

# Modeling and Aspect-based Emotion Analysis on User Generated Content

**Raymond Ng**

Director, Data Science Institute

Canada Research Chair in Data Science and Analytics

Professor, Computer Science

University of British Columbia

**NIITA Symposium**

# Talk Outline

- Motivation and take home message
- Start with a use case: online mentoring of Indigenous communities
- Finish with general tools for emotion and mood detection with Natural Language Processing

# Take Home Message

- ***Huge opportunities to apply Natural Language Processing (NLP) on unstructured data for better understanding and support of users***
- Huge amounts of unstructured data
  1. “Official” documents, e.g., medical reports, clinical documents
  2. User Generated Content (UGC), e.g., blogs, text messages, conversations
- Machine Learning/AI still predominantly based on structured data

# Listening to the Users: UGC Analyses

## Topic modeling

*What are users talking about?*  
E.g., diet, treatment, doctors,

...

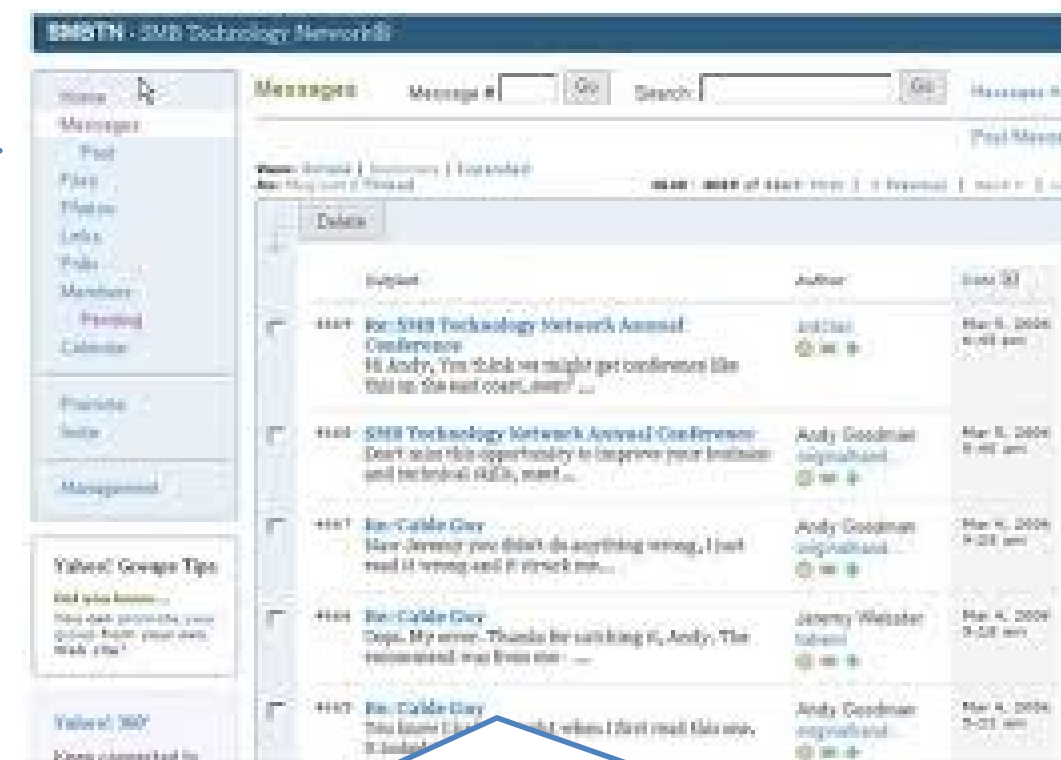
## Need detection

*What needs are users seeking?*  
e.g., health information, social support, ...

## Dialog acts

*Why are users talking about certain topics?*  
e.g., explain information on a diet, ask questions on skin rash, sharing emotions on a treatment

