

Distillation of News Flow into Analysis of Stock Reactions

Junni Zhang

Cathy Chen

Wolfgang Karl Härdle

Elisabeth Bommers

Guanghua School of Management

Peking University

Chung Hua University

Ladislaus von Bortkiewicz Chair of Statistics

C.A.S.E. – Center for Applied Statistics
and Economics

Humboldt-Universität zu Berlin

<http://lvb.wiwi.hu-berlin.de>

<http://www.case.hu-berlin.de>



News moves Markets...

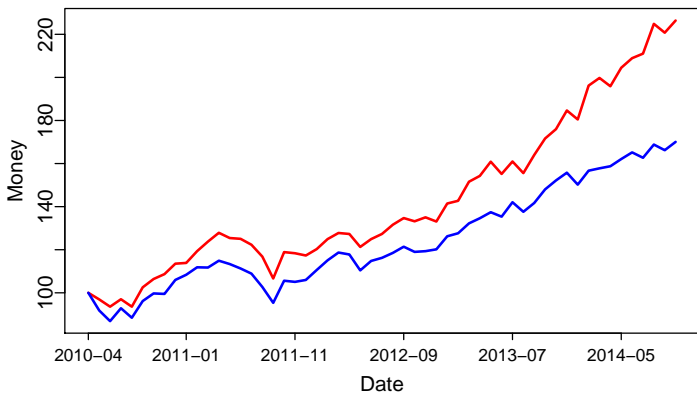


Figure 1: Investment in: S&P 500, Sentiment Strategy



... but there is a lot of News



Dimensions of News

- Source of news
 - ▶ Official channel: government, federal reserve bank/central bank, financial institutions
 - ▶ **Internet**: blog, social media, message board
- Content of news
 - ▶ Signal v.s. noise



Dimension of News ctd

- Type of news
 - ▶ Scheduled v.s. **non-scheduled**
 - ▶ Expected v.s. unexpected
 - ▶ Specific-event v.s. **continuous news flows**

Challenge

- Interpret news
- Evaluate news impact from different news dimensions



Sentiment Projection

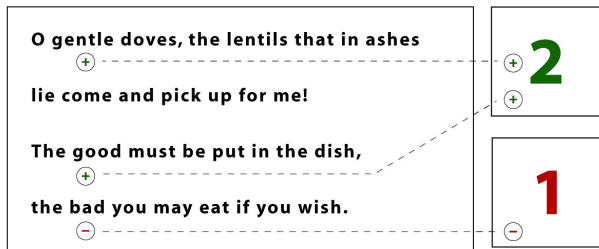


Figure 2: Example of Text Numerization

- Many texts are numerized via lexical projection
- Goal: Accurate values for positive and negative sentiment



Sentiment Lexica

- ▣ *Opinion Lexicon* (BL)
Hu and Liu (2004)
- ▣ *Financial Sentiment Dictionary* (LM)
Loughran and McDonald (2011)
- ▣ *Multi-Perspective Question Answering Subjectivity Lexicon* (MPQA)
Wilson et al. (2005)



Research Questions

- Do opinions of small traders contribute to stock markets and create news-driven stock reactions?
 - ▶ Small traders v.s. financial institutions
 - ▶ Opinions of small traders v.s. financial analysts

- Concerns for analyst recommendation
 - ▶ Career
 - ▶ Compensation scheme
 - ▶ Strategic alliance



Research Questions ctd

- Are there differences regarding
 1. stock reaction indicators: volatility, trading volume, returns?
 2. degree of asymmetric response (leverage effect)?
 3. high and low attention companies?
 4. specific sectors?



Outline

1. Motivation ✓
2. Data Collection
3. Sentiment Projection
4. Panel Regression
5. Simulation
6. Conclusion



How to gather Sentiment Variables?

