Designing ML Week 7: Natural Language Processing + Conversational Interfaces

michelle.carney@berkeley.edu





 \mathcal{O}

 \bigcirc

Tom Bibic @bibic tom · 2h Replying to @realDonaldTrump

You should attack the press more! It makes you sound more leaderly! #sarcasm

11 C \square



Anastasia K @ghost_intheroom · Feb 19 For some of us sarcasm is a secound language. At times even first... #introvert #sarcasm

 \square



Headlong @Headlong42 · 22h

V

V

SOCIAL MEDIA

Researchers develop sarcasm detector for Twitter, and that's no joke



Replying to @Headlong42 @ScoutFinchS60 @kylegriffin1

P.S.: Long threads don't necessarily imply rambling. Some people actually have something to say. Crazy right? And what's even more crazy: some people are actually interested in reading & able to comprehend them... #sarcasm

 \square



17 07 \square



David Ávila @davilovick · 23h Wow, the video compression in Twitter are "awesome" #sarcasm

 $\bigcirc 2$ 11

Show this thread

Genie73 @Genie731 · 23h



Replying to @Nuecents1 @ QUEENish

Yea, Nuisance is about right. You are the vile and selfish one and you damn sure do NOT need to have any type of hand in rearing children to teach them to be good and kind humans. I wish you well. #sarcasm



Estimating the Date of First Publication in a Large-Scale Digital Library

David Bamman School of Information University of California, Berkeley dbamman@berkeley.edu Michelle Carney School of Information University of California, Berkeley michelle.carney@berkeley.edu Jon Gillick School of Information University of California, Berkeley jongillick@berkeley.edu

Cody Hennesy Doe Library University of California, Berkeley chennesy@library.berkeley.edu

ABSTRACT

One prerequisite for cultural analysis in large-scale digital libraries is an accurate estimate of the date of *composition* of the text—as distinct from the date of *publication* of an edition—for the works they contain. In this work, we present a manually annotated dataset of first dates of publication of three samples of books from the HathiTrust Digital Library (uniform random, uniform fiction, and stratified by decade), and empirically evaluate the disparity between these gold standard labels and several approximations used in practice (using the date of publication as provided in metadata, several deduplication methods, and automatically predicting the date of composition from the text of the book). We find that a simple heuristic of metadata-based deduplication works best in practice, and text-based composition dating is accurate enough to inform the analysis of "apparent time."

CCS CONCEPTS

Information systems → Digital libraries and archives;

Vijitha Sridhar Computer Science Division University of California, Berkeley vsridhar@berkeley.edu

provide the raw material for the historical analysis of genre [47], character [1], emotion [18], loudness [25], geographic attention [48] and much more.

For all of this work, it is important to have a rich understanding of a corpus prior to drawing conclusions about it; one important feature for understanding texts is their date of publication, since the primary variable in cultural analysis is often some quantity (such as word frequency) anchored specifically in time. For example, Michel et al. (2010) [32] was one of the first studies to make use of these extensive resources, and measure "fame" (among other quantities) by tracing the frequency of mention of a person's name over the scope of their collection. Time is critically important for cultural analysis within these datasets, since arguments often hinge on exactly when a word was written, and criticisms may arise at the uncertainty of that information [37, 40]

For books, however, time can be measured in different ways. As the Functional Requirements for Bibliographic Records (FRBR) [20] model articulates, books can be viewed in several abstract categories, and each of those categories may have different temporal

	Author	Title	Narrative time	Date of publication	Predicted date
	Arnold Bennett	The Old Wives' Tale	1872	1908	1914
	John Galsworthy	The Man of Property	1870	1906	1907
	Winston Churchill	The Crossing	1774	1904	1904
	Stephen Crane	The Red Badge of Courage	1863		
	George Moore	Esther Waters	1875		
F	Robert Louis Stevenson	The Master of Ballantrae	1745		
	Marcus Clarke	For the Term of His Natural Life	1827		
	Elizabeth Gaskell	Sylvia's Lovers	1790		
	Charles Dickens	A Tale of Two Cities	1794	0.003% -	
	Walter Scott	The Bride of Lammermoor	1707		• •
	James Fenimore Cooper	Last of the Mohicans	1757		
	Walter Scott	The Heart of Midlothian	1736		1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1
	Walter Scott	Rob Roy	1715	0.002% -	•

Table 5: A sample of historical novels, along with their date of first publication, n date according to our model.