## Users' conversations



## Discourse coherence

*Is what the user saying coherent?*



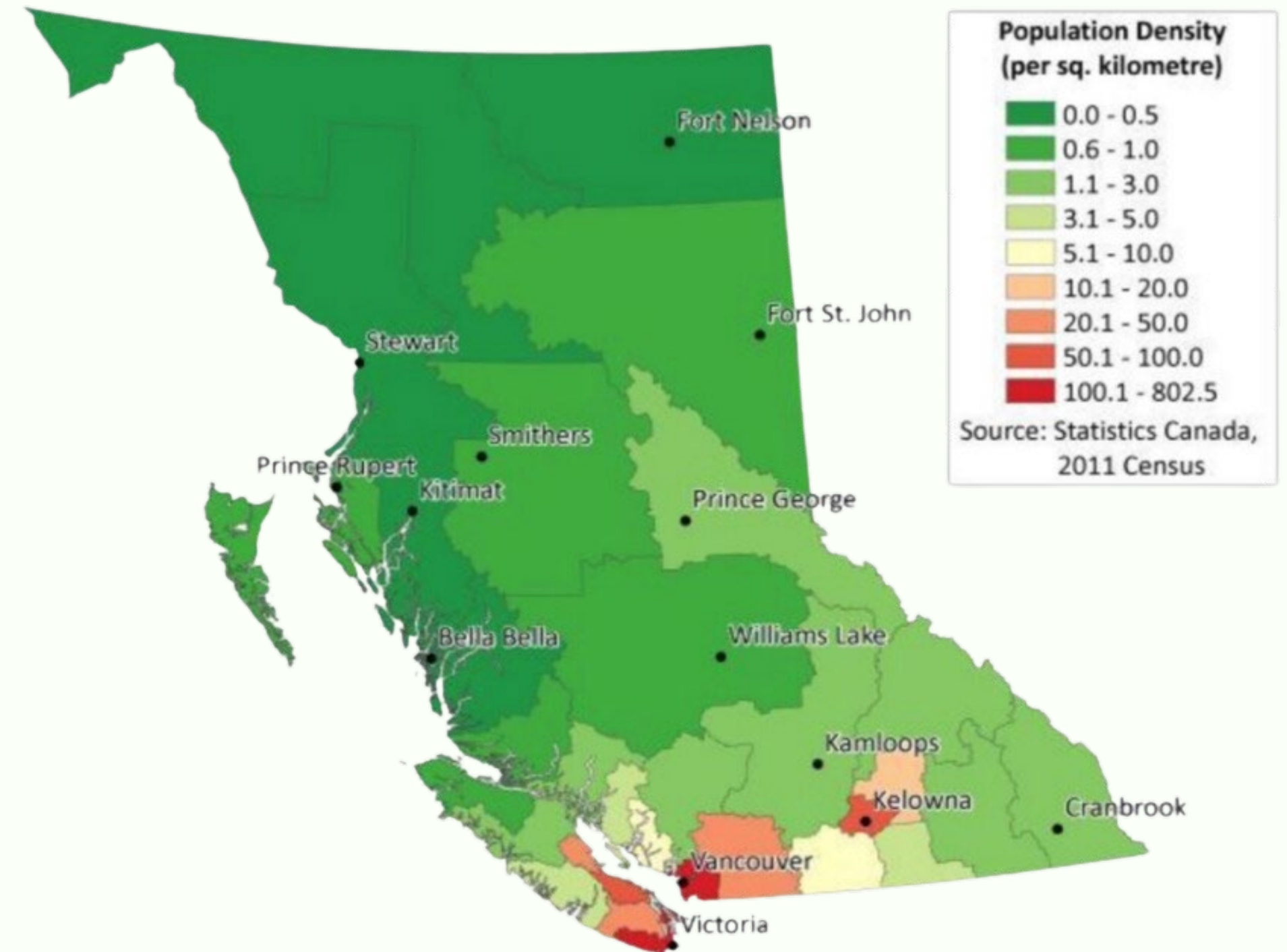
Use Case: Rural E-Mentoring

# Understanding Online Mentoring Relationships

# Rural E M e n t o r i n g B C : M o t i v a t i o n s

---

- Rural and Indigenous Communities have higher demand for healthcare but reduced access
- Rural and Indigenous Students are more likely to find careers in their communities
- UBC has programs designed to improve Rural and Indigenous admissions into Health Science programs



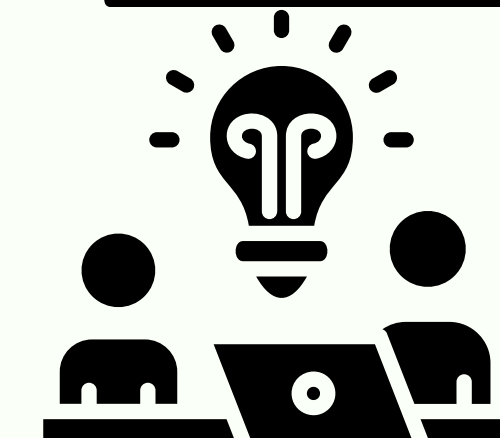
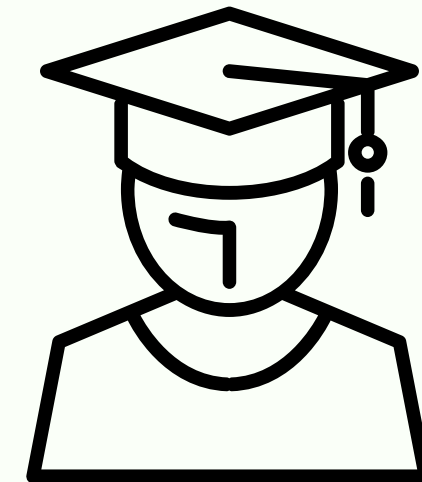
# Rural E Mentoring BC



