

# DETECTING DECEPTION

Lecture by: Myle Ott<sup>1</sup>

Incl. joint work with: Claire Cardie,<sup>1,2</sup> Yejin Choi,<sup>1</sup> Jeff Hancock<sup>2,3</sup>

Depts of C.S.,<sup>1</sup> I.S.,<sup>2</sup> Comm.<sup>3</sup>

Cornell University, Ithaca, New York

# Background

- Language use varies:
  - By **location**
    - soda vs. pop vs. coke
    - “koo” vs. “coo” (Eisenstein et al., 2010; 2011)
    - Also Johnstone (2010), Mei et al. (2006; 2007), Labov et al. (2006), Tagliamonte (2006), ...

# Background

- Language use varies:
  - By **genre**
    - British National Corpus: Koppel et al. (2002), Rayson et al. (2001), Biber et al. (1999), ...
    - Web: Mehler et al. (2010), Rehm et al. (2008), ...
    - Twitter: Westman and Freund (2010), ...

# Background

- Language use varies:
  - By the author's **gender**
    - British National Corpus: Koppel et al. (2002), ...
    - Blogs: Mukherjee and Liu (2010), ...
    - Twitter: Burger et al. (2011), ...
    - Cross-topic/domain: Sarawgi et al. (2011)

# Background

- Language use varies:
  - By the author's **beliefs, feelings, opinions**
    - Opinion mining and sentiment analysis:  
Pang and Lee (2008), ...
    - Belief annotation and tagging:  
Prabhakaran et al. (2010), Diab et al. (2009), ...
    - Detecting hedges: CoNLL 2010 Shared Task, ...

# Background

- Language use varies:
  - By whether the author is being **truthful** or **deceptive**
  - Studies have considered deception involving:
    - Emotional states: Ekman and Friesen (1969), ...
    - Views on social issues, e.g., death penalty: Newman et al. (2003), Mihalcea and Strapparava (2009), ...
    - Online dating profiles: Hancock et al. (2007), ...
    - Online product reviews: Ott et al. (2011; 2012), ...
    - ...

# Outline

- Briefly go over **a few important studies** and meta-analyses of deception:
  - Bond and DePaulo (2006)
  - Newman et al. (2003)
  - Vrij (2008)
- Case study on **detecting deceptive online reviews** of hotels: Ott et al. (2011)

# Bond and DePaulo (2006)

- Meta-analysis of over 200 studies of deception
- Finds that **human judges are relatively bad at detecting deception**, with an average accuracy of just 54%
- Poor performance due in part to **truth-bias**
  - Human judges are more likely to erroneously judge something as truthful than erroneous judge something as deceptive



# Newman et al. (2003)

- Hundreds of **true and false verbal and written samples** from undergraduates across three subjects: stance on abortion, feelings about friends, and a mock crime
- Language analyzed using the **Linguistic Inquiry and Word Count (LIWC)** software, developed by James Pennebaker (a co-author of the study)

# Newman et al. (2003)

- LIWC
  - Counts instances of ~4,500 **keywords**
    - Regular expressions, actually
  - Keywords are divided into 80 **psycholinguistically-motivated dimensions** across 4 broad groups
  - Reports means and standard deviations

# Newman et al. (2003)

- LIWC
  - Linguistic processes
    - e.g., average number of words per sentence
  - Psychological processes
    - e.g., talk, happy, know, feeling, eat
  - Personal concerns
    - e.g., job, cook, family
  - Spoken categories
    - e.g., yes, umm, blah

# Newman et al. (2003)

- LIWC
  - Linguistic processes
    - e.g., average number of words per sentence
  - Psychological processes
    - e.g., talk, happy, know, feeling, eat
  - Personal concerns
    - e.g., job, cook, family
  - Spoken categories
    - e.g., yes, umm, blah

# Newman et al. (2003)

- LIWC
  - Linguistic processes
    - e.g., average number of words per sentence
  - Psychological processes
    - e.g., talk, happy, know, feeling, eat
  - Personal concerns
    - e.g., job, cook, family
  - Spoken categories
    - e.g., yes, umm, blah

