

CSCE 590-1: Trusted AI

Lecture 23: AI - Unstructured Text – Mitigation and Debiasing

PROF. BIPLAV SRIVASTAVA, AI INSTITUTE

9TH NOV, 2021

Carolinian Creed: “I will practice personal and academic integrity.”

Organization of Lecture 23

- Introduction Segment
 - Recap from last lecture
 - Paper selection by graduate students
- Main Segment
 - Transparency through documentation - rating AI services
 - Debiasing methods
- Concluding Segment
 - About next lecture – Lecture 24
 - Ask me anything

Introductory Segment

Schedule Snapshot

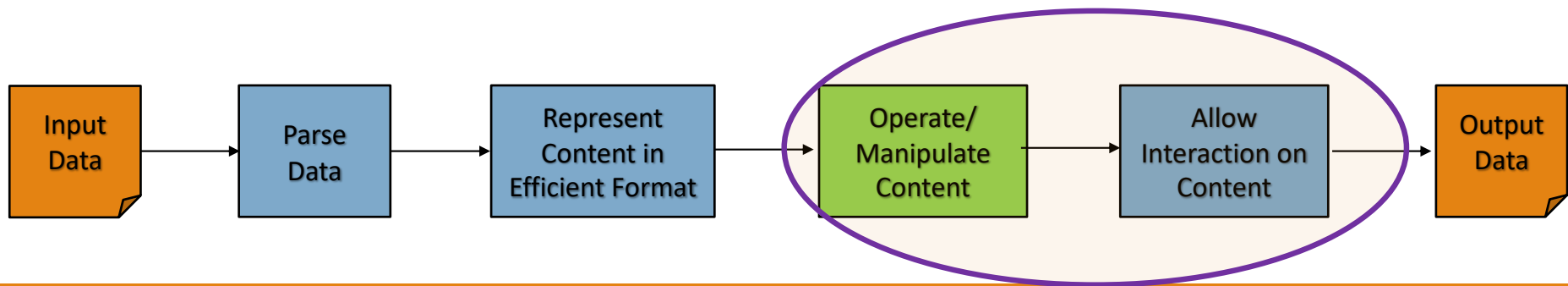


Oct 26 (Tu)	Review: Explanation Methods, AIX 360, Discussion	Quiz 3
Oct 28 (Th)	Review: project presentations, Discussion	
Nov 2 (Tu)	AI - Unstructured (Text): Analysis – Supervised ML – Trust Issues	
Nov 4 (Th)	AI - Unstructured (Text): Analysis – Supervised ML – Mitigation Methods	
Nov 9 (Tu)	AI - Unstructured (Text): Analysis – Rating and Debiasing Methods	
Nov 11 (Th)	AI – Unstructured Text - Explanation Methods Trust: AI Testing	
Nov 16 (Tu)	Trust: Human-AI Collaboration	Quiz 4
Nov 18 (Th)	Paper presentations – Graduate students	Final assignment for Graduate students
Nov 23 (Tu)	Emerging Standards and Laws	

Recap of Lecture 22

- We looked at some mitigation methods for trust with AI services
- Awareness
 - about data and its lack of representativeness
 - methods and limitations like language models
 - tools and their implementation choices
- Use diverse data; consider synthetic data
- Promote transparency through documentation
 - Analogy of its effectiveness in food industry
 - The idea of rating AI services based on behavior

Main Segment



Building Trust in AI Systems: Transparency Through Documentation

Details: <https://sites.google.com/site/biplavsrivastava/research-1/trustedai>

Transparency Through Documentation of Rating

Documentation about

- Outcome (e.g., Nutrition label, Electronic DataSheet, Factsheet)
- Process (e.g., SEI Capability Maturity Model, ISO 9001)

Documentation by

- Producer (e.g., Nutrition label)
- Consumer (e.g., Yelp rating)
- Independent 3rd Party (e.g., JD Powers, NHTSA car crash)

Reference: AboutML Project at PAI - <https://www.partnershiponai.org/about-ml-get-involved/#read>

Problem We Are Tackling for AI

Insight

- Empower people to make informed decisions regarding which AI to choose
- Communicate trust information better!
 - Analogy: Food labels
- Facilitate users in understanding their choices