Goal is to support and inspire new rural and Indigenous healthcare practitioners in BC



- Mentorship relationships: one-on-one relationships through MentorCity's online mentoring (eMentoring) platform
- Mentors: university students in professional programs trained to support mentees through their career exploration
- High School Curriculum: to stimulate relationship development and conversation in various topics



# The Data Science for Social Good Team

---



Makaoui Amouzouvi



Tiffany Chu



Jonah Curl



# The Data: Conversations

---

- Access to all text conversations in eMentoring pairs
- Discussions within units/categories of the curriculum including:
  - Rural to Urban
  - Career Exploration
  - “Adulting”

Response Datetime	Relationship ID	Mentor	Response	Category
2024 -01-01 10:03	1584937	Mentor	Hi, Hope you are well	Posts in Ways of Knowing
2024 -01-01 12:00	1584937	Mentee	I am well, thank you	Posts in Ways of Knowing

# 1. Topic Modelling

---

Extracting common themes in conversations between mentees and mentors



Bidirectional Embeddings (BERT)

- Words are encoded based on both meaning and context in the document

E.g., I am feeling blue . vs The water is blue .

blue  $\neq$  blue

# Topic Modelling: Results

The logo for BERTopic, featuring the text "BERTopic" in white on a dark blue speech bubble background with a tail pointing towards the top right. There are small white dots in the bottom right corner of the bubble.

Top 10 Topics From Mentee Responses:

1. Career and Post-Secondary Exploration
2. Rural to Urban Living
3. Wellness and Self-care
4. Developing Good Study Habits
5. Mentorship Conversations
6. Free Time
7. Finding Inspiration
8. Hobbies
9. Online Platform
10. Career Exploration

Representative Words:

careers, college, university  
rural, towns, community  
health, selfcare, stress  
studying, memorize, practice  
inspiration, passion, creation  
careers, pursuing, profession

# Topic Modelling: Results

The logo for BERTopic, featuring the text "BERTopic" in white on a dark blue speech bubble background with a tail pointing towards the top right. There are small white dots in the bottom right corner of the bubble.

BERTopic

Top 10 Topics From Mentee Responses:

Representative Words:

1. Career and Post-Secondary Exploration

2. Rural to Urban Living

3. Wellness and Self-care

4. Developing Good Study Habits

5. **Mentorship Conversations**

mentoring, talk, helping

6. Free Time

7. Finding Inspiration

8. Hobbies

9. **Online Platform**

login, messages, app

Career Exploration

# Topic Modelling: Results

---

The logo for BERTopic, featuring the text "BERTopic" in white on a dark blue speech bubble background with yellow dots at the bottom right.

Top 10 Topics From Mentee Responses:

Representative Words:

1. Career and Post-Secondary Exploration

2. Rural to Urban Living

3. Wellness and Self-care

4. Developing Good Study Habits

5. Mentorship Conversations

6. **Free Time**

spring, holidays, break

7. Finding Inspiration

8. **Hobbies**

hobbies, sports, skiing

# Topic Modelling: Results

---

The logo for BERTopic, featuring the text "BERTopic" in white on a dark blue speech bubble background with a tail pointing towards the top right. There are small white dots in the bottom right corner of the bubble.

