or -1)

Percentage of Records with Neutral vs. Strong Sentiment Score
(Strong Sentiment Scores Broken Out by Call Center)



Preliminary Results

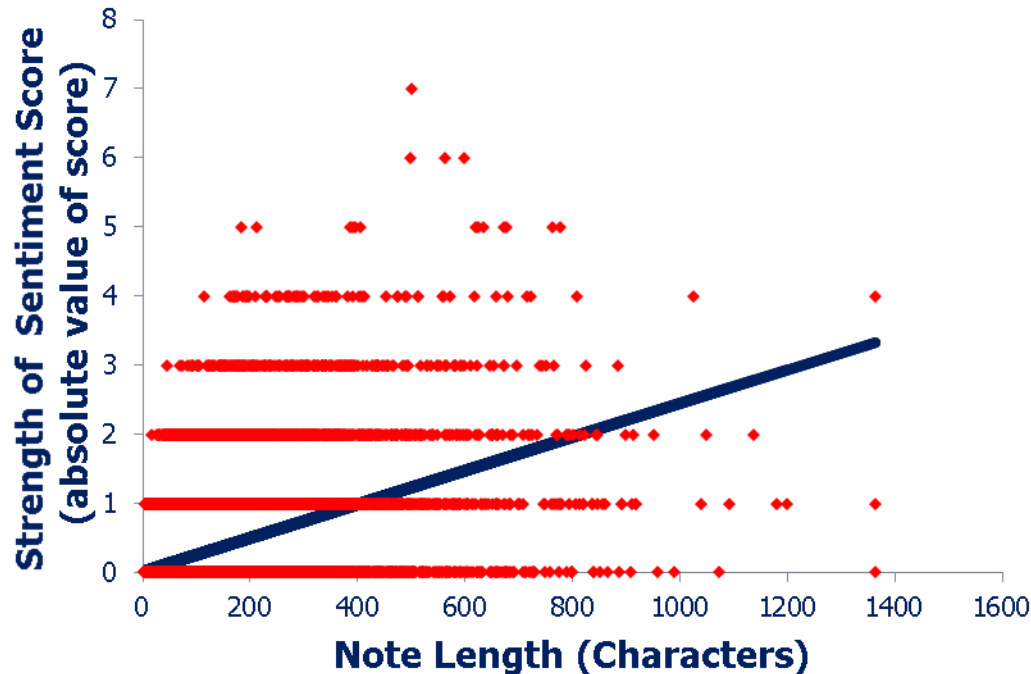
Sentiment Score Distribution by Call Center, excluding "neutral" scores of 1, 0, and -1



Preliminary Results

- Longer notes have better chance of getting a “strong” sentiment score

Strength of Sentiment Score vs. Length of Note



Regression Equation:

y = score strength

x = note length

$y = 0.0146 + 0.00244x$

95% Confidence Interval (Length

Coefficient): (0.00239, 0.00248)

R-Square: 0.18396

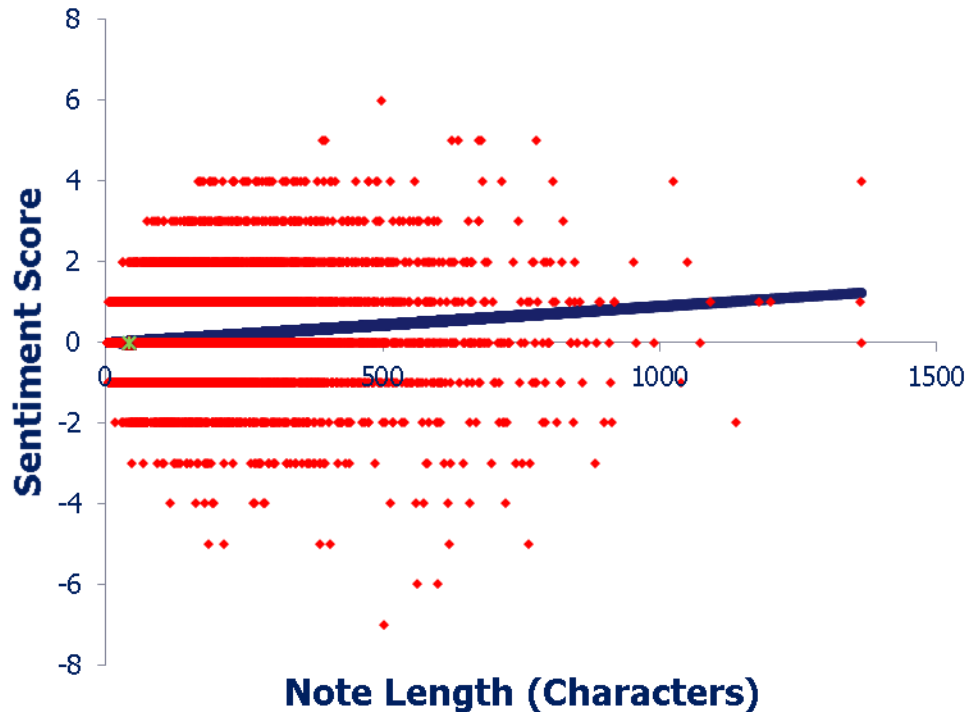
◆ Absolute Value Sentiment Score

— Predicted Absolute Value Sentiment Score

Preliminary Results

- Longer notes also have a better chance of getting a positive sentiment score

Sentiment Score vs. Note Length



Regression Equation:

$y = \text{sentiment score}$

$x = \text{note length}$

$y = -0.0225 + 0.0009124x$

95% Confidence Interval (Length

Coefficient): (0.0008648, 0.0009601)

R-Square :0.02253

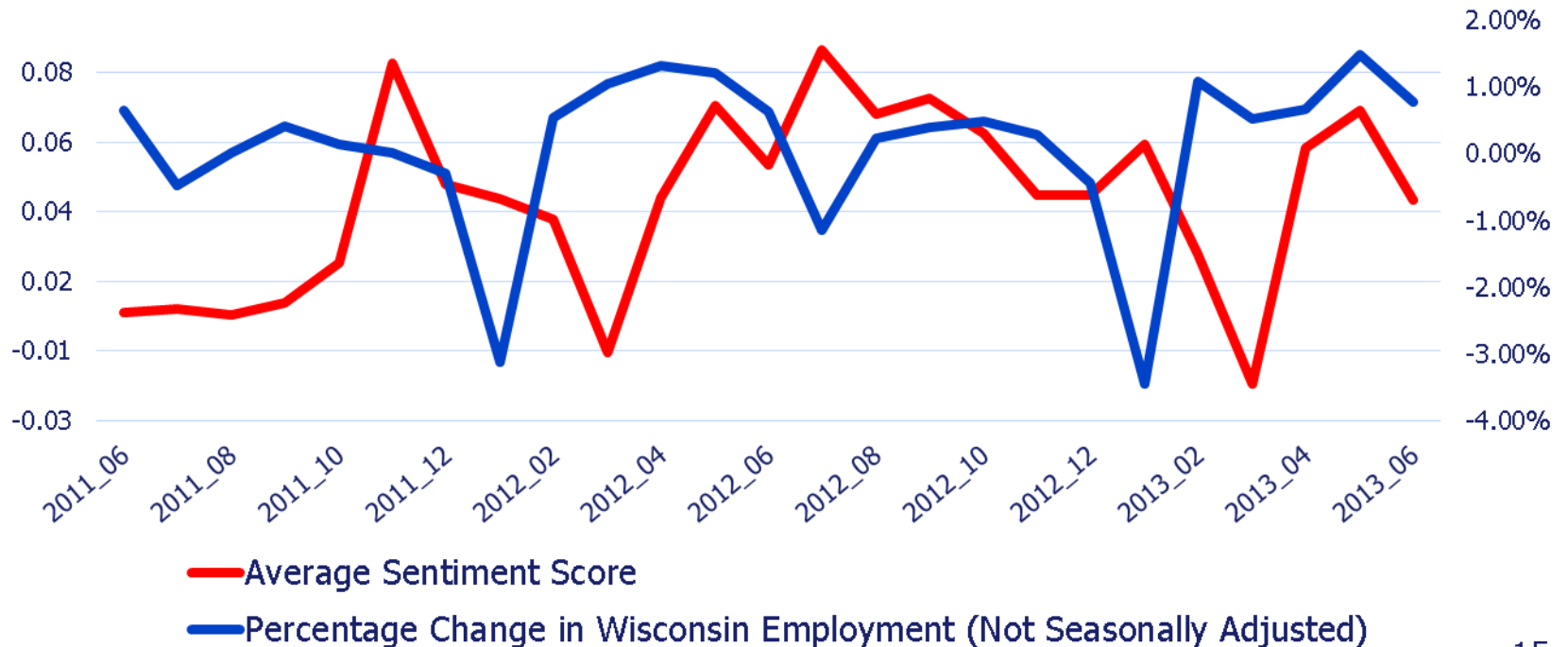
◆ Sentiment Score

— Predicted Sentiment Score

Preliminary Results

- Average score compared to change in employment over time

Average Sentiment Score vs. Percentage Change in Wisconsin Employment (Not Seasonally Adjusted)



Conclusions

- Our preliminary sentiment score proof of concept was successful
- Sentiment analysis allowed us to quantify information from existing TopCATI notes
 - ▶ Old notes had previously been taking up storage space with unusable qualitative information
- Sentiment analysis algorithms can be applied to existing qualitative BLS data at low cost/resources

Going Forward

- Expand this research to other states and time periods for further analysis
- More targeted/complex algorithms
 - ▶ Creating a CES specific “positive” and “negative” word list
 - ▶ Using a scaling system (degree of negativity/positivity of each word)
- Seasonally adjust Sentiment Scores to compare with seasonally adjusted CES data
 - ▶ Or compare yearly average of sentiment scores and CES Data

Contact Information

Matthew Burgess

Economist

OEUS, CES Data Collection Branch

www.bls.gov/ces

202-691-6519

Burgess.Matthew@bls.gov



Appendix

- Analysis was done on a sample of about 61,000 notes from Wisconsin
- Stats by call center:

Averages by Call Center								
Call Center	Min Score	Avg Score	Max Score	Avg Length	Max Length	Strong Sentiment	Total Record Count	% with Strong Sentiment
Atlanta	-5	0.02	4	74	958	163	7404	2.2%
Dallas	-6	0.08	5	86	1200	455	13451	3.4%
Kansas City	-7	0.05	5	69	1364	618	28581	2.2%
Fort Walton Beach	-4	0.01	6	58	1049	229	11686	2.0%
Total	-7	0.04	6	71	1364	1465	61122	2.4%