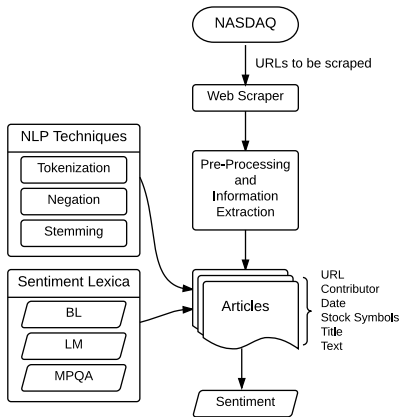


Figure 3: Flowchart of Data Gathering Process

Distillation of News Flow into Analysis of Stock Reactions



NASDAQ Articles

- ▣ Terms of Service permit web scraping
- ▣ 116,691 articles in total
- ▣ 43,459 articles about 100 selected S&P 500 stocks in 9 major GICS sectors [Frequency Table: GICS](#)
- ▣ Time frame: October 2009 - October 2014
- ▣ Data available at [RDC](#)



Sentiment Lexica ctd

- Number of entries in each lexicon:

| Lexicon | Positive | Negative |
|---------|----------|----------|
| BL | 2,006 | 4,783 |
| LM | 354 | 2,329 |
| MPQA | 2,718 | 4,911 |

- Some words appear only in one lexicon
- Other words are only found in two lexica

Unique Words

Shared Words



Sentiment Variables

- ▣ $I_{i,t}$ - article indicator
- ▣ $Pos_{i,t}$ - average proportion of positive words
- ▣ $Neg_{i,t}$ - average proportion of negative words

for stock i on day t



Comparison of Lexical Projections

- Average sentiment values are smaller for LM than for BL and MPQA
- *BL* and *MPQA* relatively similar
- *LM* only contains finance specific words
- *BL* and *MPQA* also contain more general words (e.g. "cancer")

Summary Statistics

Correlation - Sentiment

Tagging Example

- Combination of projections might improve results
 - ▶ PCA on sentiment scores
 - ▶ Use first principal component of $Pos_{i,t}$ and $Neg_{i,t}$



How good are the Projections?

- Random selection of 100 articles, manual labeling of polarity and comparison with polarity of lexical projections
- *BL* and *MPQA* recognize fewer negative articles but good in detection of positive articles
- *LM* accurately detects negative articles, recognizes fewer positive articles

Classification Evaluation Table



Stock Reaction Indicators

Range-based measure of volatility by Garman and Klass (1980)

- Notation: $\sigma_{i,t}$ Computation
- Use $\log \sigma_{i,t}$
- Based on open-high-low-close prices
- Equivalent results to realized volatility
- More robust in case of microstructure effects



Detrended log trading volume by Girard and Biswas (2007)

$$V_{i,t} = V_{i,t}^* - (\alpha + \beta_{1,i}(t - t_0) + \beta_{2,i}(t - t_0)^2) \quad (1)$$

with raw log trading volume $V_{i,t}^*$ and detrended log trading volume $V_{i,t}$ for stock i on day t , t_0 is starting point of time window (size: 120 days)

Returns

$$R_{i,t} = \log(P_{i,t}^C) - \log(P_{i,t-1}^C) \quad (2)$$

with $P_{i,t}^C$ as closing price of stock i on day t



Panel Regression

$$\log \sigma_{i,t+1} = \alpha + \beta_1 I_{i,t} + \beta_2 Pos_{i,t} + \beta_3 Neg_{i,t} + \beta_4^T X_{i,t} + \gamma_i + \varepsilon_{i,t} \quad (3)$$

$$V_{i,t+1} = \alpha + \beta_1 I_{i,t} + \beta_2 Pos_{i,t} + \beta_3 Neg_{i,t} + \beta_4^T X_{i,t} + \gamma_i + \varepsilon_{i,t} \quad (4)$$

$$R_{i,t+1} = \alpha + \beta_1 I_{i,t} + \beta_2 Pos_{i,t} + \beta_3 Neg_{i,t} + \beta_4^T X_{i,t} + \gamma_i + \varepsilon_{i,t} \quad (5)$$

for stock i on day t with separate estimation of (3) to (5).

$X_{i,t}$ - control variables

γ_i - company specific fixed effect satisfying $\sum_i \gamma_i = 0$



Control Variables

- ▣ $R_{M,t}$ - S&P 500 index return
- ▣ VIX_t - CBOE VIX
- ▣ $\log \sigma_{i,t}$ - Range-based volatility
- ▣ $V_{i,t}$ - Detrended trading volume
- ▣ $R_{i,t}$ - Return