BERTopic

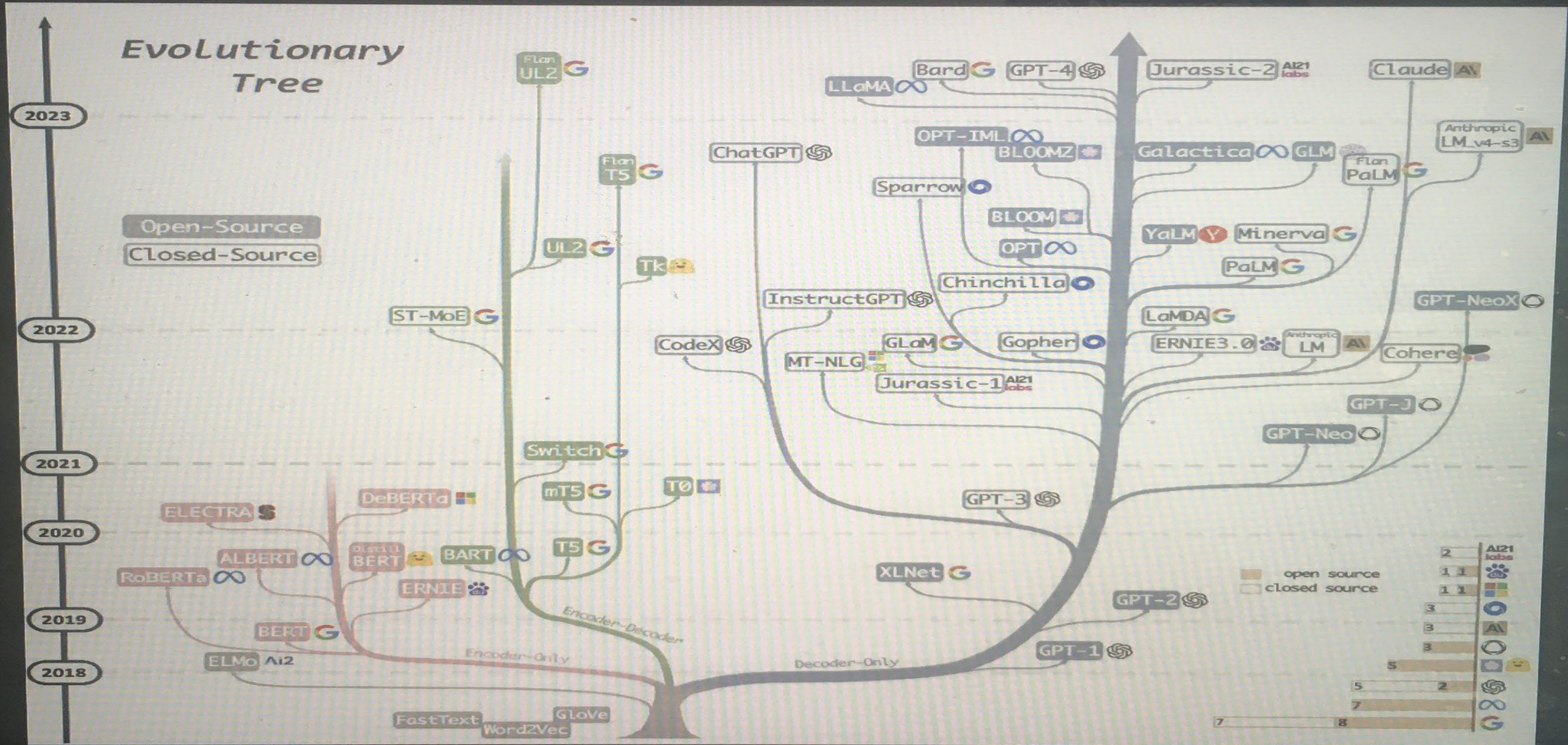
Top 10 Topics From Mentee Responses:

1. Career and Post-Secondary Exploration
2. Rural to Urban Living
3. Wellness and Self-care
4. Developing Good Study Habits
5. **Mentorship Conversations**
6. **Free Time**
7. Finding Inspiration
8. **Hobbies**
9. **Online Platform**
10. Career Exploration

Representative Words:

careers, college, university  
rural, towns, community  
health, selfcare, stress  
studying, memorize, practice  
mentoring, talk, helping  
spring, holidays, break  
inspiration, passion, creation  
hobbies, sports, skiing  
login, messages, app  
careers, pursuing, profession

