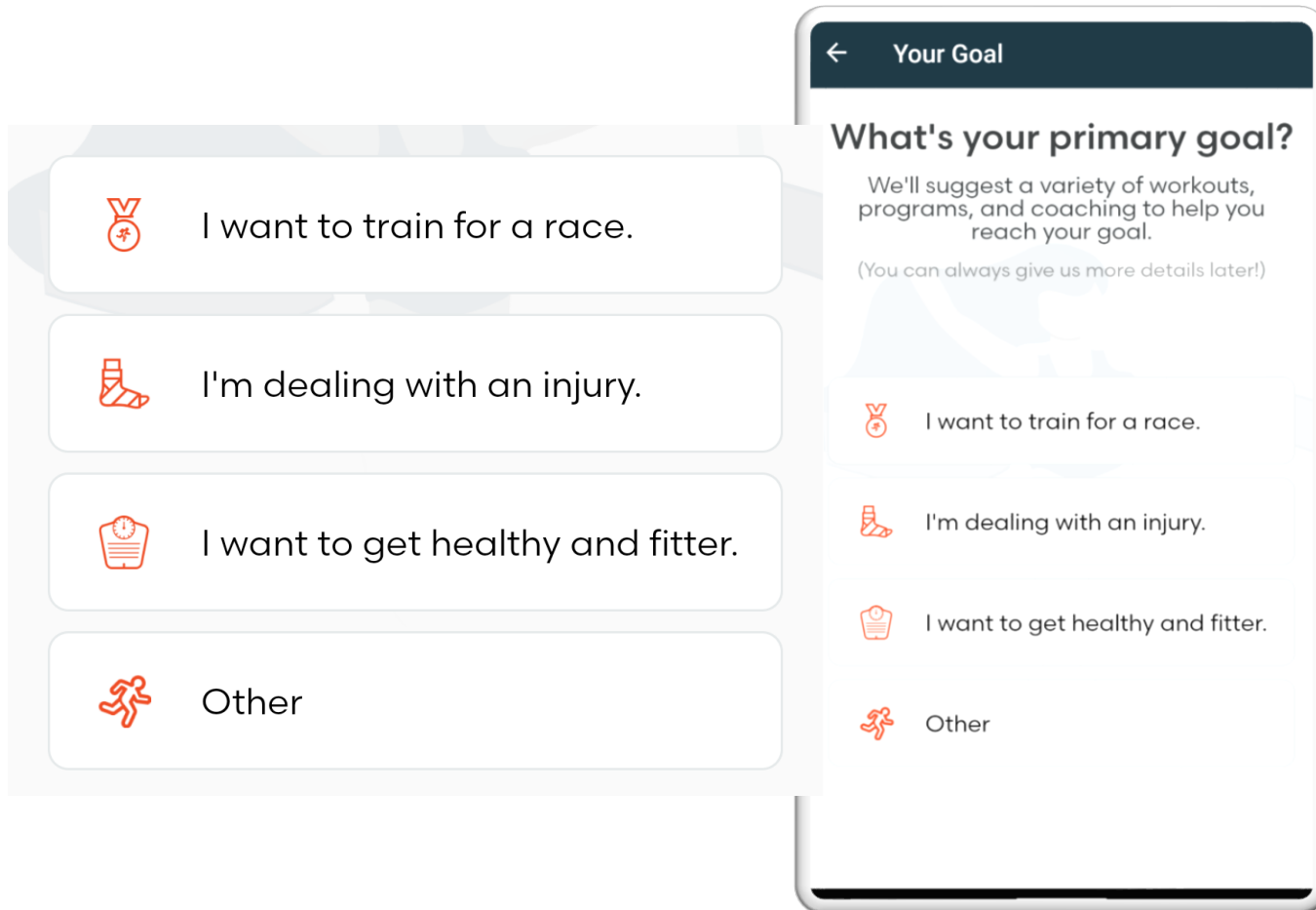# An app for running exercise

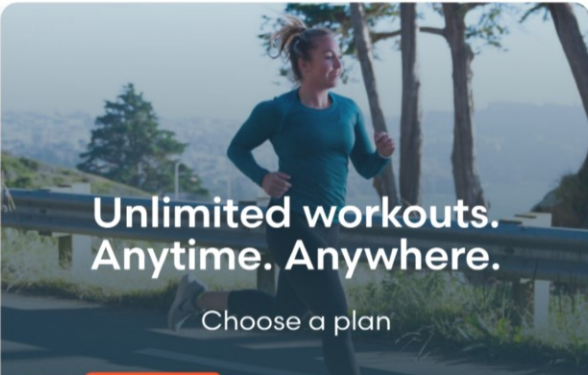# Imagine Eric uses this app for running exercise