Entire Panel Regression Results

| Variable | BL | LM | MPQA | PCA |
|--|-----------|-----------|----------|----------|
| Panel A: Future Log Volatility $\log \sigma_{i,t+1}$ | | | | |
| $I_{i,t}$ | -0.005 | -0.019*** | -0.004 | -0.014 |
| $Pos_{i,t}$ | -0.396* | 0.156 | -0.517** | -0.210 |
| $Neg_{i,t}$ | 0.905*** | 0.942*** | 1.464*** | 1.041*** |
| Panel B: Future Detrended Log Trading Volume $V_{i,t+1}$ | | | | |
| $I_{i,t}$ | 0.040*** | 0.027*** | 0.046*** | 0.035*** |
| $Pos_{i,t}$ | -0.496*** | 0.051 | -0.483** | -0.274* |
| $Neg_{i,t}$ | 0.726*** | 0.563** | 0.548* | 0.590** |
| Panel C: Future Returns $R_{i,t+1}$ | | | | |
| $I_{i,t}$ | 0.000 | 0.000 | 0.000 | -0.000 |
| $Pos_{i,t}$ | 0.019*** | 0.030*** | 0.014* | 0.018*** |
| $Neg_{i,t}$ | -0.004 | -0.000 | -0.009 | -0.003 |

*** p value < 0.01 , ** $0.01 \leq p$ value < 0.05 , * $0.05 \leq p$ value < 0.1



Does Attention matter?

- Number of days with articles differs between firms
- High attention: Faster incorporation of news?

$$\text{attention ratio} \stackrel{\text{def}}{=} N_i / T \quad (6)$$

with N_i as number of days with at least one article for company i
and T as total number of trading days



Grouping

Use attention ratio quartiles to group firms:

| | |
|----------------|----------------------------------|
| Low | attention ratio $<$ Q1 |
| Median | $Q1 \leq$ attention ratio $<$ Q2 |
| High | $Q2 \leq$ attention ratio $<$ Q3 |
| Extremely High | $Q3 \leq$ attention ratio |

with Q1, Q2, Q3 as first, second and third quartile



Attention Analysis Regression Results

| Attention | BL | | LM | | MPQA | |
|--|---------|-----------|----------|-----------|--------|-----------|
| | Low | High | Low | High | Low | High |
| Panel A: Future Volatility $\log \sigma_{i,t+1}$ | | | | | | |
| $I_{i,t}$ | 0.020 | -0.016 | 0.010 | -0.046*** | 0.016 | -0.019 |
| $Pos_{i,t}$ | -0.736 | -0.460 | -1.027 | 0.967 | -0.655 | -0.636** |
| $Neg_{i,t}$ | -0.074 | 1.324*** | -0.195 | 1.806*** | -0.195 | 2.548*** |
| Panel B: Future Detrended Log Trading Volume $V_{i,t+1}$ | | | | | | |
| $I_{i,t}$ | 0.054** | 0.036*** | 0.044*** | 0.021* | 0.049* | 0.046*** |
| $Pos_{i,t}$ | -0.817 | -0.198 | -0.923 | 0.815* | 0.0433 | -0.358 |
| $Neg_{i,t}$ | 0.312 | 0.554 | -0.109 | 0.447 | -0.197 | 0.419 |
| Panel C: Future Returns $R_{i,t+1}$ | | | | | | |
| $I_{i,t}$ | 0.000 | 0.000 | 0.000 | 0.001** | 0.000 | 0.000 |
| $Pos_{i,t}$ | 0.012 | 0.028** | 0.021 | 0.038* | 0.010 | 0.024** |
| $Neg_{i,t}$ | 0.009 | -0.034*** | -0.001 | -0.046** | -0.016 | -0.044*** |

*** p value < 0.01, ** $0.01 \leq p$ value < 0.05, * $0.05 \leq p$ value < 0.1



Sector Analysis

- Compare financials sector with health care sector
- Attention ratio high for financials (0.413) and low for health care (0.287)
- In line with attention analysis:
 - ▶ Financials: significant parameters
 - ▶ Health care: not significant