# Newman et al. (2003)

- LIWC
  - Linguistic processes
    - e.g., average number of words per sentence
  - Psychological processes
    - e.g., talk, happy, know, feeling, eat
  - Personal concerns
    - e.g., job, cook, family
  - Spoken categories
    - e.g., yes, umm, blah

# Newman et al. (2003)

- Results showed that deceptive samples have:
  - Reduced first-person singular (**psychological distancing**)
    - Liars avoid taking ownership of their lies, either to “dissociate” or due to a lack of personal experience
  - Increased **negative** emotion words
    - Possibly due to discomfort and guilt about lying
  - Reduced complexity and less exclusive language
    - Possibly due to increased **cognitive load**

# Vrij (2008)

- Comprehensive review of the current state of deception detection research
- In addition to the previous findings:
  - Meta-analysis of 30 studies shows that deceivers have difficulty encoding **spatial** and **temporal** information into their deceptions



# Outline

- Briefly go over **a few important studies** and meta-analyses of deception:
  - Bond and DePaulo (2006)
  - Newman et al. (2003)
  - Vrij (2008)
- Case study on **detecting deceptive online reviews** of hotels: Ott et al. (2011)

# Finding Deceptive Opinion Spam by Any Stretch of the Imagination


Myle Ott,<sup>1</sup> Yejin Choi,<sup>1</sup> Claire Cardie,<sup>1</sup> and Jeff Hancock<sup>2</sup>  
Dept. of Computer Science,<sup>1</sup> Communication<sup>2</sup>  
Cornell University, Ithaca, NY

# Motivation

- Consumers increasingly rate, review and research products online
- Potential for opinion spam
  - Disruptive opinion spam
  - Deceptive opinion spam

**Portland Marriott Downtown**  1

Hotel class   
1401 SW Naito Parkway, Portland, OR 97201

 **Reviews you can trust**

---

1-10 of 51 reviews << 1 2 ... 6 >>

Sort by [ Date ▼ ] English first ▾

---

  
nitropin...  
Auburn, WA  
9 reviews

**“A great riverfront getaway via Amtrak and free Streetcar!”**  


Date of review: Apr 22, 2011


As other reviewers have stated, yes the rooms are small but don't let that detour you from staying here. I'm still giving this hotel 5 stars based on the quality and level of service we received from everybody here. We payed a little extra online for the breakfast package and it was well worth it. The breakfast was a full...  
[more](#) ▼

# Motivation

- Consumers increasingly rate, review and research products online
- **Potential for opinion spam**
  - Disruptive opinion spam
  - Deceptive opinion spam

**Portland Marriott Downtown**  1

Hotel class   
1401 SW Naito Parkway, Portland, OR 97201

 **Reviews you can trust**

---

1-10 of 51 reviews << 1 2 ... 6 >>

Sort by [ Date ▼ ] English first ▾

---

  
nitropin...  
Auburn, WA  
9 reviews

**“A great riverfront getaway via Amtrak and free Streetcar!”**  
  
Date of review: Apr 22, 2011

As other reviewers have stated, yes the rooms are small but don't let that detour you from staying here. I'm still giving this hotel 5 stars based on the quality and level of service we received from everybody here. We payed a little extra online for the breakfast package and it was well worth it. The breakfast was a full...  
[more](#) ▼

# Motivation

- Consumers increasingly rate, review and research products online
- Potential for opinion spam
  - Disruptive opinion spam
  - Deceptive opinion spam

★★★★★ **Great Customer Service!!**, April 7, 2011

By [akaempf](#)  - [See all my reviews](#)

**Amazon Verified Purchase** (What's this?)

**This review is from: [Apple iPad 2 MC984LL/A Tablet \(64GB, Wifi + AT&T 3G, White\) NEWEST MODEL \(Personal Computers\)](#)**

"WE SHIP TECH" is a great reliable company. I ordered the iPad2 late 3/30 @ 10:50pm and received the iPad2 4/1. When I wrote an email to them on the 3/31 they responded in about 20 min max. It's so hard to find great customer service and not get scammed these days that "We Ship Tech" is a breath of fresh air!! I would surely use them again and highly recommend them to anyone who expects great products & service. Thank you We Ship Tech!!!!