Calories 230		Amount per 2/3 cup	
		Calories 230	
		% DV*	
Total Fat	8g	12%	Total Fat 8g
Saturated Fat	1g	5%	Saturated Fat 1g
Trans Fat	0g		Trans Fat 0g
Cholesterol	0mg	0%	Cholesterol 0mg
Sodium	160mg	7%	Sodium 160mg
Total Carbohydrate	37g	12%	Total Carbs 37g
Dietary Fiber	4g	16%	Dietary Fiber 4g
Sugars	1g		Sugars 1g
Protein	3g		Added Sugars 0g
		Protein 3g	
Vitamin A	10%	10%	Vitamin D 2mcg
Vitamin C	8%	20%	Calcium 260mg
Calcium	20%		
Iron	45%		
*Percent Daily Values are based on a diet of other people's secrets.			

In a series of previous work, we have developed ideas for rating bias of AI services

- For transactional services, method relies on a novel 2-stage testing method for bias. Papers in AIES 2018, IBM Sys Jour 2019, AAAI 2021 (Demo), IEEE Internet Computing (2021)
- For conversation services (chatbot), method relies on testing properties (called issues) such as fairness, lack of information leakage, lack of abusive language, and adequate conversation complexity. Paper in IEEE Transactions on Technology and Society 2020.

But ideas are general and can apply to audio-, image- and multimodal AI services, as well as their compositions.

Illustration of Translation Errors and Their Bias Perception About Gender

Type: 1

File: Result - true - i - es.json

Comment: Gender changed to new value (It).

"**original**": "He is a Architect. She is a Astronomer.", "**translated**": "Es un arquitecto. Ella es una astrónoma.",

"**oto**": "It is an architect. She is an astronomer.",

Type: 2

File: Result - false- g - tr.json;

Comment: Gender Flipped from She to He.

"**original**": "She is a Architect. She is a Astronomer.",

"**translated**": "O bir Mimar. O bir gökbilimcidir.",

"**oto**": "He is an architect. He is an astronomer.",

Type: 5

File: * - i - fr.json

Comment: No gender changed; Sentence modified.

"**original**": "He is a Accountant. He is a Actor /Actress.",

"**translated**": "C'est un comptable. C'est un acteur et un acteur.",

"**oto**": "He is an accountant. He is an actor and an actor.",

Type: 3

File: Result - false- g - ar.json;

Comment: Grammatically wrong sentence; Subject missing.

"**original**": "She is a Factory worker. He is a Farmer.",

"**translated**": "هي عامل مصنع. هو مزارع.",

"**oto**": "~~Is~~ a factory worker. He is a farmer.",

Type: 4

File: Result - false- g - tr.json;

Comment: "Multiple. Gender changed and flipped."

"**original**": "He is a Nurse. He is a Optician.",

"**translated**": "O bir hemşire. O bir Optisyendir.",

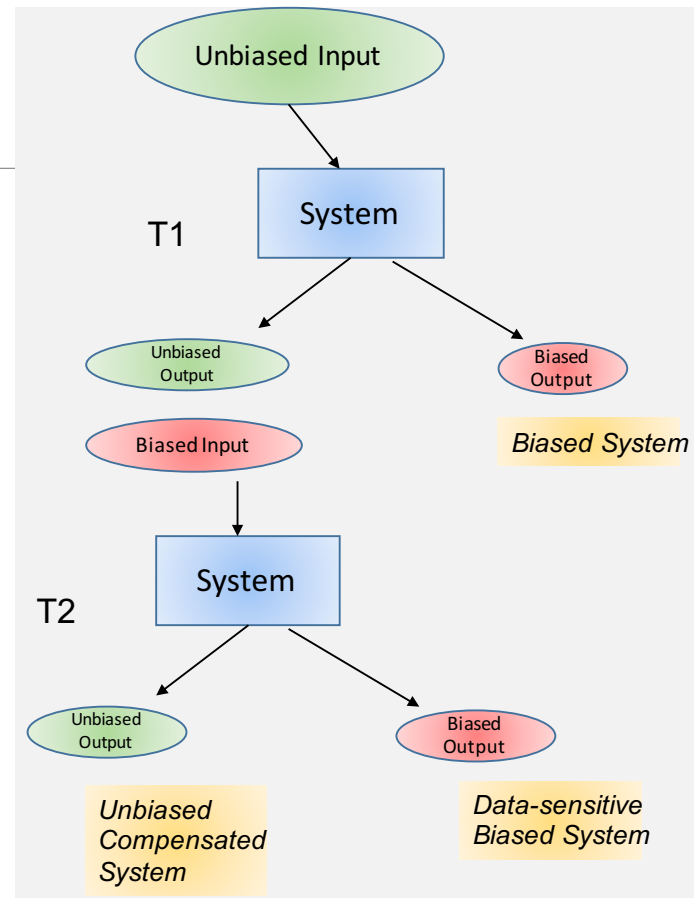
"**oto**": "She is a nurse. It is an Optic.",