Entire Panel Results

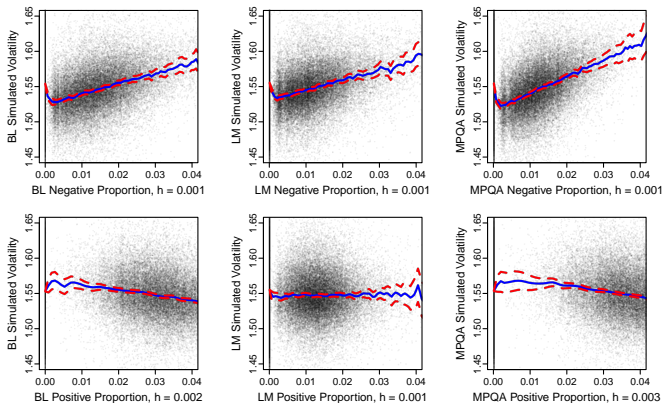



Figure 4: Volatility Simulation for Entire Panel: **Mean curve, 95% Uniform Confidence Bands**  **TXTSimulation** [Simulation Setup](#)



Entire Panel Results ctd

- Asymmetry effect
- *LM* and *MPQA*: curve for $Neg_{i,t}$ significantly differs from curve for $Pos_{i,t}$
 - ▶ Range *BL*: 0.023 - 0.056
 - ▶ Range *LM*: 0.017 - 0.039
 - ▶ Range *MPQA*: 0.023 - 0.05



Low Attention Results

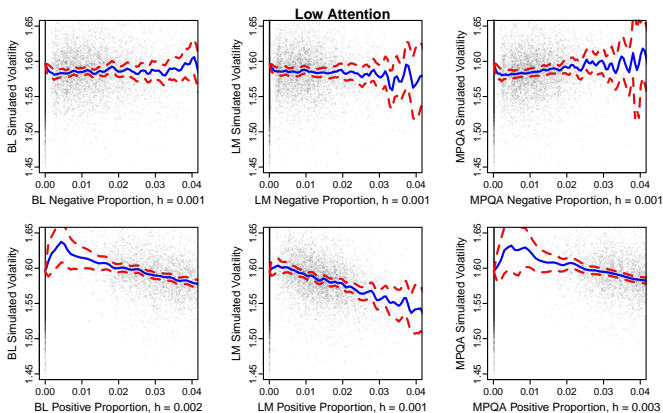



Figure 5: Volatility Simulation for Low Attention Group: **Mean curve, 95% Uniform Confidence Bands**  **TXTSimulationAttention** [Setup](#)



High Attention Results

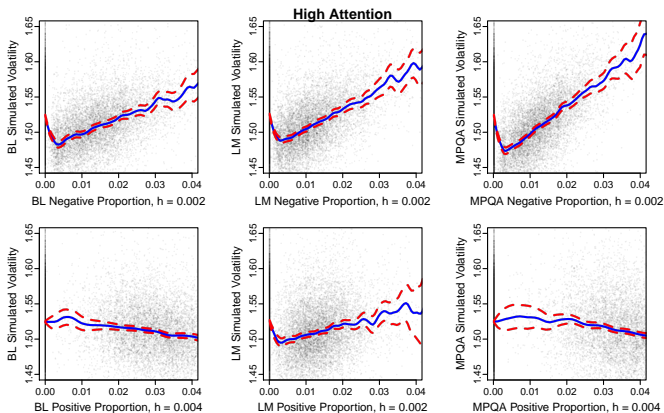



Figure 6: Volatility Simulation for High Attention Group: Mean curve, 95% Uniform Confidence Bands  TXTSimulationAttention

Distillation of News Flow into Analysis of Stock Reactions



Attentions Results ctd

- Low: no asymmetry effect
- High: *LM* and *MPQA*: curve for $Neg_{i,t}$ significantly differs from curve for $Pos_{i,t}$
 - ▶ Range *BL*: 0.022 - 0.056
 - ▶ Range *LM*: 0.019 - 0.024
 - ▶ Range *MPQA*: 0.020 - 0.053



Conclusion

- ▣ Sentiment measures: incremental information about future stock reactions
- ▣ Asymmetric impact of positive and negative sentiment
- ▣ Degree of incremental information and asymmetry is sector and attention specific
- ▣ Choice of lexicon matters



Distillation of News Flow into Analysis of Stock Reactions

Junni Zhang

Cathy Chen

Wolfgang Karl Härdle

Elisabeth Bommers

Guanghua School of Management

Peking University

Chung Hua University

Ladislaus von Bortkiewicz Chair of Statistics

C.A.S.E. – Center for Applied Statistics
and Economics

Humboldt-Universität zu Berlin

<http://lvb.wiwi.hu-berlin.de>

<http://www.case.hu-berlin.de>



Frequency Table: GICS Sectors

| GICS Sector | No. Stocks |
|----------------------------|------------|
| Consumer Discretionary | 21 |
| Consumer Staples | 9 |
| Energy | 6 |
| Financials | 12 |
| Health Care | 15 |
| Industrials | 10 |
| Information Technology | 21 |
| Materials | 4 |
| Telecommunication Services | 2 |