# Motivation

- Consumers increasingly rate, review and research products online
- Potential for opinion spam
  - Disruptive opinion spam
  - **Deceptive opinion spam**

★★★★★ **Works Just as expected**, May 14, 2007

By [Laurie B. Cook](#)  - [See all my reviews](#)

REAL NAME

**This review is from:** Belkin F5U301 CableFree 4-Port USB 2.0 Hub with Dongle (Electronics)

Supplies good range and does provide true wireless USB. Software worked right out of the box. I have been recommending this nifty little device to all my friends. Very useful device.

# Motivation

Which of these two hotel reviews is deceptive opinion spam?

Date of review: Jun 9, 2006

4 people found this review helpful

I have stayed at many hotels traveling for both business and pleasure and I can honestly say that The James is tops. The service at the hotel is first class. The rooms are modern and very comfortable. The location is perfect within walking distance to all of the great sights and restaurants. Highly recommend to both business travellers and couples.

Date of review: Jun 9, 2006

4 people found this review helpful

My husband and I stayed at the James Chicago Hotel for our anniversary. This place is fantastic! We knew as soon as we arrived we made the right choice! The rooms are BEAUTIFUL and the staff very attentive and wonderful!! The area of the hotel is great, since I love to shop I couldn't ask for more!! We will definitely be back to Chicago and we will for sure be back to the James Chicago.

# Motivation

Which of these two hotel reviews is deceptive opinion spam?

Answer:

Date of review: Jun 9, 2006

4 people found this review helpful

My husband and I stayed at the James Chicago Hotel for our anniversary. This place is fantastic! We knew as soon as we arrived we made the right choice! The rooms are BEAUTIFUL and the staff very attentive and wonderful!! The area of the hotel is great, since I love to shop I couldn't ask for more!! We will definatly be back to Chicago and we will for sure be back to the James Chicago.



# Overview

- Motivation
- **Gathering Data**
- Human Performance
- Classifier Performance
- Conclusion

# Gathering Data

- Label existing reviews
  - Can't manually do this
  - Duplicate detection (Jindal and Liu, 2008)
- Create new reviews
  - Mechanical Turk

# Gathering Data

- Label existing reviews
  - Can't manually do this
  - Duplicate detection (Jindal and Liu, 2008)
- Create new reviews
  - Mechanical Turk

# Gathering Data

- Label existing reviews
  - Can't manually do this
  - Duplicate detection (Jindal and Liu, 2008)
- Create new reviews
  - Mechanical Turk

# Gathering Data

- Label existing reviews
  - Can't manually do this
  - Duplicate detection (Jindal and Liu, 2008)
- Create new reviews
  - Mechanical Turk

# Gathering Data

- Label existing reviews
  - Can't manually do this
  - Duplicate detection (Jindal and Liu, 2008)
- Create new reviews
  - Mechanical Turk

# Gathering Data

- Mechanical Turk
  - 20 hotels
  - 20 reviews / hotel
  - Offer \$1 / review
  - 400 reviews

# Gathering Data

- Mechanical Turk
  - 20 hotels
  - 20 reviews / hotel
  - Offer \$1 / review
  - 400 reviews

Home → United States → Illinois (IL) → Chicago → Chicago Hotels → James Chicago

## James Chicago

Hotel class ★★★★★  
55 East Ontario, Corner of Rush and Ontario, Chicago, IL 60611  
877.526.3755 [Hotel website](#) [E-mail hotel](#)

### What travelers say about James Chicago

- Great location (33)
- Room service (20)
- Very nice (18)
- Trader joe (16)
- Boutique hotel (15)
- Magnificent mile (14)
- Very good (13)
- Michigan avenue (13)
- Comfortable bed (10)
- Friendly and helpful (8)

### Reviews you can trust

Filter traveler reviews [Write a Review](#)

Trip type	Traveler rating
<input checked="" type="radio"/> All reviews (449)	<input checked="" type="radio"/> All (449)
<input type="radio"/> Business reviews (94)	<input type="radio"/> Excellent (278) 
<input type="radio"/> Couples reviews (194)	<input type="radio"/> Very good (116) 
<input type="radio"/> Family reviews (28)	<input type="radio"/> Average (23) 
<input type="radio"/> Friends reviews (60)	<input type="radio"/> Poor (19) 
<input type="radio"/> Solo travel reviews (62)	<input type="radio"/> Terrible (13) 