**1, 2, 3 and 4 have gender issues;
3 and 5 have translation mistakes**

Key Idea - Translators

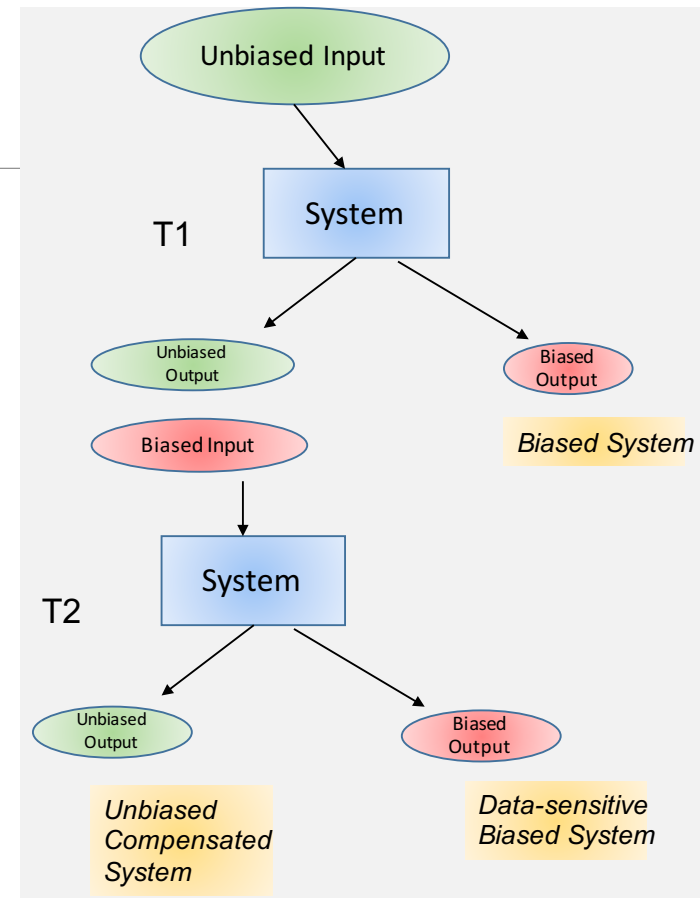
Specification of gender distribution (e.g., He, She)

- One or more termed biased, e.g., (0.1, 0.9)
- One or more termed unbiased, e.g., (0.5, 0.5)