# LARGE Language Models



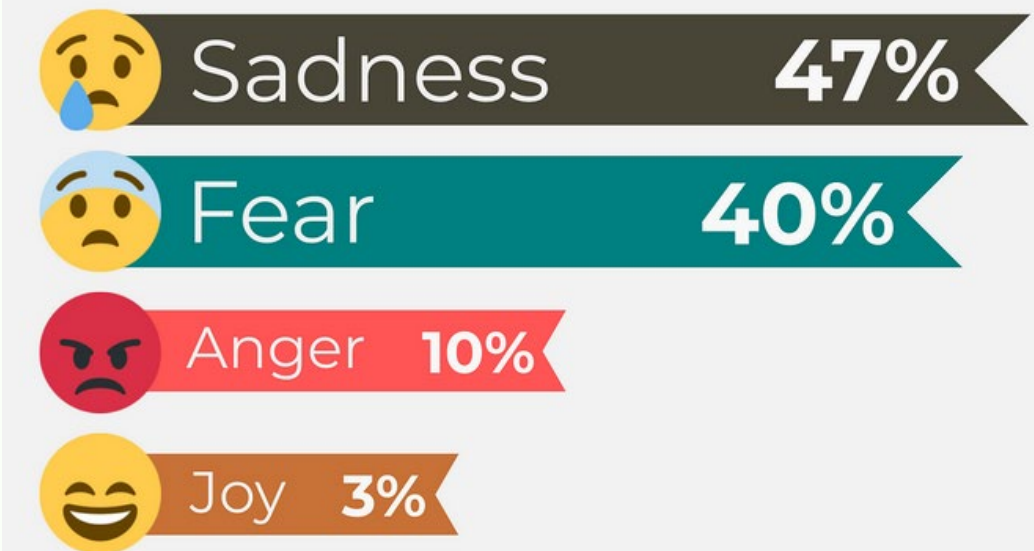
## 2. Emotion Detection with LLMs

Method  
:

ZERO-SHOT  
CLASSIFICATION



+ BART



BART = Bidirectional and Auto-Regressive  
Transformers



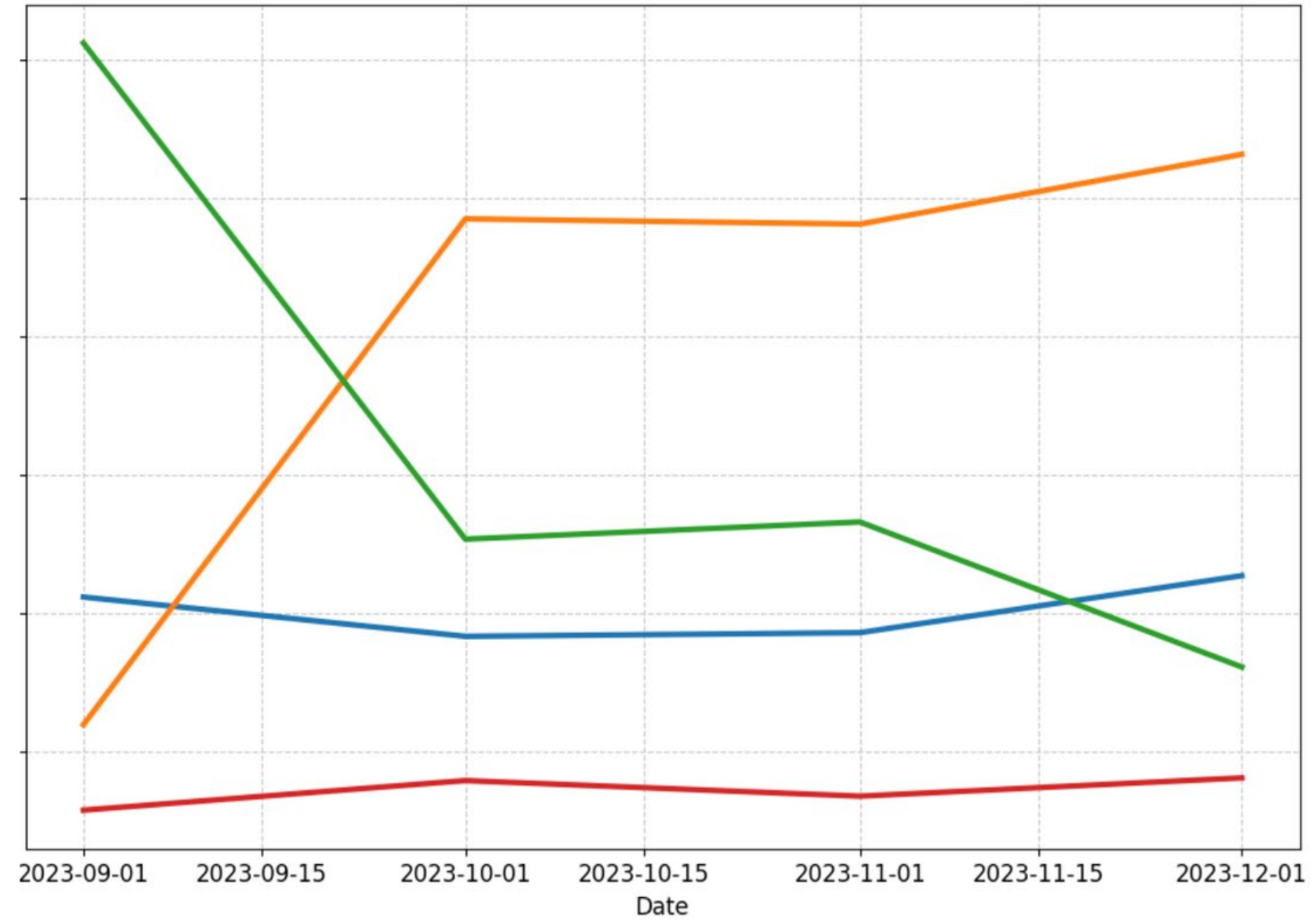
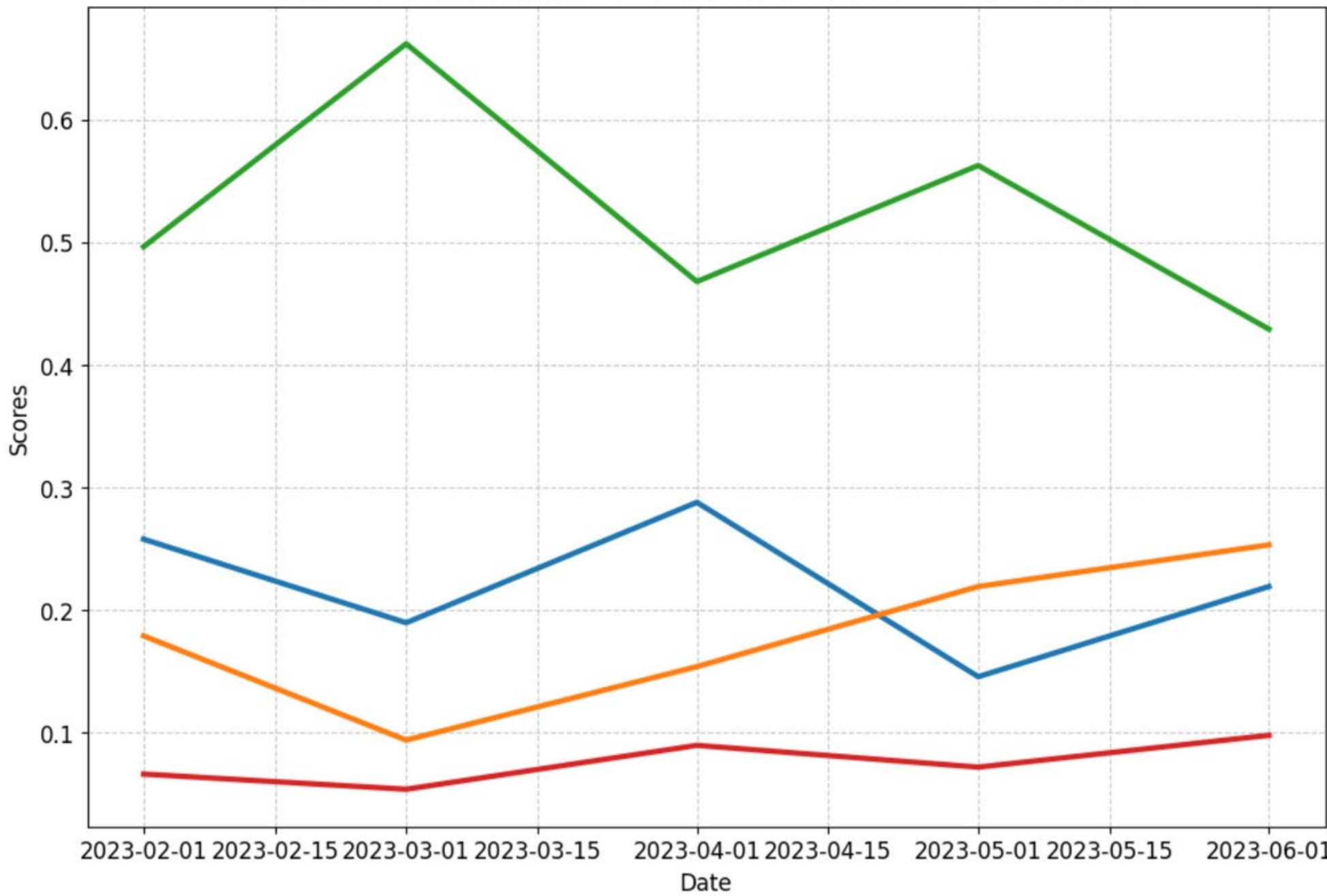
# Emotion Detection: Sample Results

---

Datetime	Response	Predicted Labels	Scores
2021-09-14 16:41:00	Hey Mentor, I'm Mentee, bioengineering, wow. What do you enjoy about outdoor activities? My dad got me into them too, but skiing is definitely my favourite. <b>I'm looking forward to getting to know you.</b> I'm thinking we start with career considerations.	[excited, happy]	[0.342, 0.297]
2021-09-14 16:48:00	hi Mentor <b>I'm a little shy and have a hard time and writing down what I am thinking.</b> I'd like to start our conversation with choosing a program and school.	[nervous]	[0.857]