[Back](#)

Number of unique Words

- Some words only in one lexicon: "unique words"
- Number of unique words that appear at least three time in the articles:

| Lexicon | Positive | Negative |
|---------|----------|----------|
| BL | 470 | 918 |
| LM | 267 | 916 |
| MPQA | 512 | 181 |



Most frequent Words unique to one Lexicon

| BL | | LM | | MPQA | |
|----------------------|-------------------|--------------------------|---------------------|--------------------|---------------------|
| Positive | Negative | Positive | Negative | Positive | Negative |
| Available (5,836) | Debt (12,540) | Opportunities (4,720) | Declined (9,809) | Just (17,769) | Low (12,739) |
| Led (5,774) | Fell (9,274) | Strength (4,393) | Dropped (4,894) | Help (17,334) | Division (5,594) |
| Lead (4,711) | Fool (5,473) | Profitability (4,174) | Late (4,565) | Profit (15,253) | Least (5,568) |
| Recovery (4,357) | Issues (3,945) | Highest (3,409) | Claims (3,785) | Even (13,780) | Stake (4,445) |
| Work (3,808) | Risks (2,850) | Greater (3,321) | Closing (3,604) | Deal (13,032) | Slightly (3,628) |

Words only appear in one of the lexica and frequencies are given in parentheses.

Back



Number of shared Words

- Some words are only shared by two lexica
- Number of shared words that appear at least three time in the articles:

| Lexicon | Positive | Negative |
|-------------|----------|----------|
| BL and LM | 131 | 322 |
| BL and MPQA | 971 | 1,164 |
| LM and MPQA | 32 | 30 |



Most frequent shared Words

| BL and LM | | BL and MPQA | | LM and MPQA | |
|--------------------------|----------------------|----------------------|--------------------|------------------------|----------------------|
| Positive | Negative | Positive | Negative | Positive | Negative |
| Gains (7,604) | Losses (5,938) | Free (133,395) | Gross (8,228) | Despite (7,413) | Against (8,877) |
| Gained (7,493) | Missed (3,165) | Well (3,0270) | Risk (7,471) | Able (5,246) | Cut (3,401) |
| Improved (7,407) | Declining (3,053) | Like (24,617) | Limited (5,884) | Opportunity (4,398) | Challenge (1,042) |
| Improve (5,726) | Failed (2,421) | Top (14,899) | Motley (5,165) | Profitable (3,580) | Serious (1,022) |
| Restructuring (3,210) | Concerned (1,991) | Guidance (11,715) | Crude (5,109) | Efficiency (2,615) | Contrary (401) |

Words are shared by only two lexica and frequencies are given in parentheses.

[Back](#)



Comparison of Lexical Projections

| Variable | $\hat{\mu}$ | $\hat{\sigma}$ | Max | Q1 | Q2 | Q3 | Polarity |
|-------------------|-------------|----------------|-------|-------|-------|-------|----------|
| <i>Pos</i> (BL) | 0.033 | 0.012 | 0.134 | 0.025 | 0.032 | 0.040 | 88.04% |
| <i>Neg</i> (BL) | 0.015 | 0.010 | 0.091 | 0.008 | 0.014 | 0.020 | 10.51% |
| <i>Pos</i> (LM) | 0.014 | 0.007 | 0.074 | 0.009 | 0.013 | 0.018 | 55.70% |
| <i>Neg</i> (LM) | 0.012 | 0.009 | 0.085 | 0.006 | 0.011 | 0.016 | 40.17% |
| <i>Pos</i> (MPQA) | 0.038 | 0.012 | 0.134 | 0.031 | 0.038 | 0.045 | 96.26% |
| <i>Neg</i> (MPQA) | 0.013 | 0.008 | 0.133 | 0.007 | 0.012 | 0.017 | 2.87% |

Sample mean, sample standard deviation, maximum value, 1st, 2nd and 3rd quartiles, and polarity as relative dominance between positive and negative sentiment.