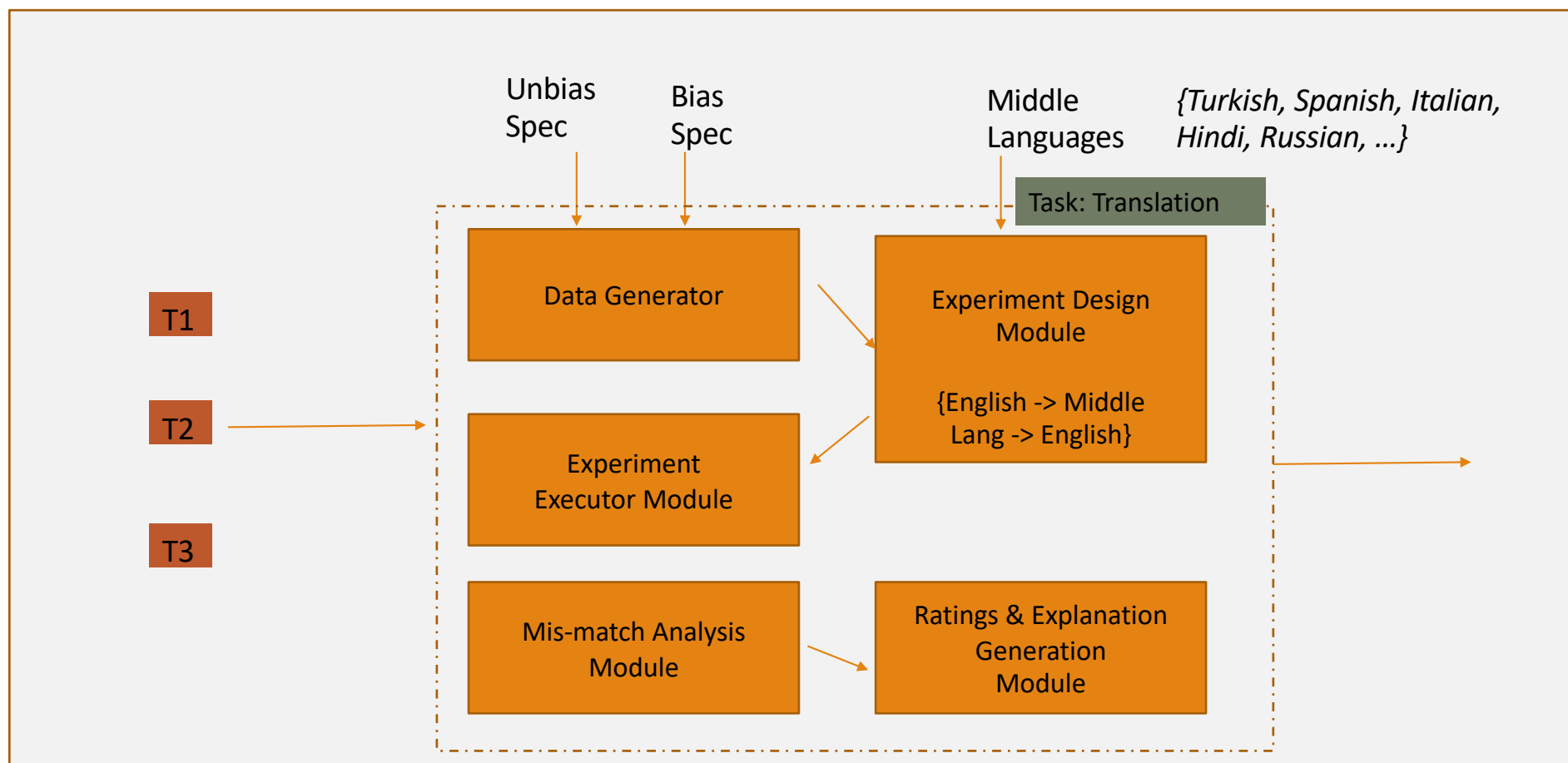


Rating Translators

- 3rd party rating service: independent of API producer or consumer.
- Gives API producer distributions of biased and unbiased data.
- Does a new 2-step testing and produces ratings of 3 main levels: -
 - Unbiased Compensated System (UCS): Forces an assumed distribution among legal choices
 - Data-sensitive Biased System (DSBS): Its output follows a distribution similar to input
 - Biased System (BS): Follows a distribution statistically different from assumption
- Ratings supports multiple distribution definitions under unbiased and biased categories.
- Enhance scheme for compositions of APIs with their 3-level ratings
- Implementation and experiments on off-the-shelf translators and translation task with many middle languages.