# Imagine Eric uses this app for running exercise

# Will Eric subscribe after the free trial

# Imagine if you have user text data during the free trial

# What would you do as a data scientist

# What would you do as a data scientist

- Text:
  - Characters, words



Total characters: 150
Total words: 36

# What would you do as a data scientist

- Text:
  - Characters, words
  - Likes



> I finished my 30 day challenge ! I didn't know how to really video my push-ups and hollow body. But that's ok. I was really happy with my results. 30 minutes went from 3.04 miles to 3.41 👏💥😄💯

👍 12    💬 2

Comment

3 days ago

I finished my 30DC yesterday and (think?) I submitted my stats in Dropbox... I didn't see the caveat that they should be labeled with my name and which test it was 🙍 ahh, well 😁 my numbers did not improve for pushups

Home    Library    Commu...    Profile

Total characters: 150
Total words: 36
Total likes: 12

# What would you do as a data scientist

- Text:
  - Characters, words
  - Likes

- Sentiment:



I finished my 30 day challenge ! I didn't know how to really video my push-ups and hollow body. But that's ok.  I was really happy with my results. 30 minutes went from 3.04 miles to 3.41 👏💥😊💯

👍 12    💬 2

Comment

I finished my 30DC yesterday and (think?) I submitted my stats in Dropbox... I didn't see the caveat that they should be labeled with my name and which test it was 🧏 ahh, well 😁 my numbers did not improve for pushups

Total characters: 150
Total words: 36
Total likes: 12
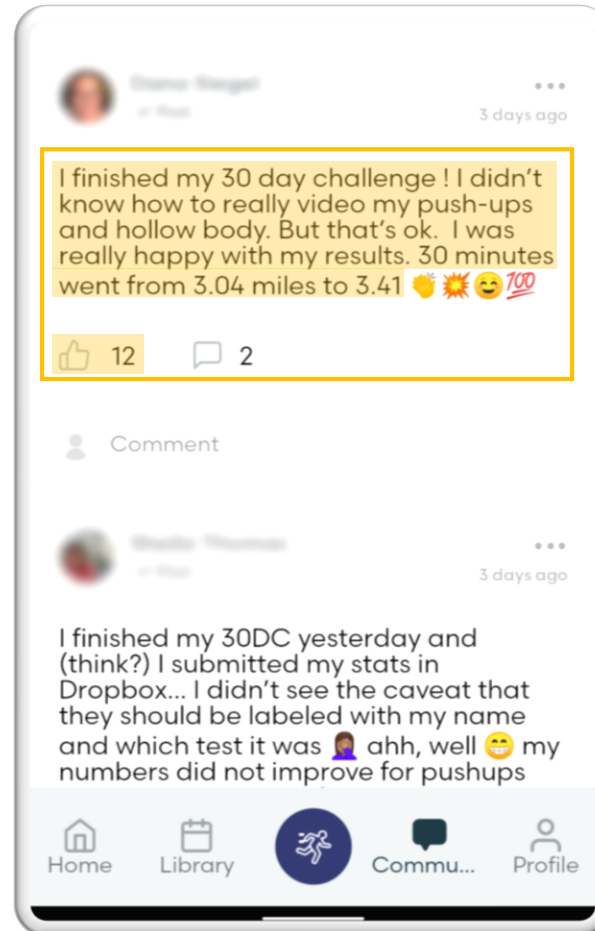
# What would you do as a data scientist

- Text:
  - Characters, words
  - Likes

- Sentiment:
  - Text and Emoji



Total characters: 150
Total words: 36
Total likes: 12

Emojis:
👏
💥
😊
💯

# What would you do as a data scientist

- Text:
  - Characters, words
  - Likes

- Sentiment:
  - Text and Emoji



Total characters: 150
Total words: 36
Total likes: 12

Emojis:
👏 >>> Clapping Hands
💥 >>> Collision
😊 >>> Smiley Face
💯 >>> Hundred Points

