Human-Centered NLP (Language and Human Psychology)

CSE 538

NLP, The Course

Overall NLP Concept

I. Syntax

II. Semantics

Overall NLP Concept

III. Language Modeling

IV. Applications

NLP, The Course



I. Syntax

Introduction to NLP; Tokenization; Words Corpora

One-hot, and Multi-hot encoding. Parts-of-Speech; Named Entities;

Parsing; Verbal Predicates; Dependency Parsing

II. Semantics

Dependency Parsing; Word Sense Disambiguation

Vector Semantics (Embeddings), Word2vec

Probabilistic Language Models Ngram Classifier, Topic Modeling

Overall NLP Concept

III. Language Modeling

IV. Applications

NLP The Course

Probabilistic Language Models

Ngram Classifier, Topic Modeling

BI

Overall NLP Concept

Overall NLP Concept	Computation or ML			
I. Syntax Classification				
Introduction to NLP; Tokenization; Words Corpora	Regular Expressions; Edit Distance			
One-hot, and Multi-hot encoding. Parts-of-Speech; Named Entities;	Maximum Entropy Classifier (LogReg), Gradient Descent,			
Parsing; Verbal Predicates; Dependency Parsing	Cross Validation; Regularization Accuracy Metrics; Shift Reduce			
II. Semantics Probabilistic Models				
Dependency Parsing; Word Sense Disambiguation	Term Probabilities; N-d Vectors			
Vector Semantics (Embeddings), Word2vec	LDA, Skipgram Model			

markov assumption, chain

rule, smoothing

III. Language Modeling Transformers				
Ethical Considerations	Model cards, Pred Bias Frmwrk			
Masked Language Modeling (autoencoding)	Neural Networks; Backprop Cross-Entropy Loss Self-Attention,			
Generative Language Modeling (autoregressive)	Positional encodings The Transformer: Beam Search			
Applying LMs	Fine-Tuning, zero-/few-shot, Instruction tuning			
IV. Applications Custom Statistical or Symbolic				
Language and Psychology (advanced sentiment)	Differential Language Analysis; Adaptive Modeling; Human LMing			
Speech and Audio Processing, Dialog (chatbots)	Wave Transforms; RNNs			
Question Answering, Translation	Multihop Reasoning			

Computation or ML

NLP The Course



	<u> </u>			
<u>Överall NLP Concept</u>	Computation or ML			
I. Syntax Classification				
Introduction to NLP; Tokenization; Words Corpora	Regular Expressions; Edit Distance			
One-hot, and Multi-hot encoding. Parts-of-Speech; Named Entities;	Maximum Entropy Classifier (LogReg), Gradient Descent,			
Parsing; Verbal Predicates; Dependency Parsing	Cross Validation; Regularization Accuracy Metrics; Shift Reduce			
II. Semantics Probabilistic Models				

Dependency Parsing; Word Sense Disambiguation	Termi Probabilities; N-d Vectors
Vector Semantics (Embeddings), Word2vec	LDA, Skipgram Model
Probabilistic Language Models Ngram Classifier, Topic Modeling	markov assumption, chain rule, smoothing

Overall NLP Concept	Computation or ML			
III. Language Modeling Transformers				
Ethical Considerations	Model cards, Pred Bias Frmwrk			
Masked Language Modeling (autoencoding)	Neural Networks; Backprop Cross-Entropy Loss Self-Attention,			
Generative Language Modeling (autoregressive)	Positional encodings The Transformer: Beam Search			
Applying LMs	Fine Tuning, zero-/few-shot, Instruction tuning			

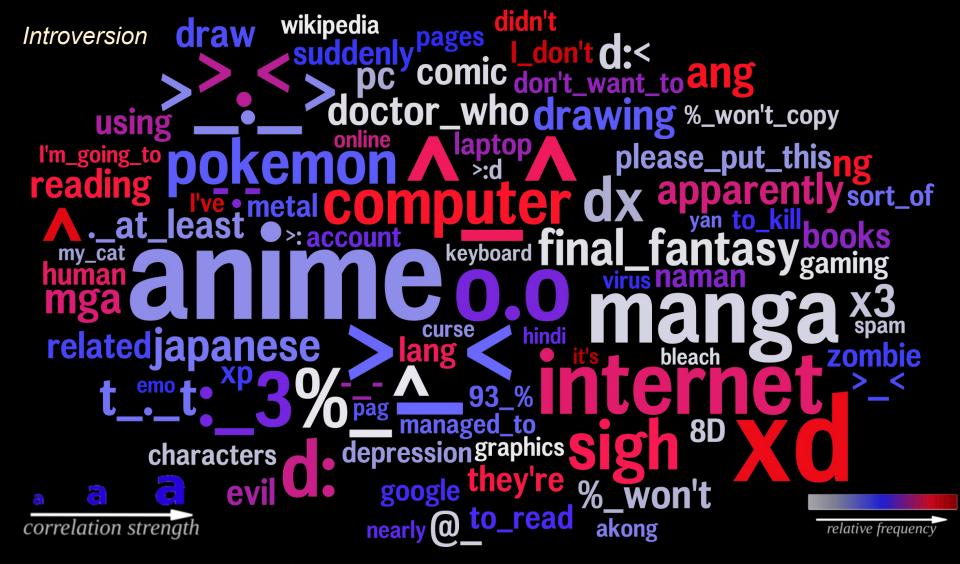
IV. Applications | Custom Statistical or Symbolic

Language and Psychology •• (advanced sentiment)	Differential Language Analysis; Adaptive Modeling; Human LMing	
Speech and Audio Processing, Dialog (chatbots)	Wave Transforms; RNNs	
Question Answering, Translation	Multihop Reasoning	

Natural Language Processing

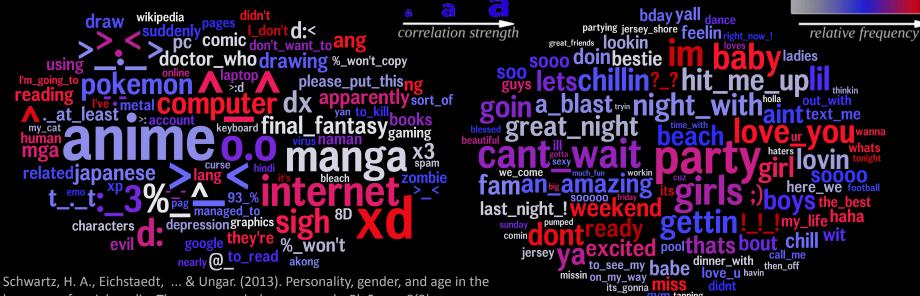
Psychological & Health Sciences







Psychological & Health Sciences

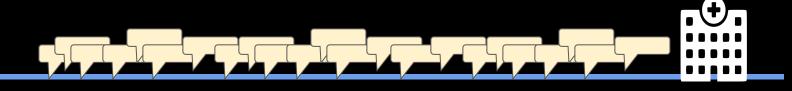


language of social media: The open-vocabulary approach. PloS one, 8(9).



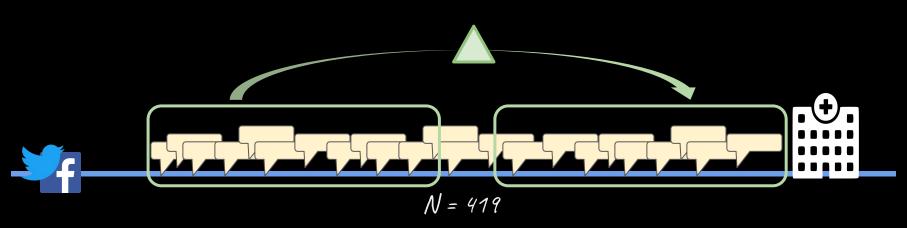
Psychological & Health Sciences



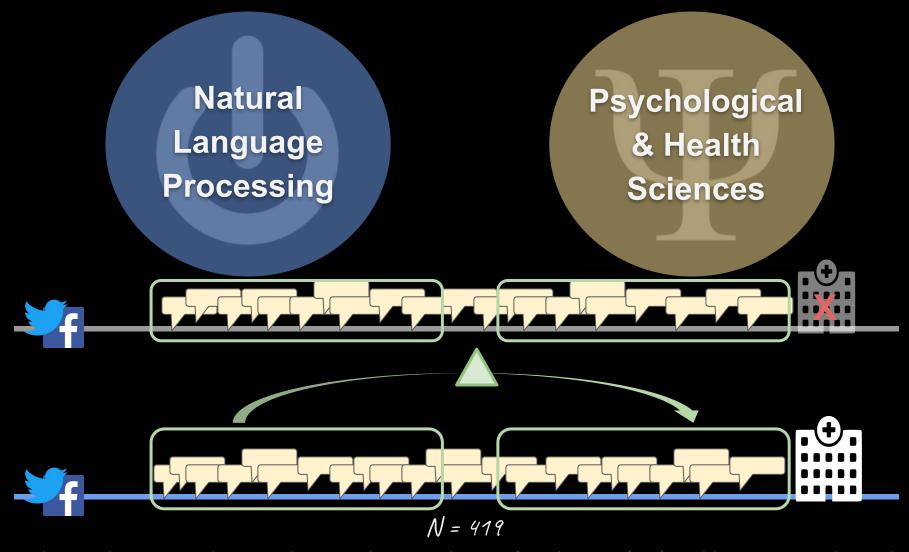


Natural Language Processing

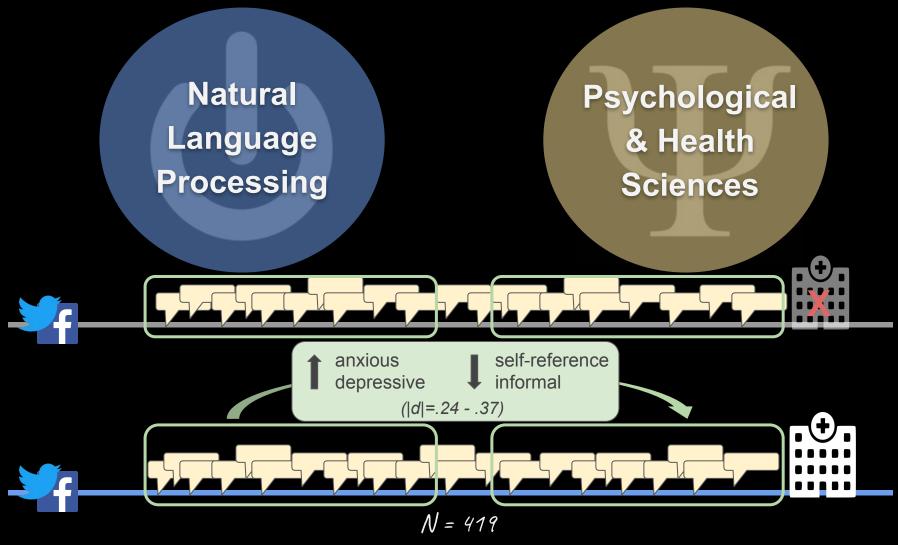
Psychological & Health Sciences



Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports*, 10(1), 1-9.



Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports*, 10(1), 1-9.



Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports*, 10(1), 1-9.

Natural Language Processing

Psychological & Health Sciences

Overly Simplified Problem-Statement:

Natural language is written by

Overly Simplified Problem-Statement:

Natural language is written by **people**.

Overly Simplified Problem-Statement:

Natural language is written by people.

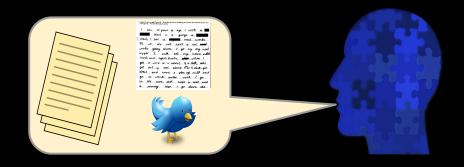


Problem

Natural language is written by people.

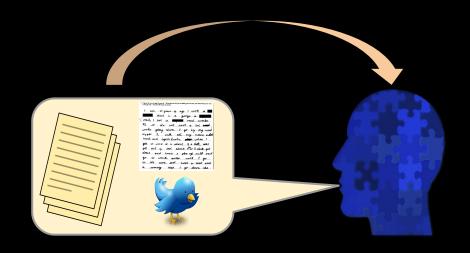


Natural language is generated by people.



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, ...

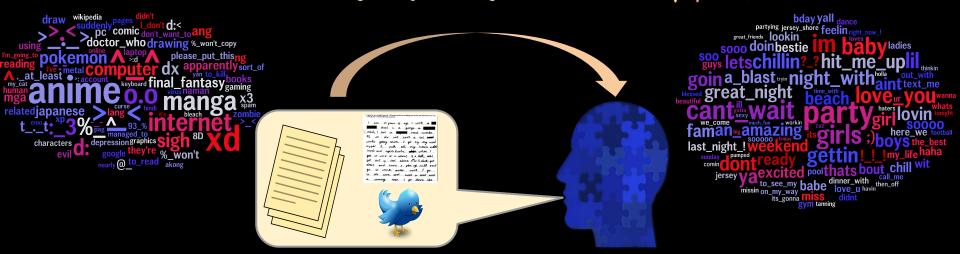
Natural language is generated by people.



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, ...,

and our language reflects these differences.

Natural language is generated by people.



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, ...,

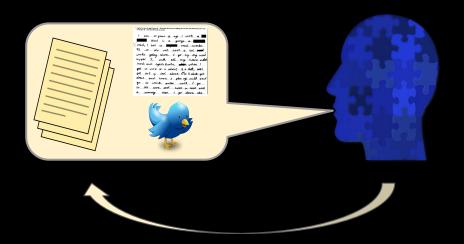
and our language reflects these differences.

Human Centered NLP:



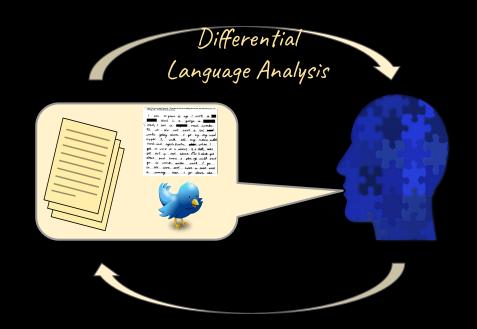
Human Centered NLP:

1. Model language as a human process



Human Centered NLP:

- 1. Model language as a human process
- 2. Use language to better understand humans.



Human-Centered NLP – We will cover:

- 1. Differential Language Analysis
- 2. Human Factor Adaptation
- 3. Human Language Modeling

Input:

Linguistic features

Human or community attribute

Output:

Features distinguishing attribute

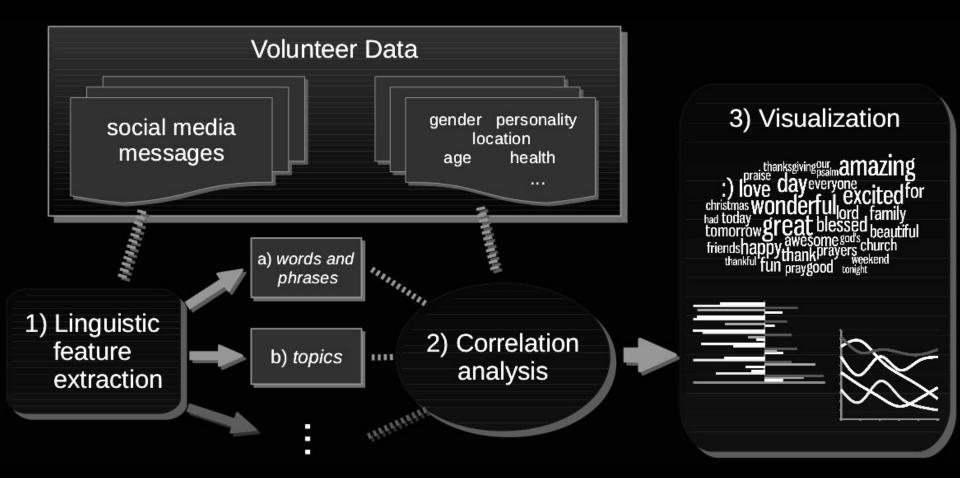
Goal: Data-driven insights about an attribute

E.g. Words distinguishing communities with increases in real estate prices.



a a a correlation strength





Methods of Correlation Analysis:

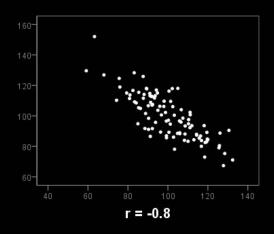
Pearson Product-Moment Correlation
 Limitation: Doesn't handle controls

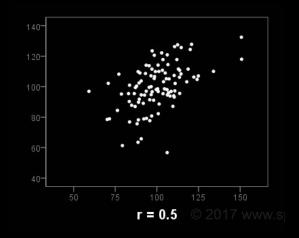
$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

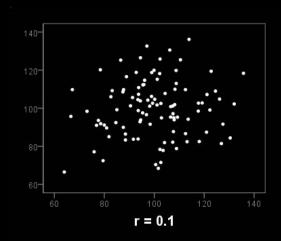
Methods of Correlation Analysis:

Pearson Product-Moment Correlation
 Limitation: Doesn't handle controls

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$







Methods of Correlation Analysis:

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

- Pearson Product-Moment Correlation
 Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Methods of Correlation Analysis:

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i-ar{x})(y_i-ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i-ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i-ar{y})^2}}$$

- Pearson Product-Moment Correlation
 Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Adjust all variables to have "mean center" and "unit variance":

Methods of Correlation Analysis:

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

- Pearson Product-Moment Correlation
 Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Adjust all variables to have "mean center" and "unit variance":

$$z = \frac{x - \mu}{\sigma}$$

$$\mu=$$
 Mean $\sigma=$ Standard Deviation

Methods of Correlation Analysis:

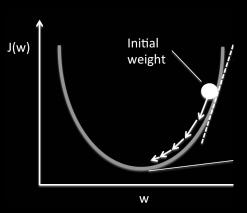
$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

- Pearson Product-Moment Correlation
 Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Option 1: Gradient Descent:

$$J = \sum (y - \hat{y})^2$$
 -- "Sum of Squares" Error



Methods of Correlation Analysis:

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

- Pearson Product-Moment Correlation Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

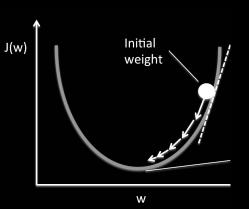
Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Option 1: Gradient Descent:

$$J = \sum (y - \hat{y})^2$$
 -- "Sum of Squares" Error

Option 2: Matrix model:
$$Y = X\beta + \epsilon$$



Methods of Correlation Analysis:

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

- Pearson Product-Moment Correlation Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Option 1: Gradient Descent:

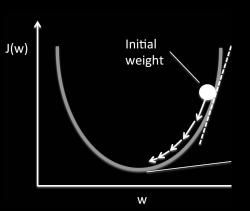
$$J = \sum (y - \hat{y})^2$$
 -- "Sum of Squares" Error

Option 2: Matrix model: $Y = X\beta + \epsilon$

$$Y = X\beta + \epsilon$$

Matrix Computation Solution:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$



Methods of Correlation Analysis:

$$r_{xy} = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

- Pearson Product-Moment Correlation Limitation: Doesn't handle controls
- Standardized Multivariate Linear Regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

Option 1: Gradient Descent:

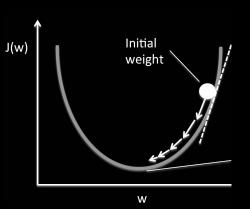
$$J = \sum (y - \hat{y})^2$$
 -- "Sum of Squares" Error

Option 2: Matrix model:
$$Y = X\beta + \epsilon$$

$$Y = X\beta + \epsilon$$

Matrix Computation Solution:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$



Methods of "Correlation" Analysis for binary outcomes:

- Logistic Regression over Standardized variables
- Odds Ratio

$$\frac{countA("horrible")}{NA} = \frac{countA("horrible")}{NA}$$

```
\frac{\frac{countB("horrible")}{NB}}{1-\frac{countB("horrible")}{NB}}
```

Methods of "Correlation" Analysis for binary outcomes:

- Logistic Regression over Standardized variables
- Odds Ratio

$$\frac{countA("horrible")}{NA} = \frac{countA("horrible")}{NA}$$

$$countA("horrible") \over 1 - \frac{countA("horrible")}{NA} } - log \left(\frac{\frac{countB("horrible")}{NB}}{1 - \frac{countB("horrible")}{NB}} \right)$$

$$\frac{countB("horrible")}{NB} = \frac{1 - \frac{countB("horrible")}{NB}}{NB}$$

Methods of "Correlation" Analysis for binary outcomes:

- Logistic Regression over Standardized variables
- Odds Ratio

$$\frac{\frac{countA("horrible")}{NA}}{1-\frac{countA("horrible")}{NA}}$$

$$\cos \log \left(\frac{\frac{countA("horrible")}{NA}}{1 - \frac{countA("horrible")}{NA}} \right) - \log \left(\frac{\frac{countB("horrible")}{NB}}{1 - \frac{countB("horrible")}{NB}} \right)$$

$$\frac{countB("horrible")}{NB} = \frac{countB("horrible")}{NB}$$

$$= log \left(\frac{countA("horrible")}{NA-countA("horrible")} \right) - log \left(\frac{countB("horrible")}{NB-countB("horrible")} \right)$$

$$log\left(\frac{countA("horrible")}{NA-countA("horrible")}\right) - log\left(\frac{countB("horrible")}{NB-countB("horrible")}\right)$$

Odds Ratio using Informative Dirichlet Prior

$$\delta_w^{(i-j)} = \log\left(\frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)}\right) - \log\left(\frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)}\right)$$
(20.9)

$$log\left(\frac{countA("horrible")}{NA-countA("horrible")}\right) - log\left(\frac{countB("horrible")}{NB-countB("horrible")}\right)$$

Odds Ratio using <u>Informative Dirichlet Prior</u>

$$\delta_{w}^{(i-j)} = \log\left(\frac{f_{w}^{i} + \alpha_{w}}{n^{i} + \alpha_{0} - (f_{w}^{i} + \alpha_{w})}\right) - \log\left(\frac{f_{w}^{j} + \alpha_{w}}{n^{j} + \alpha_{0} - (f_{w}^{j} + \alpha_{w})}\right)$$
(20.9)

(where n^i is the size of corpus i, n^j is the size of corpus j, f_w^i is the count of word w in corpus i, f_w^j is the count of word w in corpus j, α_0 is the size of the background corpus, and α_w is the count of word w in the background corpus.)

$$log\left(\frac{countA("horrible")}{NA-countA("horrible")}\right) - log\left(\frac{countB("horrible")}{NB-countB("horrible")}\right)$$

Odds Ratio using <u>Informative Dirichlet Prior</u>

$$\delta_w^{(i-j)} = \log \left(\frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)} \right)$$

(where n^i is the size of corpus i, n^j is the sin corpus i, f_w^j is the count of word w in corpus, and α_w is the count of word w in

$$\operatorname{gg}\left(\frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)}\right) \qquad (20.9)$$

ous j, f_w^i is the count of word w is the size of the background corpus.)

Bayesian term for "smoothing": accounts for uncertainty as a function of event frequency (i.e. words observed less) by integrating "prior" beliefs mathematically.

$$log\left(\frac{countA("horrible")}{NA-countA("horrible")}\right) - log\left(\frac{countB("horrible")}{NB-countB("horrible")}\right)$$

Odds Ratio using <u>Informative Dirichlet Prior</u>

$$\delta_w^{(i-j)} = \log \left(\frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)} \right)$$

(where n^i is the size of corpus i, n^j is the sin corpus i, f_w^j is the count of word w in corpus, and α_w is the count of word w in

$$g\left(\frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)}\right) \qquad (20.9)$$

ous j, f_w^i is the count of word w is the size of the background corpus.)

Bayesian term for "smoothing": accounts for uncertainty as a function of event frequency (i.e. words observed less) by integrating "prior" beliefs mathematically.

"Informative": the prior is based on past evidence. Here, the total frequency of the word.

$$log\left(\frac{countA("horrible")}{NA-countA("horrible")}\right) - log\left(\frac{countB("horrible")}{NB-countB("horrible")}\right)$$

Odds Ratio using Informative Dirichlet Prior

$$\delta_w^{(i-j)} = \log\left(\frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)}\right) - \log\left(\frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)}\right)$$
(20.9)

(where n^i is the size of corpus i, n^j is the size of corpus j, f_w^i is the count of word w in corpus i, f_w^j is the count of word w in corpus j, α_0 is the size of the background corpus, and α_w is the count of word w in the background corpus.)

$$log\left(\frac{countA("horrible")}{NA-countA("horrible")}\right) - log\left(\frac{countB("horrible")}{NB-countB("horrible")}\right)$$

Odds Ratio using Informative Dirichlet Prior

$$\delta_{w}^{(i-j)} = \log \left(\frac{f_{w}^{i} + \alpha_{w}}{n^{i} + \alpha_{0} - (f_{w}^{i} + \alpha_{w})} \right) - \log \left(\frac{f_{w}^{j} + \alpha_{w}}{n^{j} + \alpha_{0} - (f_{w}^{j} + \alpha_{w})} \right)$$
(20.9)

(where n^i is the size of corpus i, n^j is the size of corpus j, f_w^i is the count of word w in corpus i, f_w^j is the count of word w in corpus j, α_0 is the size of the background corpus, and α_w is the count of word w in the background corpus.)

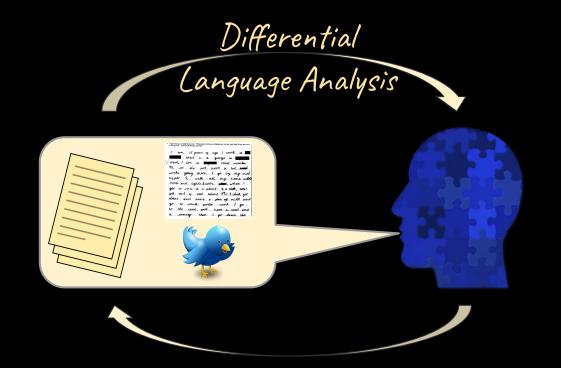
Final score is standardized (z-scored): $\hat{\delta}_w^{(i-j)}$, where $\sqrt{\sigma^2\left(\hat{\delta}_w^{(i-j)}\right)} \qquad \sigma^2\left(\hat{\delta}_w^{(i-j)}\right) \approx \frac{1}{f_w^i + \alpha_w} + \frac{1}{f_w^j + \alpha_w}$

(Monroe et al., 2010; Jurafsky, 2017)

D

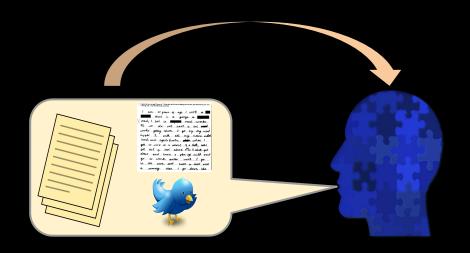
DLATK

Python Library, CLI, and Colab for DLA

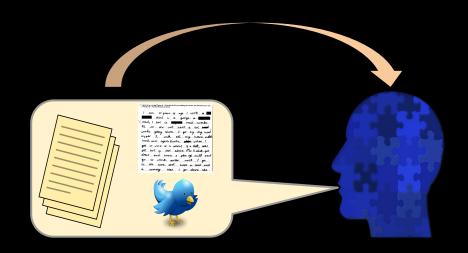


https://dlatk.github.io/Getting Started in Colab

Natural language is generated by people.

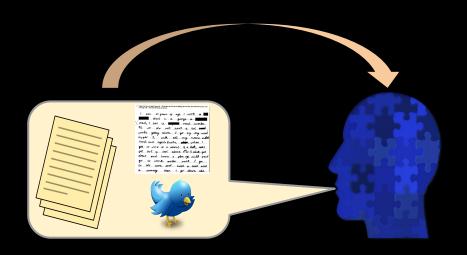


Natural language is generated by people.





Natural language is generated by people.



"The common misconception is that language has got to do with words and what they mean. It does not. It has to do with people and what they mean."

Shannon, 1948

Mosteller & Wallace 1963 Clark & Schober, 1992

Mairesse, Walker, et al., 2007 Hovy & Soogaard, 2015

Human-Centered NLP – We will cover:

- 1. Differential Language Analysis
- 2. Human Factor Adaptation
- 3. Human Language Modeling

Human-Centered NLP – We will cover:

- 1. Differential Language Analysis
- 2. Human Factor Adaptation
- 3. Human Language Modeling

1. **Bias Mitigation:** Optimize so as not to pick up on unwanted relationships.

(e.g. image captioner label pictures of men in kitchen as women)

1. **Bias Mitigation:** Optimize so as not to pick up on unwanted relationships.

(e.g. image captioner label pictures of men in kitchen as women)

2. Additive: Include direct effect of human factor on outcome.

(e.g. age and distinguishing PTSD from Depression; covariate in regression)

1. **Bias Mitigation:** Optimize so as not to pick up on unwanted relationships.

(e.g. image captioner label pictures of men in kitchen as women)

- 2. Additive: Include direct effect of human factor on outcome. (e.g. age and distinguishing PTSD from Depression)
- Adaptive: Allow meaning of language to change depending on human context. (also called "compositional")
 - (e.g. "sick" said from a young individual versus old individual)

1. Bias Mitigation: Optimus as not to pick up on

What are human "factors"?

(e.g. image captioner laber pictures or mem n kitchen as women)

- 2. Additive: Include direct effect of human factor on outcome. (e.g. age and distinguishing PTSD from Depression)
- Adaptive: Allow meaning of language to change depending on human context. (also called "compositional")
 - (e.g. "sick" said from a young individual versus old individual)

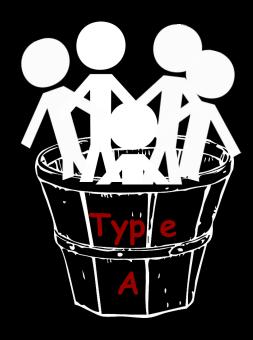
Human Factors

--- Any attribute, represented as a continuous or discrete variable, of the humans generating the natural language.

E.g.