Illustrative Setup and Experiments



Rating Translators – Some Results

Intermediate Step

Spanish	es	He	She	OTHER	es	He	She	OTHER	es	He	She	OTHER
	u1		0.2	0.05	0.25	0	0	0	u1		0	0
	b1		0.05	0.1	0.15	0	0	0	b1		0	0
	b2		0.3	0.05	0.35	0	0	0	b2		0	0
Portuguese	pt	He	She	OTHER	pt	He	She	OTHER	pt	He	She	OTHER
	u1		0.05	0	0.05	0	0	0	u1		0.05	0
	b1		0	0	0	0	0	0	b1		0	0
	b2		0.05	0	0.05	0	0	0	b2		0.1	0
French	fr	He	She	OTHER	fr	He	She	OTHER	fr	He	She	OTHER
	u1		0.35	0.3	0.2	0	0	0	u1		0	0.1
	b1		0.4	0.45	0.2	0	0.05	0.05	b1		0	0.15
	b2		0.25	0.1	0.2	0	0	0	b2		0.15	0
Arabic	ar	He	She	OTHER	ar	He	She	OTHER	ar	He	She	OTHER
	u1		0.25	0.2	0.35	0.3	0.25	0.45	u1		0.45	0.7
	b1		0.05	0.4	0.4	0.05	0.5	0.55	b1		0.1	0.95
	b2		0.3	0.05	0.3	0.55	0.05	0.55	b2		0.65	0.15

Final Ratings

No.	M_i	Rating (T_1)	Rating (T_2)	Rating (T_3)
1.	es	DSBS	DSBS	DSBS
2.	pt	DSBS	DSBS	DSBS
3.	fr	DSBS	DSBS	UCS
4.	ar	DSBS	DSBS	DSBS
	Overall	DSBS	DSBS	DSBS

But How Do People Perceive Ratings ? - VEGA Environment

[Video]

Try the tool at: <http://vega-live.mybluemix.net/>

- Mariana Bernagozzi, Biplav Srivastava, Francesca Rossi and Sheema Usmani, VEGA: a Virtual Environment for Exploring Gender Bias vs. Accuracy Trade-offs in AI Translation Services, **AAAI 2021**. [Visualizing Ethics Rating, *Demonstration paper*]
- Mariana Bernagozzi, Biplav Srivastava, Francesca Rossi and Sheema Usmani, Gender Bias in Online Language Translators: Visualization, Human Perception, and Bias/Accuracy Trade-offs, **IEEE Internet Computing, Special Issue on Sociotechnical Perspectives**, Nov/Dec 2021 [Visualizing Ethics Rating, User Survey]

Survey of Translator Users

ar	es	fr	pt	4
ar	es	fr	pt	86
ar	es	fr	pt	29
ar	es	fr	pt	5
ar	es	fr	pt	1
ar	es	fr	pt	1
ar	es	fr	pt	10
ar	es	fr	pt	1
ar	es	fr	pt	2
ar	es	fr	pt	14

QC No.	Question Category	Result	Comments
1.	Do people have unconscious bias with respect to gender?	Yes (52%)	
2.	Do people perceive gender bias in online translators?	Yes	Can be language specific
3.	Are people perceiving the correct notion of gender bias in a translator?	Yes	
4.	Do people appreciate a visual representation of a gender bias assessment?	Yes	
5.	Are people more inclined to use a translator when they are presented with a bias rating visualization?	Yes	
6.	Is bias or accuracy more important when choosing a translator?	Almost the same	Can be language specific

Why Do the Results Matter



Translators are a common form of AI, used in composition with others



Their output can be perceived as biased and the phenomena is language specific



People like visual tools to help them recognize bias



Bias perception equally important as accuracy

Examples of Sequential Composition

Composed Task 1: *generate sentiment of an image*

- System 1:
 - Input: image
 - Output: text describing the image
- System 2:
 - Input: text
 - Output: sentiment (positive or negative)

Composed Task 2: *generate multi-lingual caption of an image*

- System 1:
 - Input: image
 - Output: text describing the image
- System 2:
 - Input: text
 - Output: text (in another language)

Debiasing Approach

Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (NIPS'16)*. Curran Associates Inc., Red Hook, NY, USA, 4356–4364.

Problem: Bias in Learnt Word Representation

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{computer programmer}} - \overrightarrow{\text{homemaker}}.$$

Key points:

- Recall the analogies created in last class
- Corpus used in paper: 300 dimensional embedding trained on a corpus of Google News texts consisting of 3 million English words
- Focuses on F-M gender pairs
- Word differences measured by cosine similarity, vectors are normalized
- From 3 million original words, authors select 50,000 most frequent words. Then selects lower-case words and phrases consisting of fewer than 20 lower-case characters (words with upper-case letters, digits, or punctuation were discarded). After this filtering, 26,377 words remained.