# What would you do as a data scientist

- Text:
  - Characters, words
  - Likes

- Sentiment:
  - Text and Emoji
  - Tone
  - Factual



Total characters: 150
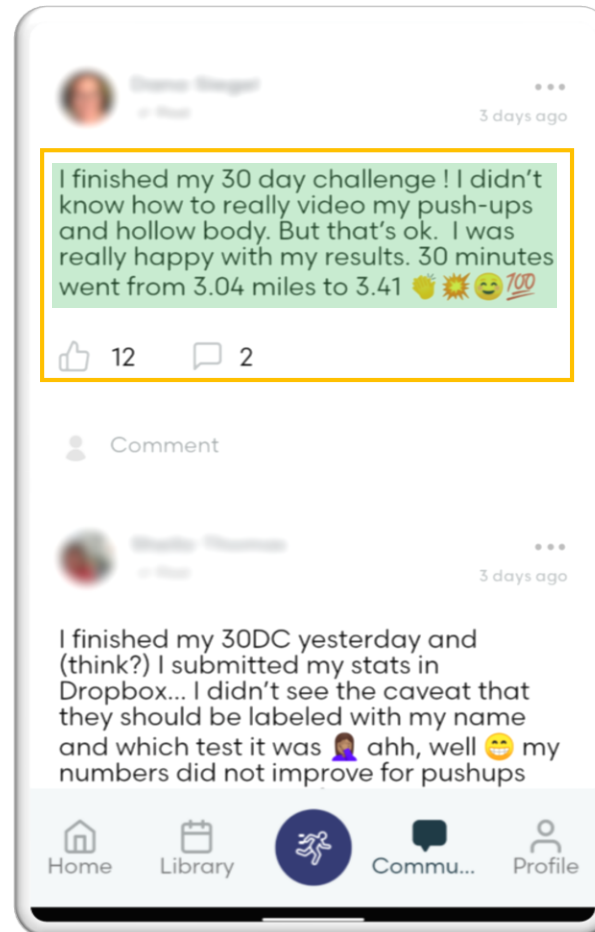Total words: 36
Total likes: 12

Emojis:
👏 >>> Clapping Hands
💥 >>> Collision
😊 >>> Smiley Face
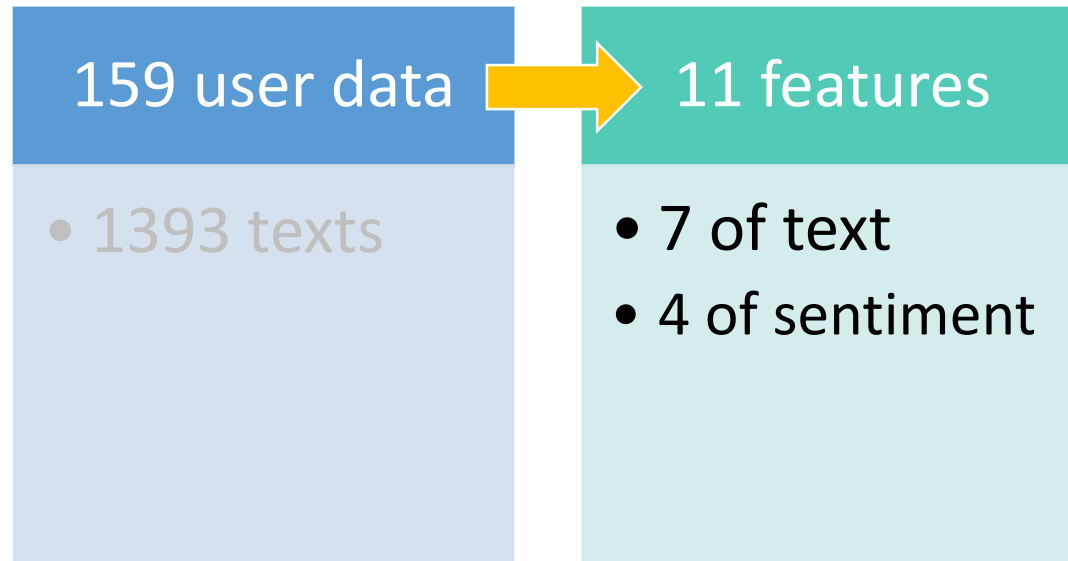💯 >>> Hundred Points