# Gathering Data

- Mechanical Turk
  - 20 hotels
  - 20 reviews / hotel
  - Offer \$1 / review
  - 400 reviews

1-10 of 449 reviews « 1 2 ... 45 »

Sort by [ Date ▼ ] [ Rating ] English first

 **“Amazing Hotel”**  
○○○○○

Date of review: Apr 25, 2011 - **New**

[emmabake...](#)   
Farnborough, UK  
2 contributions

Stayed at this hotel in May 2010. Came on business from the UK with my husband for the Snack and Candy Expo at McCormick Place and decided that this place was an easy taxi ride away but within walking distance for our spare time. Wow, the hotel was amazing, one of the best we've stayed in. Our room wasn't ready...

[more](#) ▼

# Gathering Data

- Mechanical Turk
  - 20 hotels
  - 20 reviews / hotel
  - Offer \$1 / review
  - 400 reviews



# Gathering Data

- Mechanical Turk
  - 20 hotels
  - 20 reviews / hotel
  - Offer \$1 / review
  - 400 reviews
- Average time spent:
  - > 8 minutes
- Average length:
  - > 115 words

# Gathering Data

- 400 truthful reviews
  - TripAdvisor.com
  - Lengths distributed similarly to deceptive reviews

# Overview

- Motivation
- Gathering Data
- **Human Performance**
- Classifier Performance
- Conclusion

# Human Performance

- Why bother?
  - Validates deceptive opinions
  - Baseline to compare other approaches

# Human Performance

- Why bother?
  - Validates deceptive opinions
  - Baseline to compare other approaches

# Human Performance

- Why bother?
  - Validates deceptive opinions
  - **Baseline to compare other approaches**



# Human Performance

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

			TRUTHFUL			DECEPTIVE		
		Accuracy	P	R	F	P	R	F
HUMAN	JUDGE 1	<b>61.9%</b>	57.9	87.5	<b>69.7</b>	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	<b>95.0</b>	68.8	<b>78.9</b>	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

Performed at chance  
(p-value = 0.1)

			TRUTHFUL			DECEPTIVE		
		Accuracy	P	R	F	P	R	F
HUMAN	JUDGE 1	61.9%	57.9	87.5	69.7	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	95.0	68.8	78.9	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6

Performed at chance  
(p-value = 0.5)

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

			TRUTHFUL			DECEPTIVE		
		Accuracy	P	R	F	P	R	F
HUMAN	JUDGE 1	<b>61.9%</b>	57.9	87.5	<b>69.7</b>	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	<b>95.0</b>	68.8	<b>78.9</b>	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

			TRUTHFUL			DECEPTIVE		
		Accuracy	P	R	F	P	R	F
HUMAN	JUDGE 1	<b>61.9%</b>	57.9	87.5	<b>69.7</b>	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	<b>95.0</b>	68.8	<b>78.9</b>	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6

Classified fewer than 12% of opinions as deceptive!

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

			TRUTHFUL			DECEPTIVE		
		Accuracy	P	R	F	P	R	F
HUMAN	JUDGE 1	<b>61.9%</b>	57.9	87.5	<b>69.7</b>	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	<b>95.0</b>	68.8	<b>78.9</b>	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

			TRUTHFUL			DECEPTIVE		
		Accuracy	P	R	F	P	R	F
HUMAN	JUDGE 1	<b>61.9%</b>	57.9	87.5	<b>69.7</b>	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	<b>95.0</b>	68.8	<b>78.9</b>	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6
META	MAJORITY	58.1%	54.8	92.5	68.8	76.0	23.8	36.2
	SKEPTIC	60.6%	<b>60.8</b>	60.0	60.4	60.5	<b>61.3</b>	<b>60.9</b>

- 80 truthful and 80 deceptive reviews
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

# Human Performance

		Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
HUMAN	JUDGE 1	<b>61.9%</b>	57.9	87.5	<b>69.7</b>	74.4	36.3	48.7
	JUDGE 2	56.9%	53.9	<b>95.0</b>	68.8	<b>78.9</b>	18.8	30.3
	JUDGE 3	53.1%	52.3	70.0	59.9	54.7	36.3	43.6
META	MAJORITY	58.1%	54.8	92.5	68.8	76.0	23.8	36.2
	SKEPTIC	60.6%	<b>60.8</b>	60.0	60.4	60.5	<b>61.3</b>	<b>60.9</b>

- 80 truthful and 80 deceptive
- 3 undergraduate judges
  - Truth bias
- 2 meta-judges

No more truth bias!