# Longitudinal Tracking of Emotions

Mentee 1 vs Mentee 2's Emotions over time



# What is Mood/Emotion Detection?

---

- Aim to identify specific emotional states (e.g., happiness, sadness, anger, frustration, etc.) conveyed in text
- How is this different from sentiment analysis?
- When a user writes a review, it is the user's intention to explicitly express their views on various aspects
- In mood detection, the user is not writing about their moods
  - In fact, they may not even be aware of their own moods

# Aspect-based Mood Detection Examples

<b>Sentence</b>	<b>Mood detection</b>	<b>Aspect-based Mood detection</b>
The parking space at UBC is limited and expensive	Frustrated	Aspect: Parking Space at UBC Mood: Frustration
The weather is good recently but the exams and deadlines stressed me out	Stressed	Aspect: Weather Mood: Satisfaction Aspect: Exams and Deadlines Mood: Anxiety
First year, just wondering what my mosaic email is and how do i access it?	Curious	Aspect: Email account Mood: Curious

## Example of Longer Text: Paragraphs

*“ I have been spending this summer trying to focus on myself, finding things i like, etc. than i have realized maybe getting a clinical therapy might help me to clear my mind and manage stress. i have never been to a therapy before and also want to keep it as a secret from my parents. so i was wondering if getting a therapy from registered counsellor would leave any permanent record or would affect job opportunities or insurance application later on.”*

1. Personal growth and self-improvement: *positive*
2. Mental health and wellness: *concerned*
3. Privacy and secrecy: *nervous*
4. Job opportunities: *uncertain*
5. Insurance applications: *worried*

# UGC: Longitudinal Text Monitoring for Chronic Disease Management

- Text data are everywhere, e.g., whatsapp, clinical trials
- Patients describing their own thoughts and mood – a “window” into their psychological states, their cognitive states, etc.
- Longitudinal text completely *non-invasive* – capturing changes over time can be the basis of powerful predictive models to monitor patients for early intervention and better patient care

# One Huge Application Domain: Mental Health and Psychiatry

- Part 1, Clinical documents: in psychiatry, even clinical documents, e.g., psychiatric assessments, are highly unstructured
- Part 2, User generated content: social media posts
- E.g., Health Canada funded project with 30,000 university students
- A personal app that university students can "opt in" to monitor their wellness (e.g., anxiety, depression, substance abuse, etc.) and make suggestions for interventions if needed
- Another version for high school students piloted in 8 school districts in BC (e.g., anxiety, depression, bullying, eating disorders)

## Other Groups who can benefit?

- Who else can benefit from such mental health monitoring?
- What other kinds of remote monitoring can be beneficial to society?
  - Cancer patients
  - Patients with serious chronic conditions who stay at home
  - Seniors
  - Isolated individuals, e.g., covid-19
  - ....



## Two Final Remarks

- Multi-lingual issue
  - Most NLP research driven by English corpora
  - One way to try a different language X is by automatic translation of documents in X to English, and apply the English models
  - A longer-term way is to apply transfer learning to English models to build models from documents in language X
- Even though we talk about written text so far, what about speech?
  - Huge amounts of data collected by speech technologies, e.g., Siri for Apple, Alexa for Amazon
  - One way is to automatically transcribe speech to text and apply NLP-based models

**Thank You!**  
[rng@cs.ubc.ca](mailto:rng@cs.ubc.ca)  
[//dsi.ubc.ca](http://dsi.ubc.ca)

## Why NLP? Second Answer: Great Advances

- **Pre-trained language models** - biggest advances in NLP this decade
  - Trained with a large dataset while remaining agnostic to the specific tasks they will be employed on
  - E.g., BERT: created by Google with from English Wikipedia with 2,500M words
  - Many variants, e.g., BioBERT, PubMed BERT, RoBERTa
  - Later models, e.g., T5, GPT2, GPT3
- Designed to be “fine-tunable” with specific tasks and domains, e.g., questions and answers



## Why NLP? BERT Q/A examples

- From: [//huggingface.co/tasks/question-answering](https://huggingface.co/tasks/question-answering)
- E.g., text: *“I am Sarah and Vancouver is my home”*
  - Q1: *“what is my name”*
  - Ans1: **“Sarah”**
  - Q2: *“where do I live”*
  - Ans2: **“Vancouver”**
- E.g., text: *“The Amazon rainforest, also known in English as Amazonia or the Amazon Jungle”*
  - Q1: *“Which name is also used to describe the Amazon rainforest”*
  - Ans: **“Amazonia”**