Sentiment:
Tone: Positive (1)
Factual: Rich (1)

# Processing data with machine learning models

**159 user data**

- 1393 texts
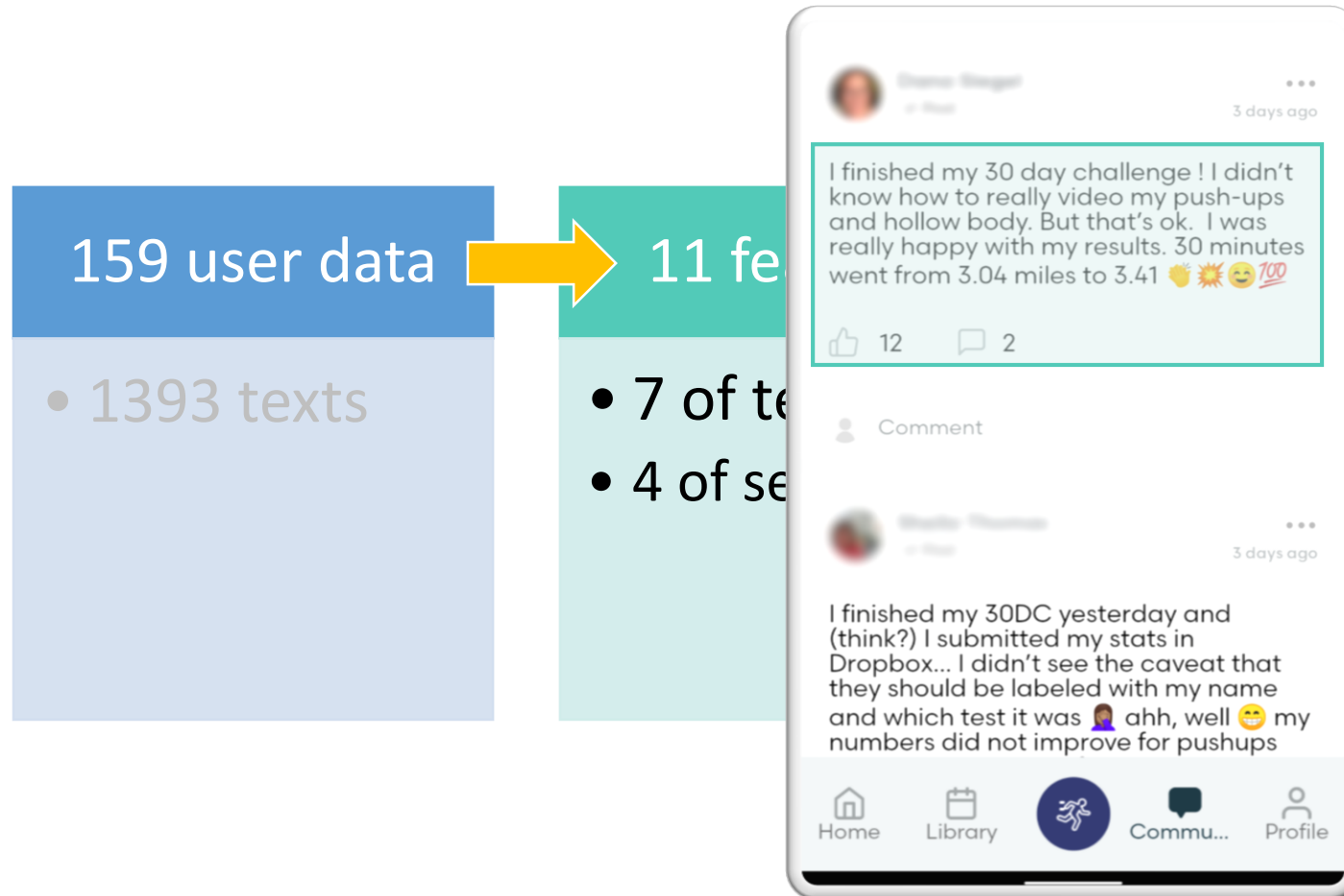
# Processing data with machine learning models

| 159 user data | 11 features |
|---|---|
| • 1393 texts | • 7 of text<br>• 4 of sentiment |

# Processing data with machine learning models

159 user data

• 1393 texts

11 fe

• 7 of te
• 4 of se

# Processing data with machine learning models

159 user data

→

11 fe...