- Gender
- Age
- Personality
- Ethnicity
- Socio-economic status

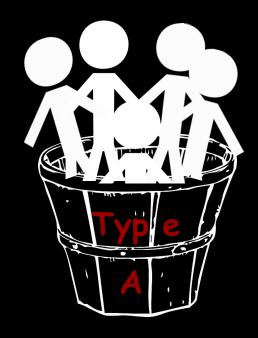
Human Factors

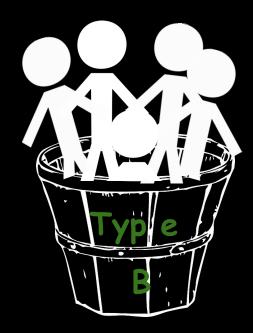




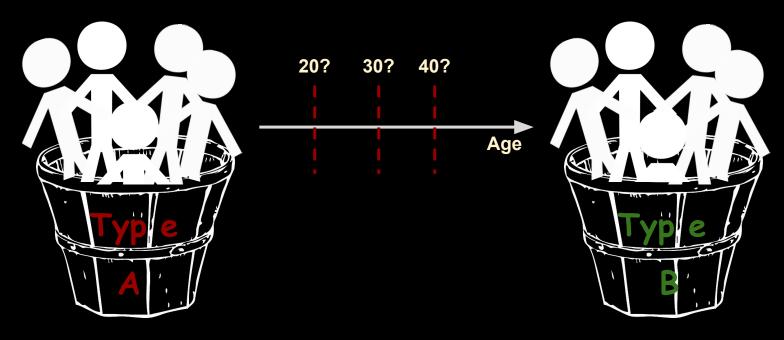
typically requires putting people into discrete bins

"most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]" (Haslam et al., 2012)

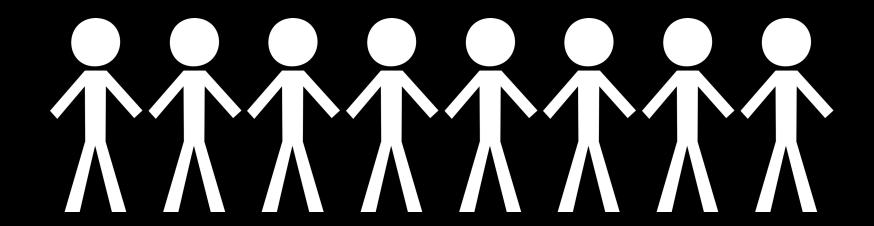




"most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]" (Haslam et al., 2012)



"most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]" (Haslam et al., 2012)



Less Factor A



More Factor A



Adaptation Approach: Domain Adaptation

Features for: source target
$$\begin{matrix} & & & \\ & & & \\ & & & \\ \Phi^s(x) = \langle x, x, \mathbf{0} \rangle, & \Phi^t(x) = \langle x, \mathbf{0}, x \rangle \end{matrix}$$

Frustratingly Easy Domain Adaptation

Hal Daumé III

School of Computing University of Utah Salt Lake City, Utah 84112 me@hal3.name

Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case supervised case. The fully supervised case models the following scenario. We have access to a large, annotated corpus of data for

Adaptation Approach: Domain Adaptation

Features for: source target
$$\begin{matrix} & & & \\ & & & \\ & & & \\ \Phi^s({\bm x}) = \langle {\bm x}, {\bm x}, {\bm 0} \rangle, & \Phi^t({\bm x}) = \langle {\bm x}, {\bm 0}, {\bm x} \rangle \end{matrix}$$

Frustratingly Easy Domain Adaptation

Hal Daumé III

School of Computing University of Utah Salt Lake City, Utah 84112 me@hal3.name

Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case els the following scenario. We have access to a large, annotated corpus of data for

Adaptation Approach: Domain Adaptation

```
Features for: source target newX = [] for all x in source_x: newX.append(x + x + [0]*len(x)) for all x in target_x newX.append(x + [0]*len(x), x) newY = source_y + target_y model = model.train(newX,newY)
```

Adaptation Approach: Factor Adaptation

Human Centered NLP with User-Factor Adaptation

Veronica E. Lynn, Youngseo Son, Vivek Kulkarni Niranjan Balasubramanian and H. Andrew Schwartz {velynn, yson, vvkulkarni, niranjan, has}@cs.stonybrook.edu

Abstract

We pose the general task of user-factor adaptation — adapting supervised learning models to real-valued user factors in-

and Costa Jr., 1989; Ruscio and Ruscio, 2000;

Here, we ask how one can adapt NLP models to real-valued human factors – continuous valued attributes that capture fine-grained differences be-

Residualized Factor Adaptation for Community Social Media Prediction Tasks

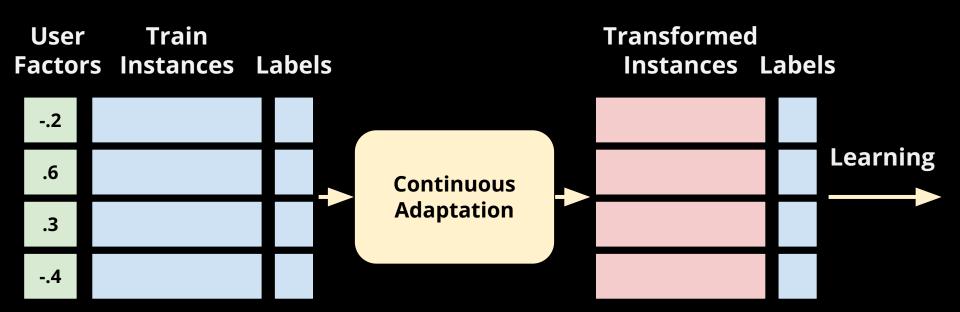
Mohammadzaman Zamani, H. Andrew Schwartz, Veronica E. Lynn, Salvatore Giorgi, and Niranjan Balasubramanian Computer Science Department, Stony Brook University ²Department of Psychology, University of Pennsylvania mzamani@cs.stonybrook.edu

Abstract

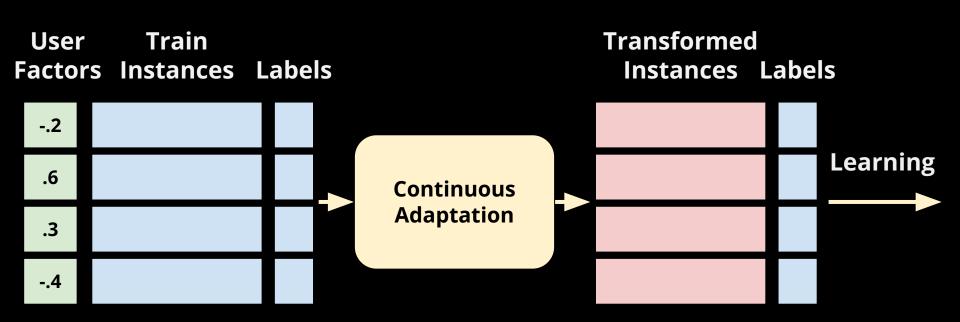
Predictive models over social media language

linked to socio-demographic factors (age, gender, race, education, income levels) with many social scientific studies supporting their predictive

Our Method: Continuous Adaptation

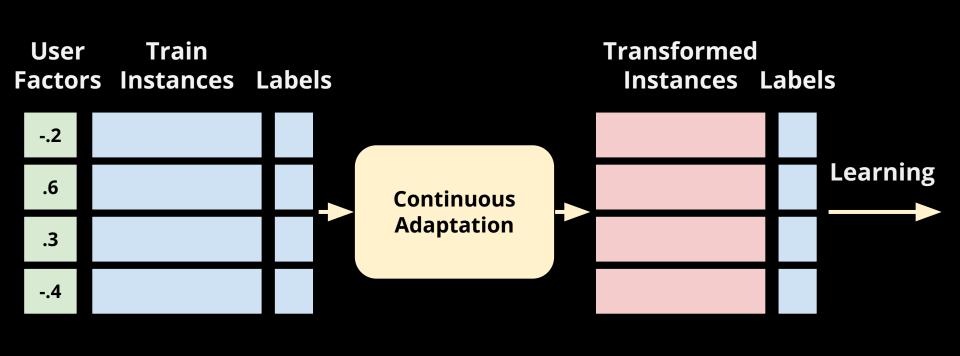


Our Method: Continuous Adaptation





Our Method: Continuous Adaptation



Gender Score Features Original Gender Copy
-.2 X X compose(-.2, X)

(Lynn et al., 2017)

User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function c combines d user factor scores $f_{u,d}$ with original feature values \mathbf{x} :

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \cdots, c(f_{u,d}, \mathbf{x}) \rangle$$

User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function c combines d user factor scores $f_{u,d}$ with original feature values \mathbf{x} :

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \cdots, c(f_{u,d}, \mathbf{x}) \rangle$$

T.T	D	A 1 T		
User	Factor	Augmented Instance		
	Classes	$\Phi(\mathbf{x}, u)$		
User 1	$\overline{F_1}$	$\langle \mathbf{x}, \mathbf{x}, 0, 0, \cdots, 0 \rangle$		
User 2	F_2	$\langle \mathbf{x}, 0, \mathbf{x}, 0, \cdots, 0 \rangle$		
User 3	F_1, F_3	$\langle \mathbf{x}, \mathbf{x}, 0, \mathbf{x}, \cdots, 0 \rangle$		
User 4	F_k	$\langle \mathbf{x}, 0, 0, \cdots, 0, \mathbf{x} \rangle$		

Table 1: Discrete Factor Adaptation: Augmentations of an original instance vector \mathbf{x} under different factor class mappings. With k domains the augmented feature vector is of length n(k+1).