# Overview

- Motivation
- Gathering Data
- Human Performance
- **Classifier Performance**
- Conclusion

# Classifier Performance

- Three feature sets
  - Genre identification
  - Psycholinguistic deception detection
  - Text categorization
- Linear SVM

# Classifier Performance

- Three feature sets
  - Genre identification
  - Psycholinguistic deception detection
  - Text categorization
- Linear SVM

# Classifier Performance

- **Genre identification**
  - 48 part-of-speech (PoS) features
  - Baseline automated approach
- **Expectations**
  - Truth similar to informative writing
  - Deception similar to imaginative writing

# Classifier Performance

- Genre identification
  - 48 part-of-speech (PoS) features
  - Baseline automated approach
- Expectations
  - Truth similar to informative writing
  - Deception similar to imaginative writing

# Classifier Performance

- Genre identification
  - 48 part-of-speech (PoS) features
  - Baseline automated approach
- Expectations
  - Truth similar to informative writing
  - Deception similar to imaginative writing

# Classifier Performance

- Genre identification
  - 48 part-of-speech (PoS) features
  - Baseline automated approach
- **Expectations**
  - Truth similar to informative writing
  - Deception similar to imaginative writing

# Classifier Performance

Approach	Features	Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
GENRE IDENTIFICATION	POS	73.0%	75.3	68.5	71.7	71.1	77.5	74.2

Finding Deceptive Opinion Spam by Any Stretch  
of the Imagination



# Classifier Performance

Approach	Features	Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
GENRE IDENTIFICATION	POS	73.0%	75.3	68.5	71.7	71.1	77.5	74.2

Outperforms human judges!  
(p-values = {0.06, 0.01, 0.001})

# Classifier Performance

TRUTHFUL/INFORMATIVE			DECEPTIVE/IMAGINATIVE			
Category	Variant	Weight	Category	Variant	Weight	
NOUNS	Singular	0.008	VERBS	Base	-0.057	
	Plural	0.002		Past tense	<b>0.041</b>	
	Proper, singular	<b>-0.041</b>		Present participle	-0.089	
	Proper, plural	0.091		Singular, present	-0.031	
ADJECTIVES	General	0.002		Third person singular, present	<b>0.026</b>	
	Comparative	0.058		Modal	-0.063	
	Superlative	<b>-0.164</b>		ADVERBS	General	<b>0.001</b>
PREPOSITIONS	General	0.064			Comparative	-0.035
DETERMINERS	General	0.009			PRONOUNS	Personal
COORD. CONJ.	General	0.094		Possessive		-0.303
VERBS	Past participle	0.053	PRE-DETERMINERS	General	<b>0.017</b>	
ADVERBS	Superlative	<b>-0.094</b>				

- Rayson et. al. (2001)
  - Informative on left, imaginative on right

# Classifier Performance

TRUTHFUL/INFORMATIVE			DECEPTIVE/IMAGINATIVE			
Category	Variant	Weight	Category	Variant	Weight	
NOUNS	Singular	0.008	VERBS	Base	-0.057	
	Plural	0.002		Past tense	<b>0.041</b>	
	Proper, singular	<b>-0.041</b>		Present participle	-0.089	
	Proper, plural	0.091		Singular, present	-0.031	
ADJECTIVES	General	0.002		Third person singular, present	<b>0.026</b>	
	Comparative	0.058		Modal	-0.063	
	Superlative ★	<b>-0.164</b>		ADVERBS	General	<b>0.001</b>
PREPOSITIONS	General	0.064		ADVERBS	Comparative	-0.035
DETERMINERS	General	0.009			PRONOUNS	Personal
COORD. CONJ.	General	0.094		PRONOUNS	Possessive	-0.303
VERBS	Past participle	0.053	PRE-DETERMINERS		General	<b>0.017</b>
ADVERBS	Superlative ★	<b>-0.094</b>				

e.g., best, finest

e.g., most

- Rayson et. al. (2001)
  - Informative on left, imaginative on right

# Classifier Performance

- Linguistic Inquire and Word Count (Pennebaker et al., 2001; 2007)
  - Counts instances of ~4,500 keywords
    - Regular expressions, actually
  - Keywords are divided into 80 dimensions across 4 broad groups