• 1393 texts

• 7 of te...

• 4 of se...



I finished my 30 day challenge ! I didn't know how to really video my push-ups and hollow body. But that's ok. I was really happy with my results. 30 minutes went from 3.04 miles to 3.41 👏💥😄💯

👍 12    💬 2

Comment

3 days ago

I finished my 30DC yesterday and (think?) I submitted my stats in Dropbox... I didn't see the caveat that they should be labeled with my name and which test it was 🧑‍🦰 ahh, well 😁 my numbers did not improve for pushups

Home    Library    Commu...    Profile

[ **7 Text Features** ]

- Total characters
- Total words
- Average character
- Average words
- Number of posts
- Number of likes
- Average likes

[ **4 Sentiment features** ]

# Processing data with machine learning models



159 user data → 11 fe...

- 1393 texts

- 7 of te...
- 4 of se...

[ 7 Text Features ]
- Total characters
- Total words
- Average character
- Average words
- Number of posts
- Number of likes
- Average likes

[ 4 Sentiment features ]
- Total tone score
- Total factual score
- Average tone score
- Average factual score

# Processing data with machine learning models



159 user data → 11 features → Prediction

- 1393 texts
- 7 of text
- 4 of sentiment
- Subscription (Y / N)
- Binary classification

# Processing data with machine learning models

| 159 user data | 11 features | Prediction | Results |
|---|---|---|---|
| • 1393 texts | • 7 of text<br>• 4 of sentiment | • Subscription (Y / N)<br>• Binary classification | • Predicted 97% subscriptions |

# How do these models work

Sentiment analysis

- Natural Language Processing (NLP)

# How do these models work

Sentiment analysis

- Natural Language Processing (NLP)
- Pre-trained BERT
- Hand labelling 60% texts

# How do these models work

Sentiment analysis

• Natural Language Processing (NLP)

• Pre-trained BERT

• Hand labelling 60% texts

Classification

• Ridge, Logistic Regression, XGBoost, Random Forest

# How do these models work

Sentiment analysis

• Natural Language Processing (NLP)

• Pre-trained BERT

• Hand labelling 60% texts

Classification

• Ridge, Logistic Regression, XGBoost, Random Forest

• Stacking

# Validation metrics

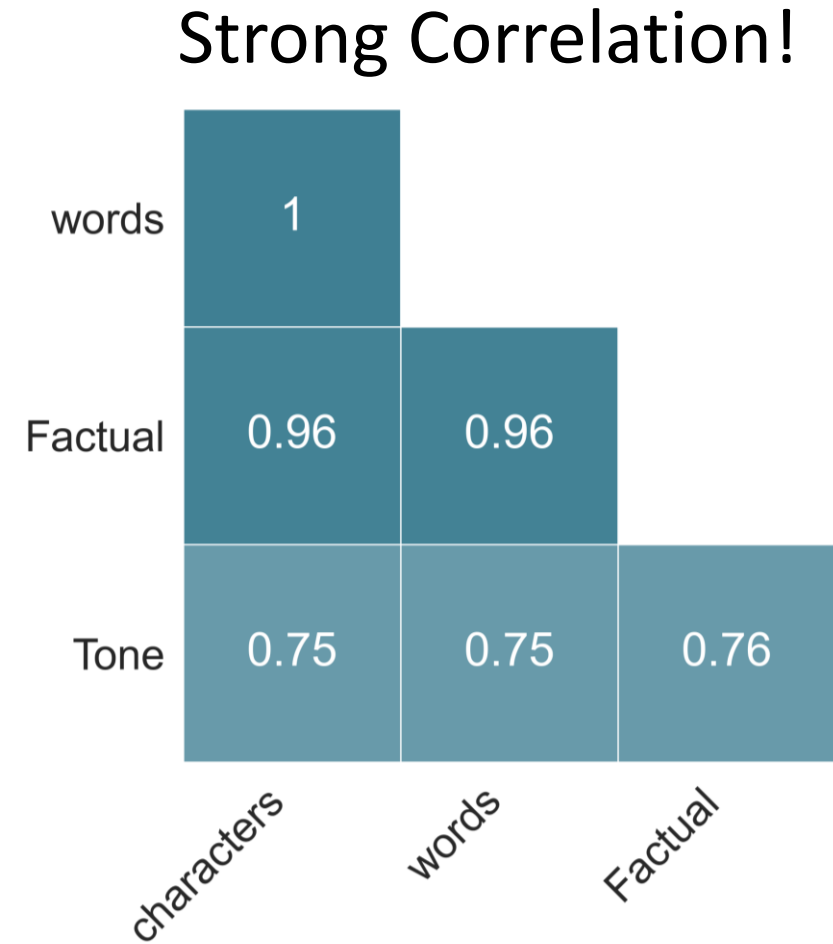| Sentiment | Accuracy |
|---|---|
| Tone | 0.851 |
| Factual | 0.776 |