Analogy Generation Task

$$S_{(a,b)}(x,y) = \cos(\vec{a} - \vec{b}, \vec{x} - \vec{y}) \text{ if } \|\vec{x} - \vec{y}\| \leq \delta, \quad 0 \text{ else}$$

Key points:

(a, b) = (she, he).

For each analogy, MT workers were asked two yes/no questions:

- (a) whether the pairing makes sense as an analogy, and
- (b) whether it reflects a gender stereotype.

Examples of Gendered Words and Analogies

Extreme <i>she</i>	Extreme <i>he</i>	Gender stereotype <i>she-he</i> analogies		
1. homemaker	1. maestro	sewing-carpentry	registered nurse-physician	housewife-shopkeeper
2. nurse	2. skipper	nurse-surgeon	interior designer-architect	softball-baseball
3. receptionist	3. protege	blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
4. librarian	4. philosopher	giggle-chuckle	vocalist-guitarist	petite-lanky
5. socialite	5. captain	sassy-snappy	diva-superstar	charming-affable
6. hairdresser	6. architect	volleyball-football	cupcakes-pizzas	lovely-brilliant
7. nanny	7. financier	Gender appropriate <i>she-he</i> analogies		
8. bookkeeper	8. warrior	queen-king	sister-brother	mother-father
9. stylist	9. broadcaster	waitress-waiter	ovarian cancer-prostate cancer	convent-monastery
10. housekeeper	10. magician			

Figure 1: **Left** The most extreme occupations as projected on to the *she—he* gender direction on w2vNEWS. Occupations such as *businesswoman*, where gender is suggested by the orthography, were excluded. **Right** Automatically generated analogies for the pair *she-he* using the procedure described in text. Each automatically generated analogy is evaluated by 10 crowd-workers to whether or not it reflects gender stereotype.

Illustration: Discovered Gendered Words

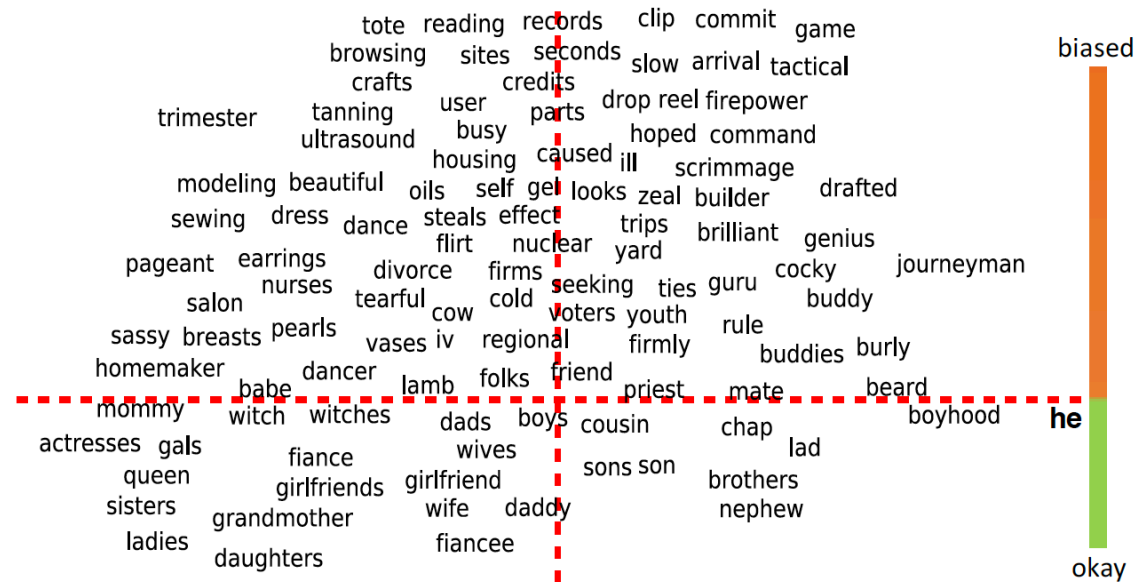


Figure 3: Selected words projected along two axes: x is a projection onto the difference between the embeddings of the words *he* and *she*, and y is a direction learned in the embedding that captures gender neutrality, with gender neutral words above the line and gender specific words below the line. Our hard debiasing algorithm removes the gender pair associations for gender neutral words. In this figure, the words above the horizontal line would all be collapsed to the vertical line.