# Classifier Performance

Approach	Features	Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
GENRE IDENTIFICATION	POS	73.0%	75.3	68.5	71.7	71.1	77.5	74.2
PSYCHOLINGUISTIC DECEPTION DETECTION	LIWC	76.8%	77.2	76.0	76.6	76.4	77.5	76.9

Finding Deceptive Opinion Spam by Any Stretch  
of the Imagination

# Classifier Performance

Approach	Features	Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
GENRE IDENTIFICATION	POS	73.0%	75.3	68.5	71.7	71.1	77.5	74.2
PSYCHOLINGUISTIC DECEPTION DETECTION	LIWC	76.8%	77.2	76.0	76.6	76.4	77.5	76.9

Outperforms PoS!  
(p-value = 0.02)

# Classifier Performance

- Text categorization (n-grams)
  - Unigrams
  - Bigrams<sup>+</sup>
    - Includes unigrams
  - Trigrams<sup>+</sup>
    - Includes unigrams and bigrams

# Classifier Performance

Approach	Features	Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
GENRE IDENTIFICATION	POS	73.0%	75.3	68.5	71.7	71.1	77.5	74.2
PSYCHOLINGUISTIC DECEPTION DETECTION	LIWC	76.8%	77.2	76.0	76.6	76.4	77.5	76.9
TEXT CATEGORIZATION	UNIGRAMS	88.4%	89.9	86.5	88.2	87.0	<b>90.3</b>	88.6
	BIGRAMS	89.6%	<b>90.1</b>	89.0	89.6	89.1	<b>90.3</b>	89.7
	LIWC+BIGRAMS	<b>89.8%</b>	89.8	<b>89.8</b>	<b>89.8</b>	<b>89.8</b>	89.8	<b>89.8</b>
	TRIGRAMS	89.0%	89.0	89.0	89.0	89.0	89.0	89.0

Finding Deceptive Opinion Spam by Any Stretch  
of the Imagination



# Classifier Performance

Approach	Features	Accuracy	TRUTHFUL			DECEPTIVE		
			P	R	F	P	R	F
GENRE IDENTIFICATION	POS	73.0%	75.3	68.5	71.7	71.1	77.5	74.2
PSYCHOLINGUISTIC DECEPTION DETECTION	LIWC	76.8%	77.2	76.0	76.6	76.4	77.5	76.9
TEXT CATEGORIZATION	UNIGRAMS	88.4%	89.9	86.5	88.2	87.0	<b>90.3</b>	88.6
	BIGRAMS	89.6%	<b>90.1</b>	89.0	89.6	89.1	<b>90.3</b>	89.7
	LIWC+BIGRAMS	<b>89.8%</b>	89.8	<b>89.8</b>	<b>89.8</b>	<b>89.8</b>	89.8	<b>89.8</b>
	TRIGRAMS	89.0%	89.0	89.0	89.0	89.0	89.0	89.0

Outperforms all  
other methods!

# Classifier Performance

LIWC+BIGRAMS	
TRUTHFUL	DECEPTIVE
-	chicago
...	my
on	hotel
location	,_and
)	luxury
allpunct <sub>LIWC</sub>	experience
floor	hilton
(	business
the_hotel	vacation
bathroom	i
small	spa
helpful	looking
\$	while
hotel_.	husband
other	my_husband

- Spatial difficulties (Vrij et al., 2009)
- Psychological distancing (Newman et al., 2003)

# Classifier Performance

LIWC+BIGRAMS	
TRUTHFUL	DECEPTIVE
-	chicago
...	my
★ on	hotel
★ location	,_and
)	luxury
allpunct <sub>LIWC</sub>	experience
★ floor	hilton
(	business
the_hotel	vacation
★ bathroom	i
★ small	spa
helpful	looking
\$	while
hotel_.	husband
other	my_husband

- Spatial difficulties (Vrij et al., 2009)
- Psychological distancing (Newman et al., 2003)