| Classification | Precision | Recall |
|---|---|---|
| Regular classifiers | 0.56 ± 0.03 | 0.89 ± 0.04 |
| Stacked | 0.58 ± 0.03 | 0.97 ± 0.03 |

# Yes, the text data can predict user decision

# Yes, the text data can predict user decision

## and …

- Strong correlation
  - Text features
  - Sentiment features



Strong Correlation!

# Yes, the text data can predict user decision
## and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

## Text features

# Yes, the text data can predict user decision
## and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

## Text features
- Good enough?

# Yes, the text data can predict user decision
# and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

## Text features
  - Good enough?

| Classification | Precision | Recall |
|---|---|---|
| Text features | 0.57 ± 0.04 | 0.90 ± 0.02 |
| Sentiment features | 0.56 ± 0.04 | 0.92 ± 0.04 |
| | | |

# Yes, the text data can predict user decision
# and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

## Text features
- Good enough!

| Classification | Precision | Recall |
|:---|:---|:---|
| Text features | 0.57 ± 0.04 | 0.90 ± 0.02 |
| Sentiment features | 0.56 ± 0.04 | 0.92 ± 0.04 |
| All features | 0.56 ± 0.03 | 0.89 ± 0.04 |

# Yes, the text data can predict user decision
# and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

## Text features
- Good enough!
- Easy to scale up

| Classification | Precision | Recall |
|---|---|---|
| Text features | 0.57 ± 0.04 | 0.90 ± 0.02 |
| Sentiment features | 0.56 ± 0.04 | 0.92 ± 0.04 |
| All features | 0.56 ± 0.03 | 0.89 ± 0.04 |

# Yes, the text data can predict user decision and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

Text features
  - Good enough!
  - Easy to scale up

Takeaways

# Yes, the text data can predict user decision and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

Text features
  - Good enough!
  - Easy to scale up

Takeaways
  - User communication: strong indicator

# Yes, the text data can predict user decision and it gets even simpler!

- Strong correlation
  - Text features
  - Sentiment features

Text features
  - Good enough!
  - Easy to scale up

Takeaways
  - User communication: strong indicator

# Yes, the text data can predict user decision and it gets even simpler!

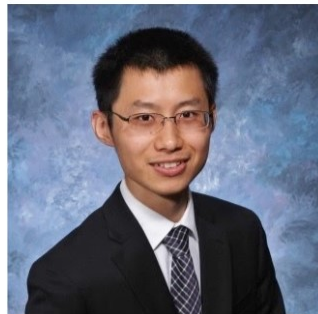- Strong correlation
  - Text features
  - Sentiment features

Text features
  - Good enough!
  - Easy to scale up

Takeaways
  - User communication: strong indicator
  - Recommend A/B testing on app features (causation)