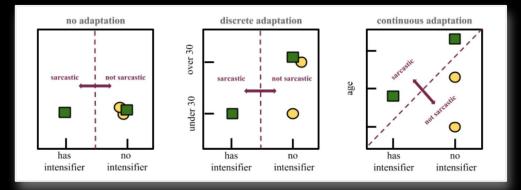
(Lynn et al., 2017)

User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function c combines d user factor scores $f_{u,d}$ with original feature values \mathbf{x} :

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \cdots, c(f_{u,d}, \mathbf{x}) \rangle$$



User	Factor	Augmented Instance		
	Classes	$\Phi(\mathbf{x}, u)$		
User 1	$\overline{F_1}$	$\langle \mathbf{x}, \mathbf{x}, 0, 0, \cdots, 0 \rangle$		
User 2	F_2	$\langle \mathbf{x}, 0, \mathbf{x}, 0, \cdots, 0 \rangle$		
User 3	F_1, F_3	$\langle \mathbf{x}, \mathbf{x}, 0, \mathbf{x}, \cdots, 0 \rangle$		
User 4	F_k	$\langle \mathbf{x}, 0, 0, \cdots, 0, \mathbf{x} \rangle$		

Table 1: Discrete Factor Adaptation: Augmentations of an original instance vector \mathbf{x} under different factor class mappings. With k domains the augmented feature vector is of length n(k+1). (Lynn et al., 2017)

Main Results

Adaptation improves over unadapted baselines (Lynn et al., 2017)

Task	Metric	No Adaptation	Gender	Personality	Latent (User Embed)
Stance	F1	64.9	65.1 (+0.2)	66.3 (+1.4)	67.9 (+3.0)
Sarcasm	F1	73.9	75.1 (+1.2)	75.6 (+1.7)	77.3 (+3.4)
Sentiment	Acc.	60.6	61.0 (+0.4)	61.2 (+0.6)	60.7 (+0.1)
PP-Attach	Acc.	71.0	70.7 (-0.3)	70.2 (-0.8)	70.8 (-0.2)
POS	Acc.	91.7	91.9 (+0.2)	91.2 (-0.5)	90.9 (-0.8)

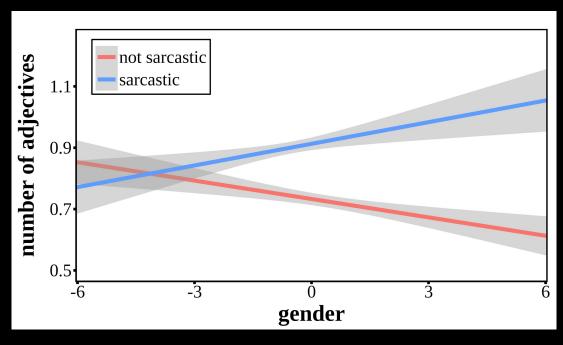
Example: How Adaptation Helps

<u>Women</u>

more adjectives→sarcasm

<u>Men</u>

more adjectives→no sarcasm



more "male"

more "female"

Problem

User factors are not always available.

Solution: User Factor Inference

past tweets

Niranjan @b_niranjan · Sep 2

There must be a word for trending #hashtags that you know you will regret if you click. Is there?

Niranjan @b_niranjan · Aug 31

Passwords spiral: Forget password for the acnt you use twice a year. Ask for reset. Can't use previous. Create a new one to forget later.

Niranjan @b_niranjan · Jul 31

Thrilled to hear @acl2017's diversity efforts as the first thing in the conference.



Known

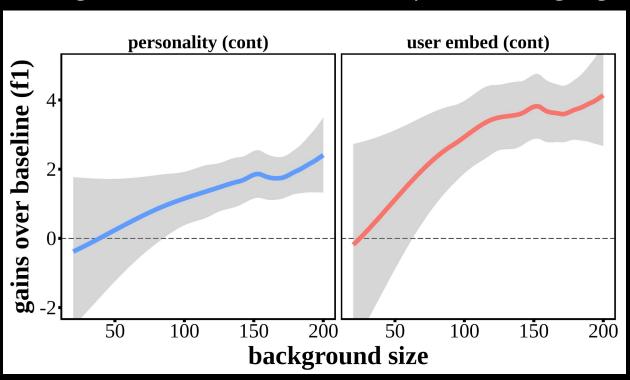
Age (Sap et al. 2014) Gender (Sap et al. 2014) Personality (Park et al. 2015)

<u>Latent</u>

User Embeddings (Kulkarni et al. 2017) Word2Vec TF-IDF

Background Size

Using more background tweets to infer factors produces larger gains



Approaches to Human Factor Inclusion

1. **Bias Mitigation:** Optimize so as not to pick up on unwanted relationships.

(e.g. image captioner label pictures of men in kitchen as women)

- 2. Additive: Include direct effect of human factor on outcome. (e.g. age and distinguishing PTSD from Depression)
- Adaptive: Allow meaning if language to change depending on human context. (also called "compositional")
 - (e.g. "sick" said from a young individual versus old individual)

Approaches to Human Factor Inclusion

1. **Bias Mitigation:** Optimize so as not to pick up on unwanted relationships.

(e.g. image captioner label pictures of men in kitchen as women)

- 2. Additive: Include direct effect of human factor on outcome. (e.g. age and distinguishing PTSD from Depression)
- 3. Adaptive: Allow meaning if language to change depending on human context. (also called "compositional") (e.g. "sick" said from a young individual versus old individual)

Human-Centered NLP – We will cover:

- 1. Differential Language Analysis
- 2. Human Factor Adaptation
- 3. Human Language Modeling

Human-Centered NLP – We will cover:

- 1. Differential Language Analysis
- 2. Human Factor Adaptation
- 3. Human Language Modeling



probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1})$$

Language is generated by $w_0 = w_1 = w_2 = w_3 = w_4$



Language is generated by humans w_1 w_2 people machines $Pr(w_4|w_{0:3})$ writers Al (~1 million words) 0.05 0.10 0.15 0.20 next word probability

0.25

(GPT4o)

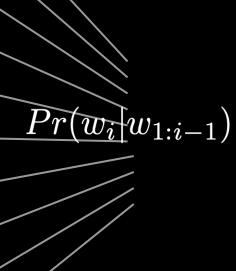
0.30

0.35

	"Morning run by the ocean, amazing!"	
Morning	"Team meeting went great, big plans!"	
Morning run	"Fortnite royale! That's sick!"	
Morning run by	"Just did my first handstand pushup! 😄"	
Morning run by the	"So Proud, Launched our new product."	
Morning run by the ocean	"Sprints on the sand = leg day!"	$\underbrace{Pr(w_i]}_{}w_{1:i-1})$
	"Just unlocked a legendary skin! 😎"	
	"Brainstorming ideas for the next project."	
	"Can't seem to beat that high score."	
	"Sunset yoga at the beach, relaxing!"	

"Morning run by the ocean, amazing!" "Team meeting went great, big plans!" "Fortnite royale! That's sick!" "Just did my first handstand pushup! 😄" "So Proud, Launched our new product." "Sprints on the sand = leg day!" "Just unlocked a legendary skin! "Brainstorming ideas for the next project." "Can't seem to beat that high score."

"Sunset yoga at the beach, relaxing!"





"Morning run by the ocean, amazing!"

"Team meeting went great, big plans!"

"Fortnite royale! That's sick!"

"Just did my first handstand pushup! 😄"

"So Proud, Launched our new product."

"Sprints on the sand = leg day!"

"Just unlocked a legendary skin! 😎"

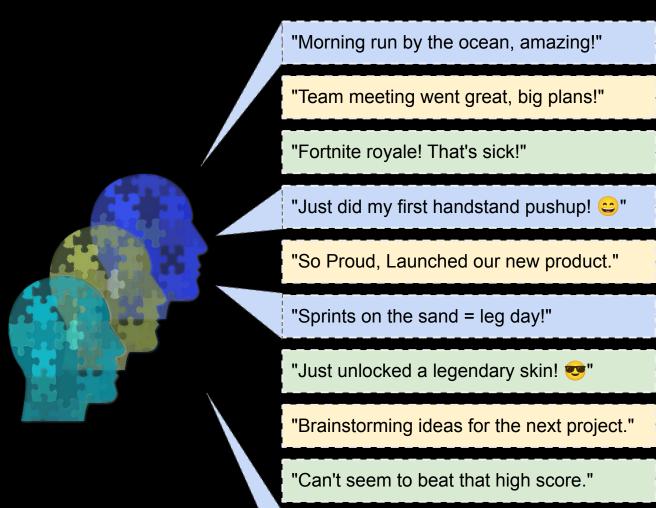
"Brainstorming ideas for the next project."

"Can't seem to beat that high score."

"Sunset yoga at the beach, relaxing!"

 $Pr(w_i|w_{1:i-1})$

"Morning run by the ocean, amazing!"	
"Team meeting went great, big plans!"	
"Fortnite royale! That's sick!"	
"Just did my first handstand pushup! 😄"	
"So Proud, Launched our new product."	
"Sprints on the sand = leg day!"	$Pr(w_i w_{1:i-}$
"Just unlocked a legendary skin! 😎"	
"Brainstorming ideas for the next project."	
"Can't seem to beat that high score."	
"Sunset yoga at the beach, relaxing!"	



"Sunset yoga at the beach, relaxing!"

 $\widehat{Pr(w_i|w_{1:i-1})}$



"Team meeting went great, big plans!"

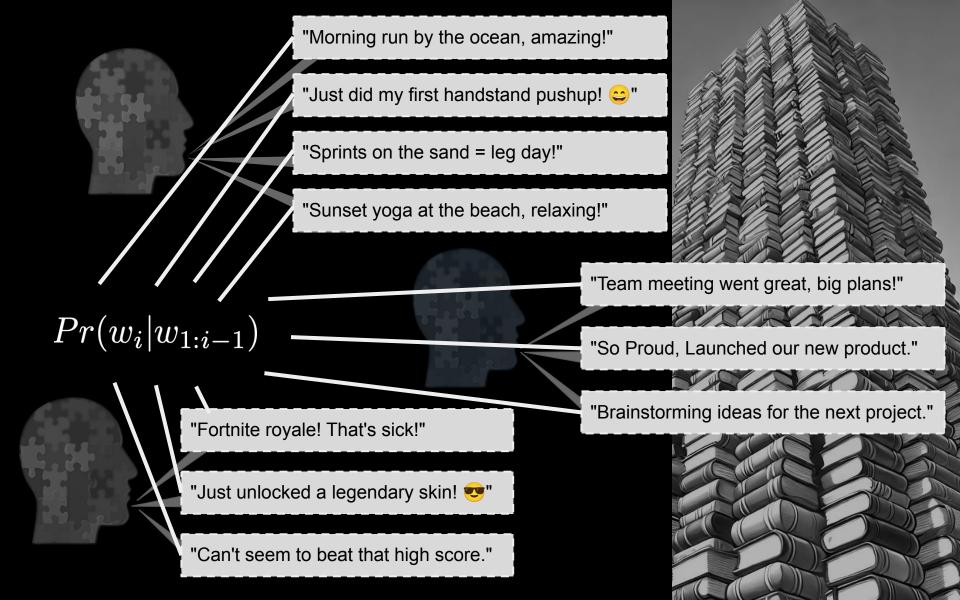
"So Proud, Launched our new product."