# Classifier Performance

LIWC+BIGRAMS	
TRUTHFUL	DECEPTIVE
-	chicago
...	my
on	hotel
location	,_and
)	luxury
allpunct <sub>LIWC</sub>	experience
floor	hilton
(	★ business
the_hotel	★ vacation
bathroom	i
small	spa
helpful	looking
\$	while
hotel_.	★ husband
other	★ my_husband

- Spatial difficulties (Vrij et al., 2009)
- Psychological distancing (Newman et al., 2003)

# Classifier Performance

LIWC+BIGRAMS	
TRUTHFUL	DECEPTIVE
-	chicago
...	my
on	hotel
location	,_and
)	luxury
allpunct <sub>LIWC</sub>	experience
floor	hilton
(	business
the_hotel	vacation
bathroom	i
small	spa
helpful	looking
\$	while
hotel_.	husband
other	my_husband

- Spatial difficulties (Vrij et al., 2009)
- Psychological distancing (Newman et al., 2003)

# Classifier Performance

LIWC+BIGRAMS	
TRUTHFUL	DECEPTIVE
-	chicago
...	★ my
on	hotel
location	,_and
)	luxury
allpunct <sub>LIWC</sub>	experience
floor	hilton
(	business
the_hotel	vacation
bathroom	★ i
small	spa
helpful	looking
\$	while
hotel_.	husband
other	my_husband

- Spatial difficulties (Vrij et al., 2009)
- Psychological distancing (Newman et al., 2003)

# Overview

- Motivation
- Gathering Data
- Human Performance
- Classifier Performance
- Conclusion

# Conclusion

- Language use varies depending on features of the text and the author
- It seems likely that whether the author is being truthful or deceptive influences their language use
- Research into detecting deception has interesting real-life applications, e.g., detecting fake reviews
- Standard n-gram text categorization can outperform human performance on this task



- Jacob Eisenstein, Brendan O'Connor, Noah A. Smith, and Eric P. Xing. 2010. A latent variable model for geographic lexical variation. In Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing (EMNLP '10). Association for Computational Linguistics, Stroudsburg, PA, USA, 1277-1287.
- Jacob Eisenstein, Noah A. Smith, and Eric P. Xing. 2011. Discovering sociolinguistic associations with structured sparsity. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1 (HLT '11), Vol. 1. Association for Computational Linguistics, Stroudsburg, PA, USA, 1365-1374.
- B. Johnstone. 2010. Language and place. In R. Mesthrie and W. Wolfram, editors, Cambridge Handbook of Sociolinguistics. Cambridge University Press.
- Qiaozhu Mei, Chao Liu, Hang Su, and ChengXiang Zhai. 2006. A probabilistic approach to spatiotemporal theme pattern mining on weblogs. In Proceedings of the 15th international conference on World Wide Web (WWW '06). ACM, New York, NY, USA, 533-542.
- Qiaozhu Mei, Xu Ling, Matthew Wondra, Hang Su, and ChengXiang Zhai. 2007. Topic sentiment mixture: modeling facets and opinions in weblogs. In Proceedings of the 16th international conference on World Wide Web (WWW '07). ACM, New York, NY, USA, 171-180.
- Labov, W., Ash, S. & Boberg, C. (2006). The atlas of North American English: phonetics, phonology, and sound change: a multimedia reference tool. Mouton de Gruyter
- Tagliamonte, S. (2006). Analysing sociolinguistic variation. Cambridge Univ Press.
- Koppel, M., Argamon, S., Shimoni, A. R.. 2002. Automatically Categorizing Written Text by Author Gender. Literary and Linguistic Computing.
- P. Rayson, A. Wilson, and G. Leech. 2001. Grammatical word class variation within the British National Corpus sampler. Language and Computers, 36(1):295-306.
- D. Biber, S. Johansson, G. Leech, S. Conrad, E. Finegan, and R. Quirk. 1999. Longman grammar of spoken and written English, volume 2. MIT Press.
- Mehler, S. Sharoff and M. Santini. 2010. Genres on the Web: Computational Models and Empirical Studies. TEXT, SPEECH AND LANGUAGE TECHNOLOGY
- Rehm, Georg; Santini, Marina; Mehler, Alexander; Braslavski, Pavel; Gleim, Rüdiger; Stubbe, Andrea; Symonenko, Svetlana; Tavano, Mirko and Vidulin, Vedrana (2008): "Towards a Reference Corpus of Web Genres for the Evaluation of Genre Identification Systems". In: Proceedings of the 6th Language Resources and Evaluation Conference (LREC 2008). Marrakech, Morocco.
- S. Westman and L. Freund. Information interaction in 140 characters or less: genres on twitter. In IliX '10, pages 323-328, 2010.
- Arjun Mukherjee and Bing Liu. 2010. Improving gender classification of blog authors. In Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, Cambridge, MA, October. Association for Computational Linguistics.
- John D. Burger, John Henderson, George Kim, and Guido Zarrella. 2011. Discriminating gender on Twitter. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP '11). Association for Computational Linguistics, Stroudsburg, PA, USA, 1301-1309.