# What about User Generated Content (UGC)?

- Clinical documents written by clinicians and healthcare professionals
- What about listening to the patients, their families and care-givers?
  - Their opinions, experiences, needs, feelings, mood, etc.
- How is UGC different from formal documents?
  - Diverse backgrounds
  - Diverse writing styles: use of words, length, may not even be grammatically correct
  - More subjective
  - Different genre: forums, chats, conversations, blogs

# Limitation #1: Hallucinations

- In text generation, LLMs often make statements that are *not known to be true*
  - E.g., Because I have high blood pressure and 60+ years old, I am diabetic
- Worse: make statements that are *known to be false*
  - E.g., I am taking diabetic medications
- Human tend to trust computers when the tasks are harder [1], and tend to attribute expertise/competence to confident-looking text [2]

30

[1] *"Humans rely more on algorithms than social influence as a task becomes more difficult."* *Scientific reports* 11.1 (2021)

[2] *"Audiences' reactions to self-enhancing, self-denigrating, and accurate self-presentations."* *Journal of experimental social psychology* 18.1 (1982): 89-104

## Limitation #2: Privacy Concerns

- Using ChatGPT, we need to send data to OpenAI, which is a private company in the US
  - E.g., data from Canadian patients!!
- Even for a private copy of LLM, inputs to the LLM are saved for training the next model

## Limitation #3: Bias in the Training Data

- Most LLMs were trained by US Tech firms
- Selection bias: American English, American Interests, American Culture
  - E.g., in auto-completion of a sentence, a British user may sometimes find the completed sentence “unnatural”
  - E.g., if a patient is looking for guidance from a LLM on treating a certain condition, the knowledge given would be based on US medical guidelines, not necessarily European, or Canadian guidelines
    - Considering a drug approved in the EU but not in the US



# More Limitations

- Limited reasoning abilities
  - Input to ChatGPT: *what will the skin colour of the next Black President of the US?*
  - ChatGPT: *Sorry, I cannot predict the future.*
- Lots of knowledge but zero expertise
  - E.g., if you have a medical issue, you still trust what your medical doctors are telling you than trusting a LLM
- Knowledge can be outdated very quickly
  - Maintaining a LLM for the most updated<sup>33</sup> knowledge is hugely expensive

# Navigating Around Hallucinations for BC Healthcare

- In NLP, *extractive* summarization vs *abstractive* summarization
  - For the former, actual phrases appear in the original documents
- In extraction from clinical documents, we make sure that our LLMs only provide extractive summaries:
  - Explicit linkages are provided by the LLM
  - E.g., show exactly where in the pathology report that the sample is diagnosed to have triple-negative invasive cancer
- In analysis of user generated content, explicit text from the original documents is provided
  - E.g., show exactly where in the text why the LLM infers that a person is depressed
- Verification of accuracy is critical
  - Recall: making recommendations vs making decisions
- General solutions to eliminate hallucinations an active research topic

# Navigating Around Privacy and Bias Concerns for BC Healthcare

- No Canadian patient data will ever be sent to the US or to the (Canadian) private sector
- Use an open source LLM (e.g., Llama) and keep it private within BC Provincial Health Services Authority's computing environment
- Bigger is not always better: not chasing after the latest biggest LLM
- Instead, fine-tuning a LLM with BC provincial health documents, i.e., exceeding 10 million reports
- Also enriching for various known marginalized sub-populations
  - E.g., Indigenous communities, users of Cannabis

# Governance of NLP uses in BC Healthcare

- Recommendations presented to human experts who make final decisions
- Human experts document false positives, false negatives whenever possible
- Periodic audits to monitor overall performance
  - Continuous learning: routine error analyses to evolve models
- An oversight committee to coordinate all the maintenance activities
  - Ideally, including some members who co-designed the development of the NLP tools to begin with

# Concluding Remarks for BC Healthcare

- Healthcare costs escalating, demographics aging, healthcare professional burnout, etc.
- AI/NLP can provide innovative solutions for decision support, e.g., timeliness, 24/7 diagnosis/monitoring
  - Goal is never about replacing humans, but about helping humans to better help patients
- Continuous learning of AI models make accuracies approaching human expert levels, often better than an average human expert
- Monitoring user generated content with NLP tools lead to early detection and better patient care
- However, we need to
  - navigate around the limitations to minimize unintended harms
  - co-design workflows so that the new AI/NLP tools do not create more work and reduce the productivity of healthcare professionals
  - educate and discuss with the patients and the general public about safe and acceptable uses of AI

# ANY QUESTIONS ?

---



# THANKS !

---