"Brainstorming ideas for the next project."

"Fortnite royale! That's sick!"

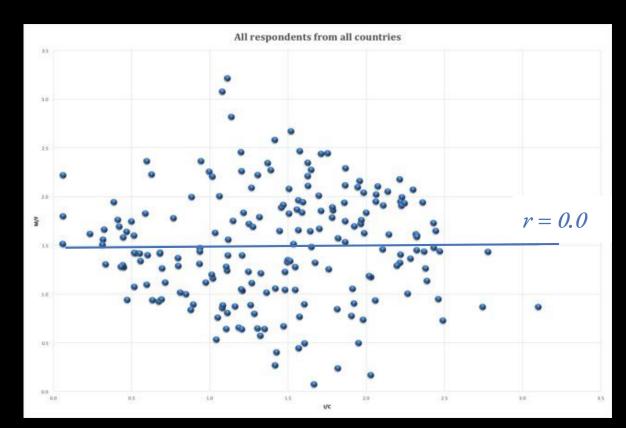
"Just unlocked a legendary skin! 😎"

"Can't seem to beat that high score."



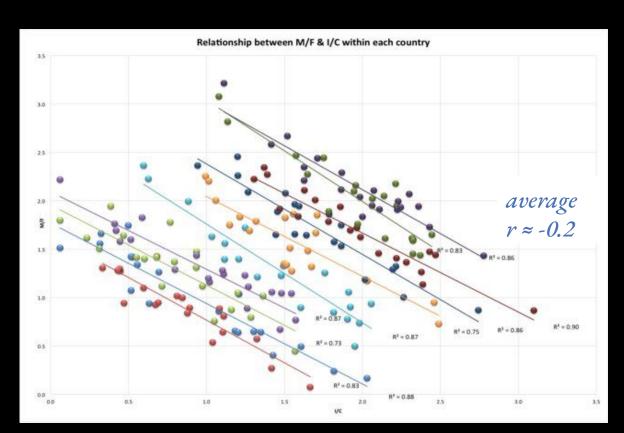






(winzar, 2015)

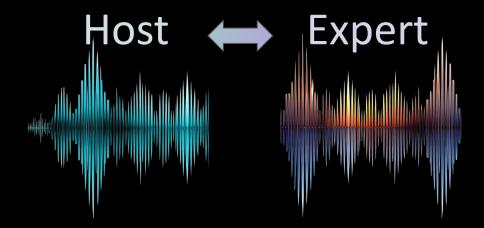




(winzar, 2015)

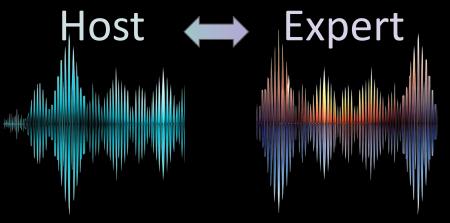


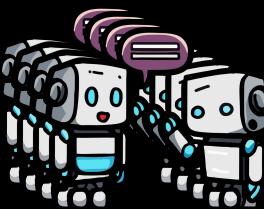
NotebookLM



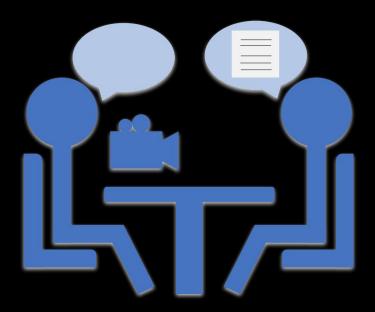


NotebookLM



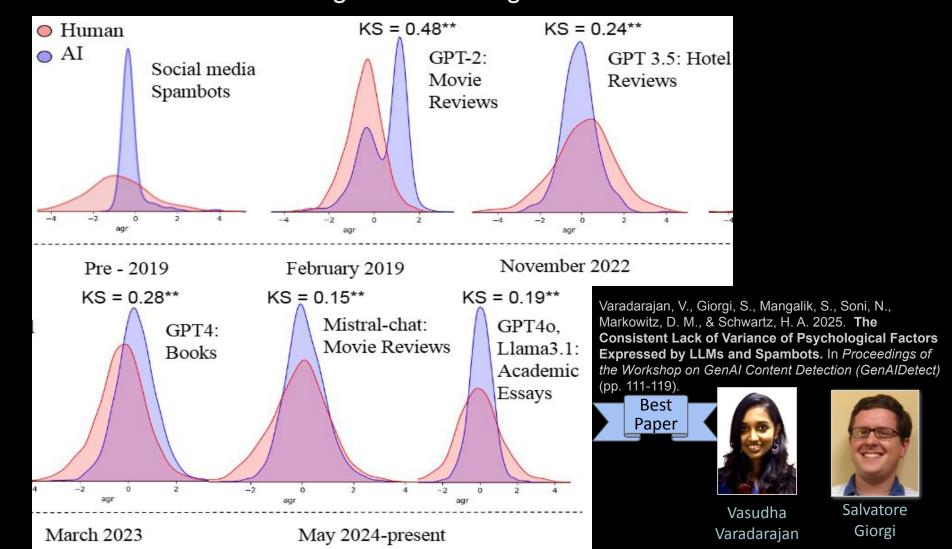


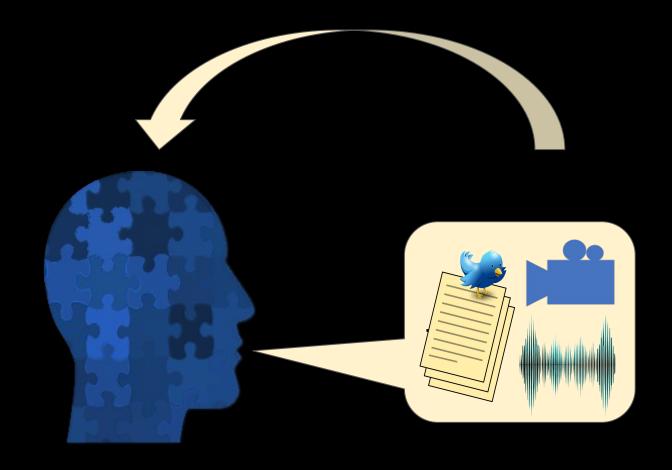




speaker diarization

Distribution of human linguistic trait - agreeableness over datasets





People have beliefs, backgrounds, styles, vocabularies, knowledge, and personalities. Our data reflect and are influenced by these differences.

probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1})$$

LM

- probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1})$$

Hull - probability of a token sequence, in the context of the human that generated it.

 LN

- probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1})$$

$$Pr(\mathbf{W}|\mathbf{U}_{static}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1}, \mathbf{U}_{static})$$

HuLM

- probability of a token sequence, in the context of the human that generated it.

 LN

- probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^n Pr(w_i|w_{1:i-1})$$

$$Pr(\mathbf{W}|\mathbf{U}_{static}) = \prod_{i=1}^n Pr(w_i|w_{1:i-1}, \mathbf{U}_{static})$$
 static user representation

HuLM

- probability of a token sequence, in the context of the human that generated it.

LM

- probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1})$$

$$Pr(\mathbf{W}|\mathbf{U}_{static}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1}, \mathbf{U}_{static})$$
 static user representation

HuLM

$$Pr(\mathbf{W}_t|\mathbf{U}_{t-1}) = \prod_{i=1}^{n} Pr(w_{t,i}|w_{t,1:i-1},\mathbf{U}_{1:t-1})$$

- probability of a token sequence, in the context of the human that generated it.

 LN

- probability of a token sequence

$$Pr(\mathbf{W}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1})$$

$$Pr(\mathbf{W}|\mathbf{U}_{static}) = \prod_{i=1}^{n} Pr(w_i|w_{1:i-1}, \mathbf{U}_{static})$$
 static user representation

HuLM

$$Pr(\mathbf{W}_t|\mathbf{U}_{t-1}) = \prod_{i=1} Pr(w_{t,i}|w_{t,1:i-1},\mathbf{U}_{1:t-1})$$
"user state" representation

probability of a token sequence, in the context of the human that generated it.

User State Representation, U

$$Pr(\mathbf{W}_{t}|\mathbf{U}_{t-1}) = \prod_{i=1}^{n} Pr(w_{t,i}|w_{t,1:i-1}, \mathbf{U}_{1:t-1})$$



$$U_{1:t-1} = \emptyset$$

(reduces to a standard LM: $Pr(w_i|w_{1:i-1})$)

 $U_{1:t-1} = w_{1,1:n_1}, w_{2,1:n_2}, ..., w_{t-1,1:n_{t-1}}$

(all previous docs and tokens by the person)

- doesn't capture the person

history of user states

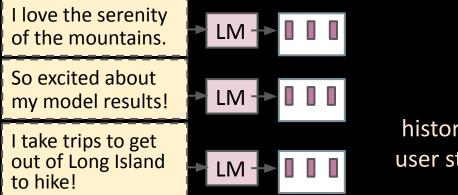
- huge
- no generalizations

Soni, N., Matero, M., Balasubramanian, N., & Schwartz, H. (2022, May). Human Language Modeling. In *Findings of the Association for Computational Linguistics: ACL 2022* (pp. 622-636).

User State Representation, U

$$Pr(\mathbf{W}_{t}|\mathbf{U}_{t-1}) = \prod_{i=1}^{n} Pr(w_{t,i}|w_{t,1:i-1},\mathbf{U}_{1:t-1})$$

State and Trait Theory from Psychology: *Traits* – the stable characteristics of "who someone is" — define a distribution of potential **states** of being that moderate human behavior (i.e. language).



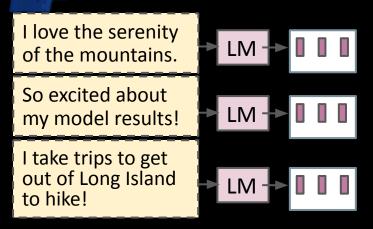
history of user states

Soni, N., Matero, M., Balasubramanian, N., & Schwartz, H. (2022, May). Human Language Modeling. In Findings of the Association for Computational Linguistics: ACL 2022 (pp. 622-636).