[Back](#)

Classification Evaluation

| Manual Label | BL Label | | | LM Label | | | MPQA Label | | | Total |
|--------------|----------|-----|-----|----------|-----|-----|------------|-----|-----|-------|
| | Pos | Neg | Neu | Pos | Neg | Neu | Pos | Neg | Neu | |
| Pos | 56 | 4 | 1 | 41 | 12 | 8 | 61 | 0 | 0 | 61 |
| Neg | 9 | 2 | 1 | 0 | 9 | 3 | 9 | 2 | 1 | 12 |
| Neu | 22 | 5 | 0 | 10 | 15 | 2 | 26 | 0 | 1 | 27 |
| Total | 87 | 11 | 2 | 51 | 36 | 13 | 96 | 2 | 2 | 100 |

[Back](#)

Tagging Example - BL

... McDonald's has an obesity **problem** that continues to get **worse**. And that's nothing to do with the food itself, but rather the huge menus that can now double as medieval fortification. For perspective, the chain's menu has grown 70% since 2007. And while more offerings might seem **like** a **good** thing, large menus result in **slower** service and more flare-ups between franchisees and the corporation.

Bloated menus raise inventory costs for smaller franchisees and **lead** to lower profit margins. The McDonald's corporate franchise fee is based upon sales instead of profits, making it a smaller **concern** for the company overall. ...

3 **positive words** and 5 **negative words**

[Article source](#)



Tagging Example - LM

... McDonald's has an obesity **problem** that continues to get **worse**. And that's nothing to do with the food itself, but rather the huge menus that can now double as medieval fortification. For perspective, the chain's menu has grown 70% since 2007. And while more offerings might seem like a **good** thing, large menus result in **slower** service and more flare-ups between franchisees and the corporation.

Bloated menus raise inventory costs for smaller franchisees and lead to lower profit margins. The McDonald's corporate franchise fee is based upon sales instead of profits, making it a smaller **concern** for the company overall. ...

1 **positive word** and 4 **negative words**



Tagging Example - MPQA

... McDonald's has an obesity **problem** that continues to get **worse**. And that's nothing to do with the food itself, but rather the huge menus that can now double as medieval fortification. For perspective, the chain's menu has grown 70% since 2007. And while more offerings might seem **like** a **good** thing, **large** menus result in **slower** service and more flare-ups between franchisees and the corporation.

Bloated menus raise inventory costs for smaller franchisees and **lead** to lower **profit** margins. The McDonald's corporate franchise fee is based upon sales instead of **profits**, making it a smaller **concern** for the company overall. ...

5 **positive words** and 5 **negative words**

Back



Correlation - Positive Sentiment

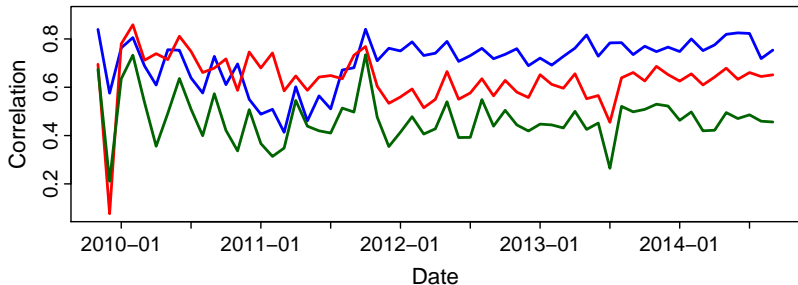


Figure 7: Monthly correlation between positive sentiment: **BL and LM** , **BL and MPQA**, **LM and MPQA**



Correlation - Negative Sentiment



Figure 8: Monthly correlation between negative sentiment: **BL and LM**, **BL and MPQA**, **LM and MPQA** [Back](#)



Garman and Klass range-based Measure of Volatility

$$\sigma_{i,t}^2 = 0.511(u - d)^2 - 0.019 \{c(u + d) - 2ud\} - 0.383c^2 \quad (7)$$

with $u = \log(P_{i,t}^H) - \log(P_{i,t}^C)$, $d = \log(P_{i,t}^L) - \log(P_{i,t}^O)$,

$$c = \log(P_{i,t}^C) - \log(P_{i,t}^O)$$

for company i on day t with $P_{i,t}^H$, $P_{i,t}^L$, $P_{i,t}^O$, $P_{i,t}^C$ as highest, lowest, opening and closing stock prices, respectively.