# Zelong (Eric) Zhang



- PhD in Computational Chemistry
- Award-winning film (US DOE), photography
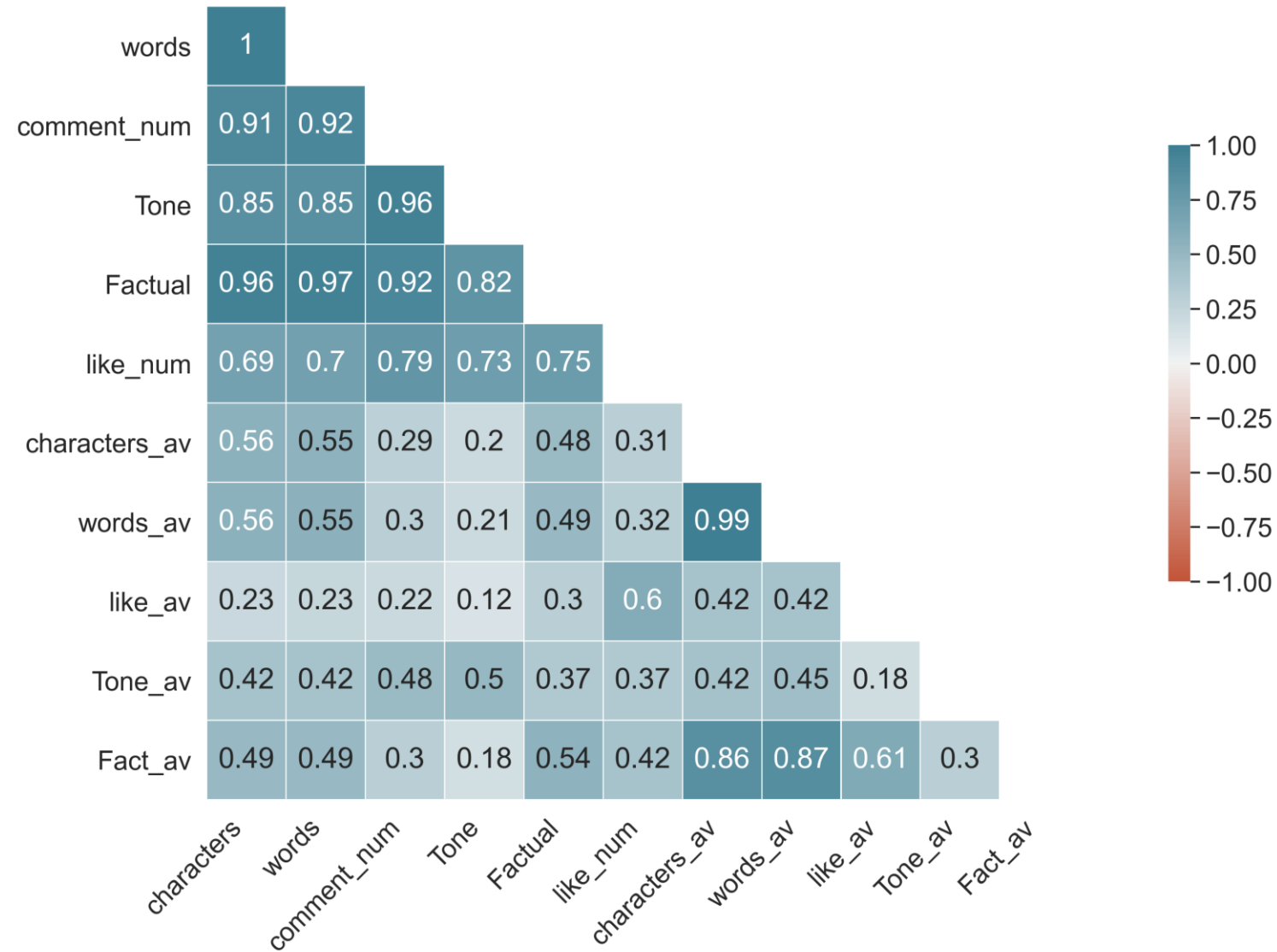- User Experience and Decision-Making