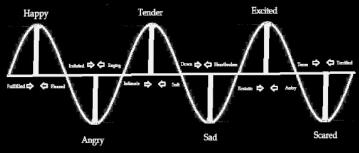
User State Representation, U

$$Pr(\mathbf{W}_{t}|\mathbf{U}_{t-1}) = \prod_{i=1}^{n} Pr(w_{t,i}|w_{t,1:i-1},\mathbf{U}_{1:t-1})$$

State and Trait Theory from Psychology: *Traits* – the stable characteristics of "who someone is" – define a distribution of potential *states* of being that moderate human behavior (i.e. language).



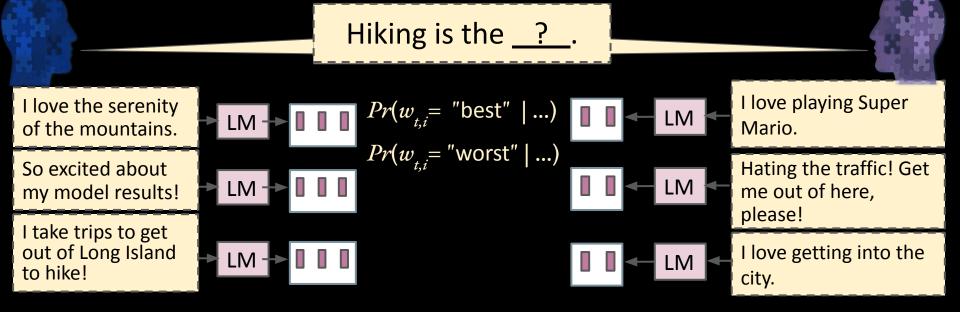
$$U_{1:t-1}$$
 = [a sequence of *states*]



(Washington Outsider, 2014)

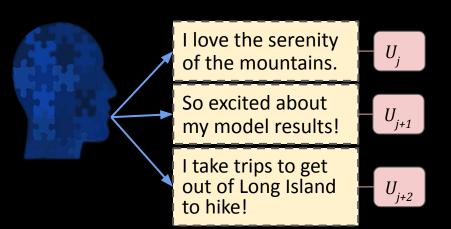
User State Representation, U

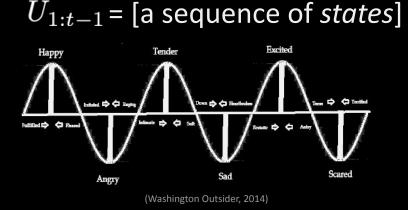
$$Pr(\mathbf{W}_{t}|\mathbf{U}_{t-1}) = \prod_{i=1}^{n} Pr(w_{t,i}|w_{t,1:i-1}, \mathbf{U}_{1:t-1})$$



User State Representation: Motivation

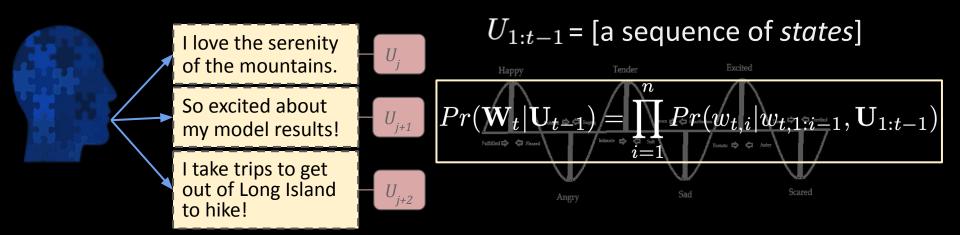
- Addressing Ecological Fallacy: Treating dependent phenomena (i.e. sequences from the same person) as if independent. (Piantadosi et al., 1988; Steel and Holt, 1996)
- Modeling the higher order structure.
- Building on ideas from human factor inclusion/adaptation (Lynn et al., 2017; Huang & Paul, 2019; Hovy & Yang, 2021) and personalized modeling. (King & Cook, 2020; Jaech & Ostendorf, 2018)





Human Language Modeling (HuLM)

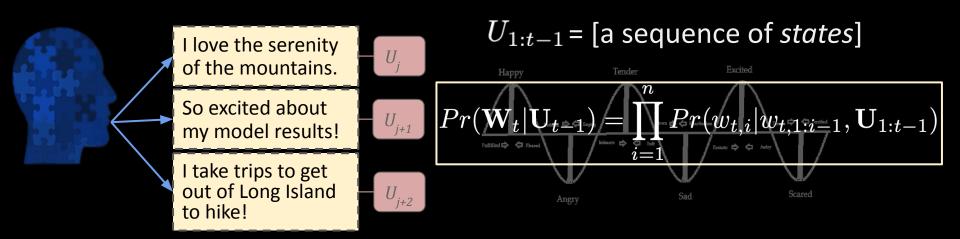
- Addressing Ecological Fallacy: Treating dependent phenomena (i.e. sequences from the same person) as if independent. (Piantadosi et al., 1988; Steel and Holt, 1996)
- Modeling the higher order structure.
- Building on ideas from human factor inclusion/adaptation (Lynn et al., 2017; Huang & Paul, 2019; Hovy & Yang, 2021) and personalized modeling. (King & Cook, 2020; Jaech & Ostendorf, 2018)

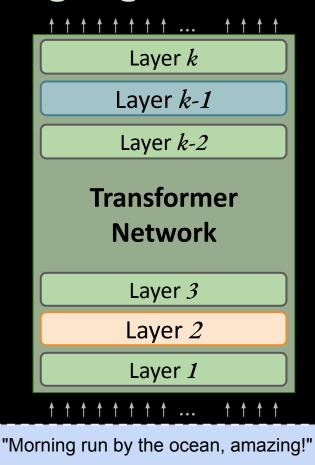


Human Language Modeling (HuLM)

Goal: Language modeling as a task grounded in the "natural" generators of language, people.

The HuLM task definition: Estimate the probability of a sequence of tokens, $w_{t,1:i'}$ conditioned on a higher-order representation, U_t , constituting the human state of being just before the sequence generation.









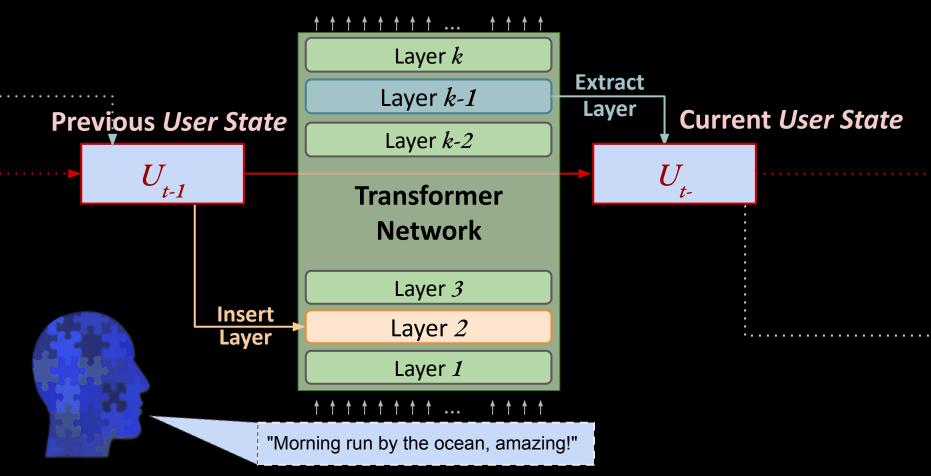


Matthew Matero

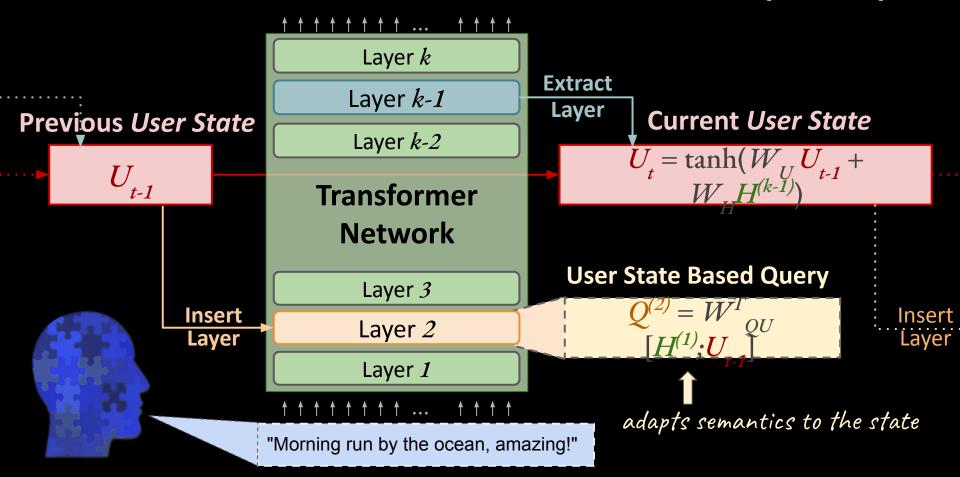


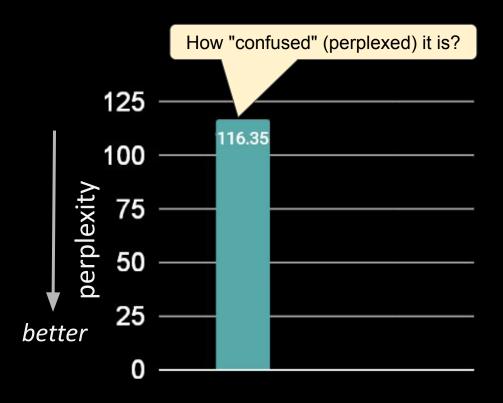
Niranjan Balasubramanian

Human-aware Recurrent Transformer (HaRT)