- Ruchita Sarawgi, Kailash Gajulapalli, and Yejin Choi. 2011. Gender attribution: tracing stylometric evidence beyond topic and genre. In Proceedings of the Fifteenth Conference on Computational Natural Language Learning, CoNLL '11, pages 78–86, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Pang, B. & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2, 1-135.
- Vinodkumar Prabhakaran, Owen Rambow, and Mona Diab. 2010. Automatic committed belief tagging. In Proceedings of the 23rd International Conference on Computational Linguistics: Posters (COLING '10). Association for Computational Linguistics, Stroudsburg, PA, USA, 1014-1022.
- Mona T. Diab, Lori Levin, Teruko Mitamura, Owen Rambow, Vinodkumar Prabhakaran, and Weiwei Guo. 2009. Committed belief annotation and tagging. In Proceedings of the Third Linguistic Annotation Workshop (ACL-IJCNLP '09). Association for Computational Linguistics, Stroudsburg, PA, USA, 68-73.
- Ekman, P., & Friesen, W. V. (1969). The repertoire of nonverbal behavior: Categories, origins, usage, and coding. *Semiotica*, 1, 49-98.
- M.L. Newman, J.W. Pennebaker, D.S. Berry, and J.M. Richards. 2003. Lying words: Predicting deception from linguistic styles. *Personality and Social Psychology Bulletin*, 29(5):665.
- R. Mihalcea and C. Strapparava. 2009. The lie detector: Explorations in the automatic recognition of deceptive language. In Proceedings of the ACL-IJCNLP 2009 Conference Short Papers, pages 309–312. Association for Computational Linguistics
- Jeffrey T. Hancock, Catalina Toma, and Nicole Ellison. 2007. The truth about lying in online dating profiles. In Proceedings of the SIGCHI conference on Human factors in computing systems (CHI '07). ACM, New York, NY, USA, 449-452. DOI=10.1145/1240624.1240697 <http://doi.acm.org/10.1145/1240624.1240697>
- Myle Ott, Yejin Choi, Claire Cardie, and Jeffrey T. Hancock. 2011. Finding deceptive opinion spam by any stretch of the imagination. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1 (HLT '11), Vol. 1. Association for Computational Linguistics, Stroudsburg, PA, USA, 309-319.
- Myle Ott, Claire Cardie, and Jeff Hancock. 2012. Estimating the prevalence of deception in online review communities. In Proceedings of the 21st international conference on World Wide Web (WWW '12). ACM, New York, NY, USA, 201-210. DOI=10.1145/2187836.2187864 <http://doi.acm.org/10.1145/2187836.2187864>
- C.F. Bond and B.M. DePaulo. 2006. Accuracy of deception judgments. *Personality and Social Psychology Review*, 10(3):214.
- A. Vrij. 2008. *Detecting lies and deceit: Pitfalls and opportunities*. Wiley-Interscience.
- N. Jindal and B. Liu. 2008. Opinion spam and analysis. In Proceedings of the international conference on Web search and web data mining, pages 219–230. ACM.
- A. Vrij, S. Leal, P.A. Granhag, S. Mann, R.P. Fisher, J. Hillman, and K. Sperry. 2009. Outsmarting the liars: The benefit of asking unanticipated questions. *Law and human behavior*, 33(2):159–166.