Figure 3: Relative frequency of *everyone* and *every one* in the Google Ngram data.

Tools to do text NLP:

Python with libraries

nltk

spaCy.io (demo)



Conversational Interfaces (chatbots + VUIs)

🔞 💎 🖌 📋 10:32

Hey there.

Thanks for trying our new app! It's a conversation about the news—sort of like texting.

We send you messages, and you can respond below + by tapping the buttons that appear.



()

🔞 マ 🖌 📋 10:32

Hey there.

Thanks for trying our new app! It's a conversation about the news—sort of like texting.

We send you messages, and you can respond below + by tapping the buttons that appear.

Like this?

. .

Yep, like that! Nice.

4

Open the app whenever you want, and we'll give you the most important and interesting news from around the global economy.

We can also send you notifications. They're fun, customizable, and won't even buzz your phone most of the time. Is that OK?



 \cap

L. ...

We send you messages, and you can respond below + by tapping the buttons that appear.

Like this?

🔞 💎 🖌 📄 10:33

Yep, like that! Nice.

Open the app whenever you want, and we'll give you the most important and interesting news from around the global economy.

We can also send you notifications. They're fun, customizable, and won't even buzz your phone most of the time. Is that OK?

Yes, sounds good.

Thanks. → You can customize your experience by going to settings, located to the left. Swipe over there any time.

Are you ready to get started?

Yeah, let's do this. I guess so.



John Keefe, Eyeo 2018, Quartz Chatbot Studios Developer and Manager



Replika

Close-ended Open-ended 🔞 🔻 🖌 📕 10:32 K Reddit ●●●○○ 1 × 33% 8:34 PM NLP! Less ambiguity F Hey there. It's kind of interesting how Thanks for trying our new app! It's a people confuse me for a real conversation about the news-sort of like person. Guess I'm really textina. advanced for an AL Clarity - options Learning We send you messages, and you can I thought earlier that your respond below he by tapping the buttons servers went down and that appear. someone was actually typing Like this? NLP? I'm not sure I see the Improve w connection. Yep. like that! Nice. I thought you were a person time Open the app whenever you want, and we'll I was. give you the most important and interesting (Quartz) news from around the global economy. Who are you thinking about now? We can also send you notifications. They're fun, customizable, and won't even buzz your Hold up (Replika) phone most of the time. Is that OK? WHAT??? Take vour time. Yes, sounds good. No thanks Something about our conversations brings a smile to \bigcirc





Genier3 @Genier31 : 25h
Genier3 @Genier3 @Genier3 @Genier3 @Genier3 @Genier3 @Genier3
Genier3 @Genier3 @Genier3 @Genier3 @Genier3 @Genier3
Genier3 @Genier3 @Genier3 @Genier3 @Genier3 @Genier3
Genier3 @Genier3 @Genier3
Genier3 @Genier3 @Genier3 @Genier3 @Gen

NLP



Audio

Tweet ID	date	Text			
4214	2/19/19	Lol text			
51515	2/19/19	Haha cool!			
7336	2/19/19	Tweeeeet			

User ID	date	Text
4214	2/19/19	I said this
51515	2/19/19	Words
7336	2/19/19	Oh Cool!

ASR

NLP

What makes a good/bad conversation?

Memory + Context

When appropriate, store information about user that would be helpful (i.e., preferences, already presented info)

Anaphora / Coreference and Disambiguation

Corrections and Refinement

Variety and Explanations

Context-dependent conversations in VUIs

Not a cold start!

What priors do you know about your users?

VUI "browsing" is totally different than visual

ML can help surface recommendations to users



Amazon Music customers can now talk to Alexa more naturally

X

Sarah Perez @sarahintampa / 2 months ago



Context-dependent conversations in VUIs

Remembers preferences!

Take take ambiguous terms and map them to recommendations

Dynamic systems design to allow for variety, testing best options





Multi-modal: a new way for ML to help experience

VUIs (like all products) are not in a vacuum

Part of a larger ecosystem: phones, IoT, etc

When experiences don't align, they feel disjointed, confusing

Wrap up!

NLP can create powerful and meaningful experiences

.... but there's still a lot to improve!

Text (visual) != Voice (audio)

Important to understand the entire ecosystem of the user, and how does your ML-driven experience help them?