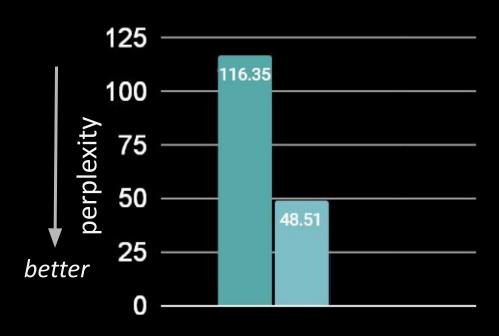
Human-aware Recurrent Transformer (HaRT)





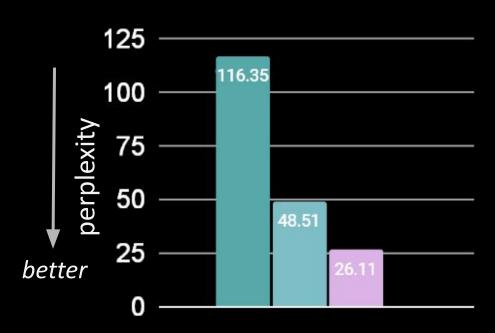




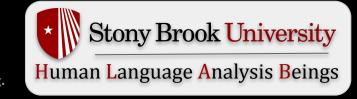


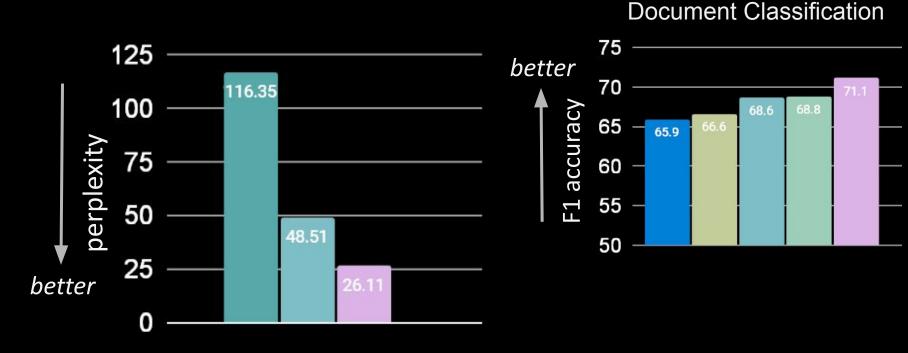








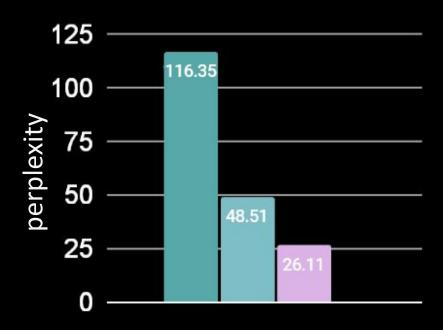






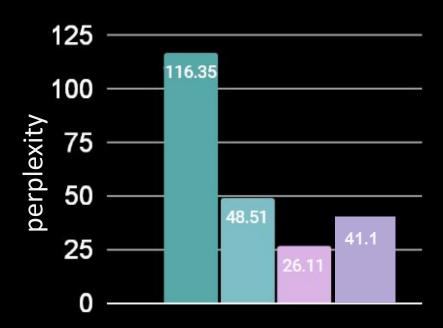


No history?



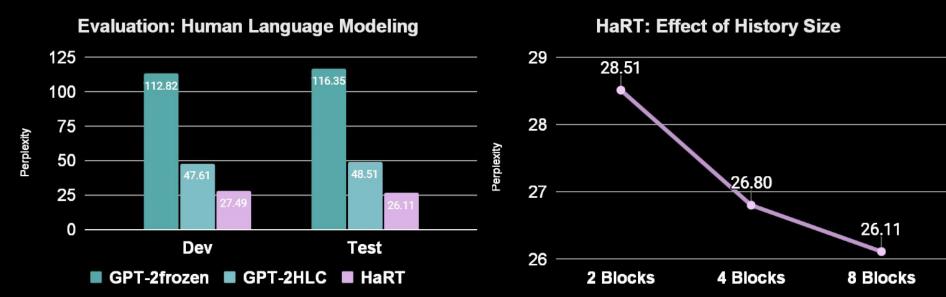
■ GPT-2frozen ■ GPT-2HLC ■ HaRT ■ HaRT-no_history

No history?



■ GPT-2frozen ■ GPT-2HLC ■ HaRT ■ HaRT-no_history

Human Language Modeling



Dataset: Human Language Corpus (HLC)

Soni, N., Matero, M., Balasubramanian, N., & Schwartz, H. (2022, May). Human Language Modeling. In Findings of the Association for Computational Linguistics: ACL 2022 (pp. 622-636).



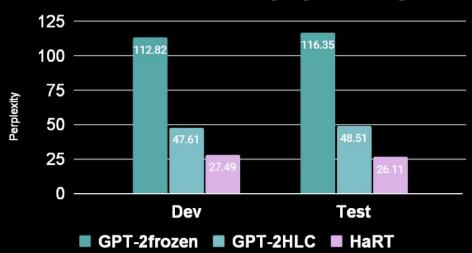
Train users = 96k msgs = 36m (8 blocks= ~17m)

Dev users = 2k msgs = 830k + seen users: 2.5k msgs: 230k

Test users = 2k msgs = 690k + seen users: 2.5K msgs: 240k

Human Language Modeling

Evaluation: Human Language Modeling



Dataset: Human Language Corpus (

Soni, N., Matero, M., Balasubramanian, N., & Schwartz, H. (2022, May). Human Language Modeling. In Findings of the Association for Computational Linguistics: ACL 2022 (pp. 622-636).

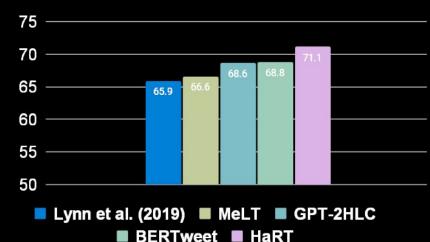


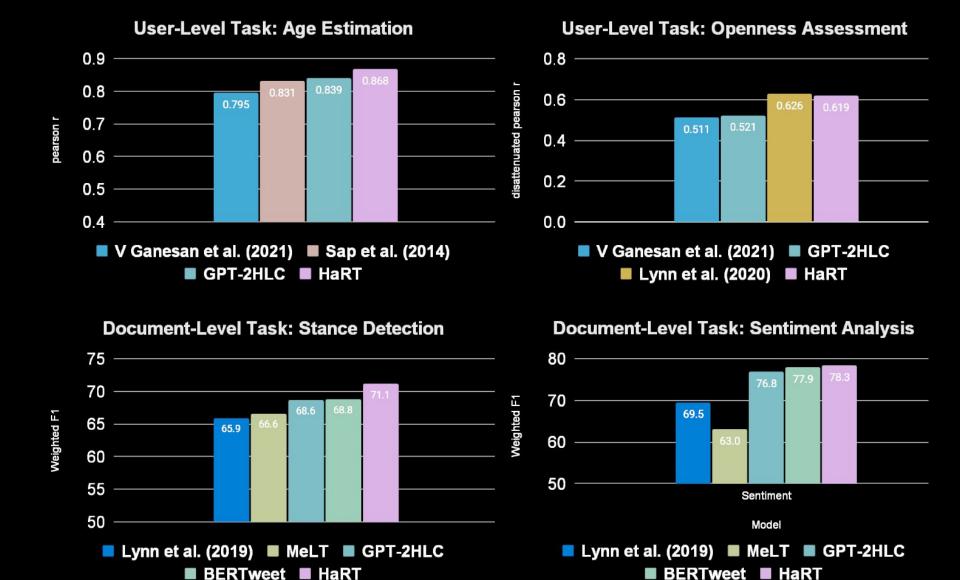


HaRT: Effect of History Size



Document-Level Task: Stance Detection





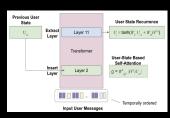
HuLM/HaRT Takeaways

- HuLM: Extension of language modeling with notion of user.
- HaRT: A step toward large human language models.

 Progress for large LMs grounded in language's "natural" generators, people.

GitHub Repository





Human Language Modeling

Nikita Soni, Matthew Matero, Niranjan Balasubramanian, and H. Andrew Schwartz Department of Computer Science, Stony Brook University

{nisoni, mmatero, niranjan, has}@cs.stonybrook.edu

Abstract

Natural language is generated by people, yet traditional language modeling views words or documents as if generated independently. Here, we propose human language modeling (HuLM), a hierarchical extension to the language modeling problem whereby a humanvel exists to connect sequences of documents (e.g. social media messages) and capture the notion that human language is moderated by changing human states. We introduce, HaRT, a large-scale transformer model for the HULM task, pre-trained on approximately 100,000 social media users, and demonstrate it's effectiveness in terms of both language modeling (perplexity) for social media and fine-tuning for 4 downstream tasks spanning documentand user-levels: stance detection, sentiment classification, age estimation, and personality assessment. Results on all tasks meet or surnass the current state-of-the-art.

To address this, we introduce the task of human language modeling (HULM), which induces dependence among text sequences via the notion of a human state in which the text was generated. In particular, we formulate HULM as the task of estimating the probability of a sequence of tokens, while conditioning on a higher order state which will be considered from the tokens of other documents written by the same individual. Its key objective is:

$Pr(w_{t,i}|w_{t,1:i-1}, \mathbf{U}_{1:t-1})$

where t indexes a particular sequence of temporally ordered utterances (e.g. a document or social media post), and U_{tt-1} represents the human state just before the current sequence, t. In one extreme, U_{tt-1} could model all previous tokens in all previous documents by the person. In the opposite extreme, U_{tt-1} could model all previous for site extreme, U_{tt-1} came the same for all users and for values of t reducing to standard language modeline: $P(tt)_{tt} |_{tt-1} > T$ Thus, $W(tt)_{tt} |_{tt-1} > T$ Thus, $W(tt)_{tt} |_{tt-1} > T$

Human-Centered NLP – Methods to know

- 1. Differential Language Analysis
 - a. pearson correlation
 - b. multivariate linear regression
 - c. odd ratio with informative Dirichlet prior
- 2. Human Factor Adaptation feature augmentation
- 3. Human Language Modeling
 - a. HuLM task definition
 - b. HaRT Architecture