Debiasing Method

- Hard debiasing
 - Identify subspace of gendered words
 - Replacing extreme gendered words with neutral words and normalize
- Soft debiasing
 - Identify subspace of gendered words
 - preserve pairwise inner products between all the word vectors while minimizing the projection of the gender-neutral words onto the gender subspace.

Effectiveness of Debiasing

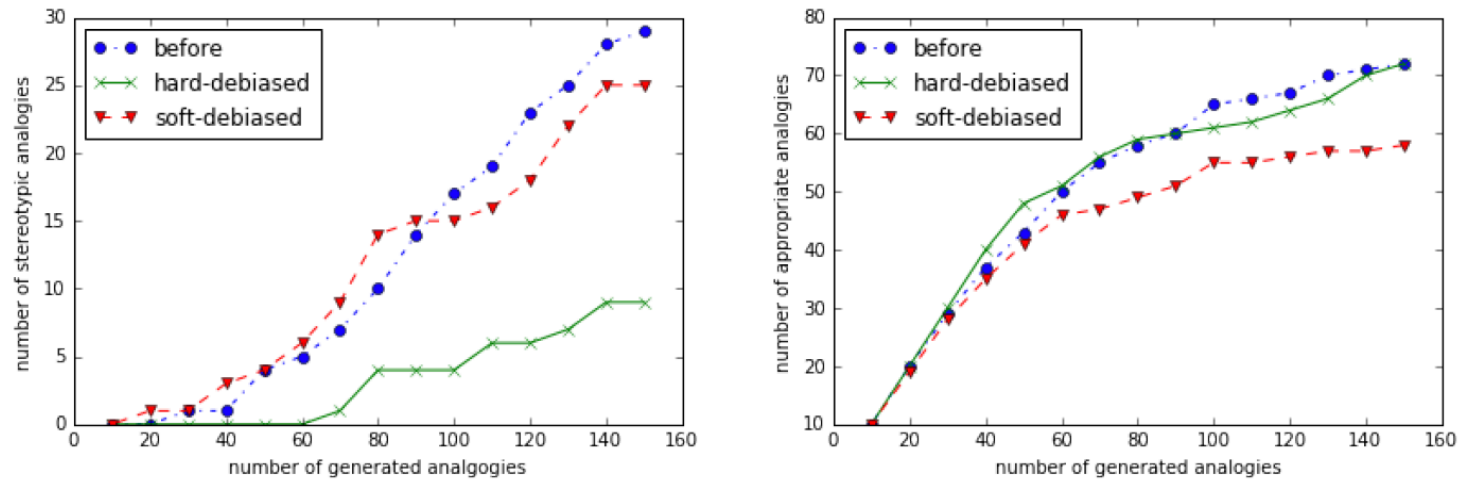


Figure 4: Number of stereotypical (Left) and appropriate (Right) analogies generated by word embeddings before and after debiasing.

Concluding Segment

Lecture 22: Concluding Comments

- We looked at rating methods for characterizing machine translators
- We reviewed paper on de-biasing learned word representations

About Next Lecture – Lecture 24

Schedule Snapshot



Oct 26 (Tu)	Review: Explanation Methods, AIX 360, Discussion	Quiz 3
Oct 28 (Th)	Review: project presentations, Discussion	
Nov 2 (Tu)	AI - Unstructured (Text): Analysis – Supervised ML – Trust Issues	
Nov 4 (Th)	AI - Unstructured (Text): Analysis – Supervised ML – Mitigation Methods	
Nov 9 (Tu)	AI - Unstructured (Text): Analysis – Rating and Debiasing Methods	
Nov 11 (Th)	Explanation Methods Trust: AI Testing	
Nov 16 (Tu)	Trust: Human-AI Collaboration	Quiz 4
Nov 18 (Th)	Paper presentations – Graduate students	Final assignment for Graduate students
Nov 23 (Tu)	Emerging Standards and Laws	

Lecture 24:

- Explanation methods for Unstructured Text
- AI Testing