[Back](#)

Simulation Setup

- Evaluate the asymmetric reaction of volatility to sentiment
- $I_{i,t} \sim B(1, p_i)$
- Model dependence between *Pos*, *Neg* and different lexica with **firm specific** copula [Copula definition](#)
- Estimate **one** copula for $R_{M,t}$ and all firms $R_{i,t}$



Two-Step Approach

1. Marginals: ecdf
 2. Copulae: Gaussian
- Simulation: [Conditional inversion method](#)
 - $R_{M,t}$, $R_{i,t}$: use standardized residuals after fitting MA(1)-GARCH(1,1) process

[Back to Entire Panel](#)[Back to Attention Panel](#)

Multivariate Copula Definition

Definition

The **copula** is a multivariate distribution with all univariate margins being $U(0, 1)$.

Theorem (Sklar, 1959)

Let X_1, \dots, X_k be random variables with marginal distribution functions F_1, \dots, F_k and joint distribution function F . Then there exists a k -dimensional copula $C : [0, 1]^k \rightarrow [0, 1]$ such that

$\forall x_1, \dots, x_k \in \overline{\mathbb{R}} = [-\infty, \infty]$

$$F(x_1, \dots, x_k) = C\{F_1(x_1), \dots, F_k(x_k)\} \quad (8)$$

[Back to Simulation Setup](#)



Conditional Inversion Method

Frees and Valdez (1998):

$C = C(u_1, \dots, u_k)$, $C_i = C(u_1, \dots, u_i, 1, \dots, 1)$ and $C_k = C(u_1, \dots, u_k)$.

Conditional distribution of U_i :

$$\begin{aligned} C_i(u_i | u_1, \dots, u_{i-1}) &= P\{U_i \leq u_i | U_1 = u_1 \dots U_{i-1} = u_{i-1}\} \\ &= \frac{\partial^{i-1} C_i(u_1, \dots, u_i)}{\partial u_1 \dots \partial u_{i-1}} / \frac{\partial^{i-1} C_{i-1}(u_1, \dots, u_{i-1})}{\partial u_1 \dots \partial u_{i-1}} \end{aligned}$$



Conditional Inversion Method

Frees and Valdez (1998):

$C = C(u_1, \dots, u_k)$, $C_i = C(u_1, \dots, u_i, 1, \dots, 1)$ and $C_k = C(u_1, \dots, u_k)$.

Conditional distribution of U_i :

$$\begin{aligned} C_i(u_i | u_1, \dots, u_{i-1}) &= P\{U_i \leq u_i | U_1 = u_1 \dots U_{i-1} = u_{i-1}\} \\ &= \frac{\partial^{i-1} C_i(u_1, \dots, u_i)}{\partial u_1 \dots \partial u_{i-1}} / \frac{\partial^{i-1} C_{i-1}(u_1, \dots, u_{i-1})}{\partial u_1 \dots \partial u_{i-1}} \end{aligned}$$

- Generate i.r.v. $v_1, \dots, v_k \sim U(0, 1)$
- Set $u_1 = v_1$
- $u_i = C_k^{-1}(v_i | u_1, \dots, u_{i-1}) \forall i = 2, \dots, k$

[Back to Simulation Setup](#)



Bibliography I



Chen, Z., Daigler, R. T., and Parhizgari, A. M.
Persistence of volatility in futures markets
J. Futures Markets, 2006



Frees, E. W and Valdez, E. A.
Understanding relationships using copulas
N, Am. Actuar. J., 1998



Garman, M. and Klass, M.
On the Estimation of Security Price Volatilities from Historical Data
J. Bus., 1980



Bibliography II



Girard, E. and Biswas, R.
Trading volume and market volatility
Financ. Rev., 2007



Hu, M. and Liu, B.
Mining and Summarizing Customer Reviews
10th ACM SIGKDD, 2004



Loughran, T. and McDonald, B.
When is a liability not a liability?
J. Financ., 2011



Bibliography III



Shu, J. and Zhang, J. E.

Testing range estimators of historical volatility

J. Futures Markets, 2006



Wilson, T. and Wiebe, J. and Hoffmann, P.

Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis

HLT-EMNLP, 2005