# Sanity check of sentiment analysis

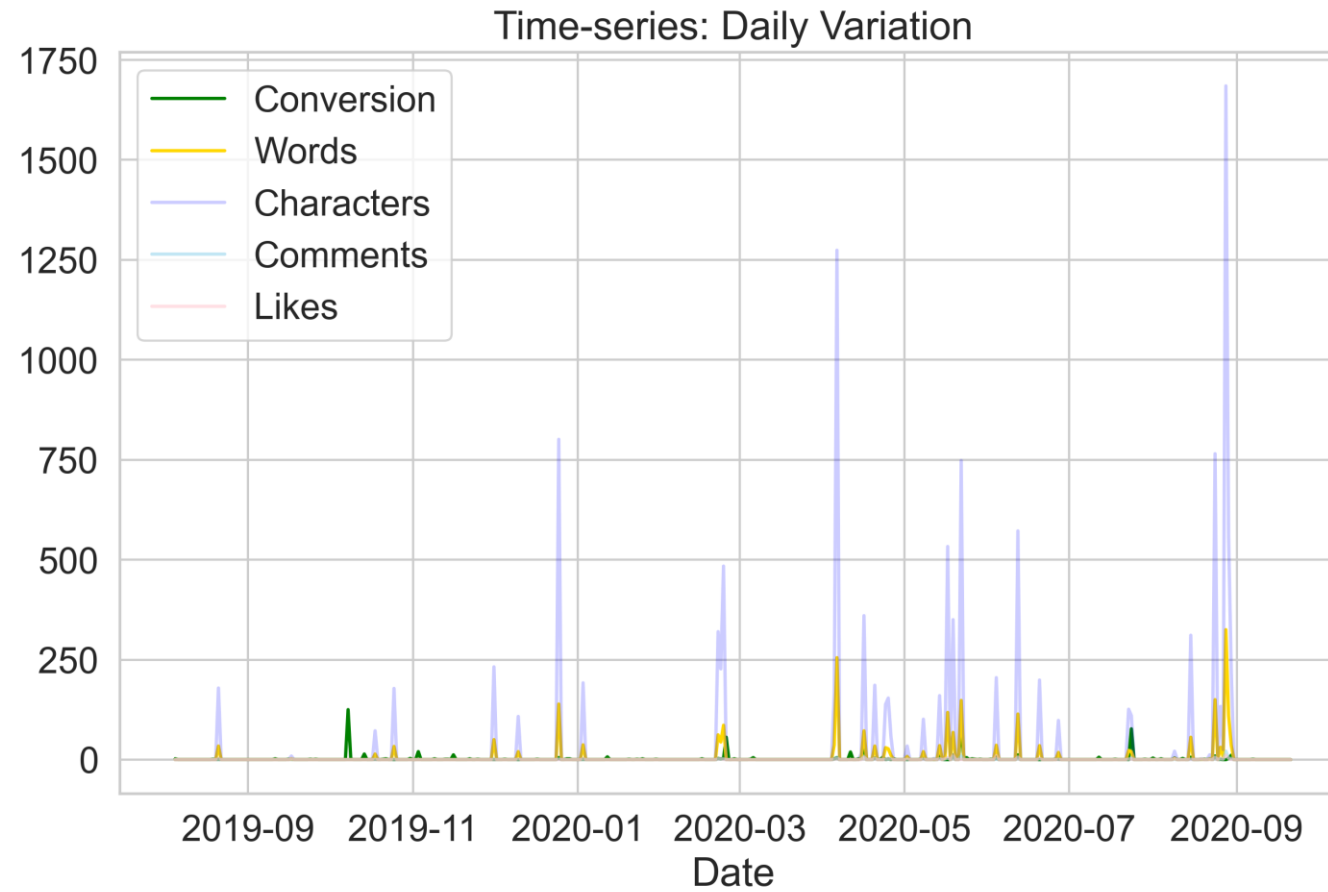| | text | Tone | Factual | VADER_Score | OTS_BERT_label | OTS_BERT_score |
|---|---|---|---|---|---|---|
| **1387** | Way to go, Michael! | positive | 0.0 | 0.0000 | positive | 0.9998 |
| **1388** | I like to torture myself!!!! 😜😜😜 | positive | 0.0 | -0.5526 | negative | 0.9992 |
| **1389** | yes!!! | positive | 0.0 | 0.5538 | positive | 0.9997 |
| **1390** | Oh dear that is swollen. Is ice helping? | neutral | 0.5 | 0.5859 | negative | 0.9983 |
| **1391** | Night run | neutral | 0.5 | 0.0000 | positive | 0.5209 |
| **1392** | Dabbling in swimming and biking. When my fitne... | positive | 0.5 | 0.3382 | negative | 0.9583 |

VADER: a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media.
OTS BERT: "distilbert-base-uncased-finetuned-sst-2-english"

# Strong Correlations between Features

# Time-series

# Time-series Analysis (ARIMAX)