



# Predicting Emoji Usage for a Recommender System



**Mariam Doliashvili**



**Lipyeow Lim**

**University of Hawaii**  **at Manoa**  
**Dec . 2017**



## Emojis, Emojis everywhere

### Problem Description

Emojitracker  
4 July 2013

emojitracker: realtime emoji use on twitter

👍	1056243340	❤️	580449867	❤️	471108721	😄	431631146	😄	345967697	😄	324085449	😄	313761984	😄	271272819
💖	268253888	😄	263924852	👍	251619883	😄	246904155	😄	216046535	😄	202974706	😄	198442656	🙏	133402319
👍	130212638	😄	127275356	👍	118696415	😄	118242683	😄	118127577	👍	110785355	🎵	109616142	👍	109265799
😄	103435721	👁️	99036466	😄	96831203	😄	90809972	😄	90783619	👍	90061405	😄	89821603	😄	88764875
💖	88592181	😄	86046163	💖	85626133	👍	84234882	😄	80873410	😄	78845347	🌱	76825366	💖	76426373
💖	77702794	😄	77361511	😄	75642628	👍	73614253	😄	73322516	💖	68324181	😄	67711119	👍	66586209
💖	66000108	💖	65685353	👍	65074643	💋	61218545	💖	59095859	👍	56502263	😄	53149554	👍	53114601
👍	50255918	😄	49645874	👍	49226507	👍	47786178	👍	47497919	👍	47339432	👍	46661267	😄	45711234
🌹	45183925	👍	43908499	👍	43165247	😄	41924782	✓	40957757	💖	40304641	👍	40098486	👍	39258281
😄	38869185	👍	38257495	👍	38077005	😄	38003656	😄	37266093	👍	37035146	💖	36911739	😄	33341665
💖	32161683	👍	31894311	😄	31766179	😄	30890414	👍	30607973	👍	30288921	😄	28674491	👍	28661904
👍	28522853	😄	28238791	👍	28016528	👍	27488797	👍	26262417	😄	26144526	👍	26135627	👍	26037534
😄	25560056	😄	24494728	👍	24025949	😄	23823858	👍	22935864	🎵	22923236	👍	22824188	👍	21582519
👍	20159625	👍	20055427	👍	19996579	💖	19792838	©	19721652	⚠️	19679234	👍	16896923	👍	18793877
😄	18352742	👍	18175625	👍	18074912	👍	17934671	👍	17096263	😄	16834180	👍	16603535	✓	16220522
👍	16091505	👍	15989889	👍	15890762	😄	15710011	👍	15674331	👍	15572185	👍	15448576	👍	15305245
👍	15094503	👍	14914755	👍	14497208	👍	14403633	😄	14384497	😄	14383954	👍	14300985	👍	14299446
👍	14207915	👍	14178182	✗	13860252	👍	13637184	👍	13332416	👍	13322356	👍	13152599	👍	12956665
👍	12932629	!	12535089	👍	12530768	♣️	12337776	👍	12219932	🌍	12152484	👍	12017688	👍	11842660
!!	11826795	✗	11620291	👍	11322781	👍	11277284	😄	11225264	💖	11011525	🌍	10962159	👍	10758004
👍	10674708	👍	10629815	👍	10273676	👍	10250701	👍	10241022	👍	10165692	👍	10129597	👍	10025090
😄	9779318	👍	9623333	👍	9505028	👍	9479433	👍	9478442	😄	9407887	👍	9396193	👍	9369400
👍	9338814	👍	9068402	👍	9061678	👍	8840198	👍	8831850	👍	8793013	😄	8717712	😄	8706435

# What can be solved?

*Emoji = Emotion ?*



## Outline

1. Use cases and complexity of the problem
- 

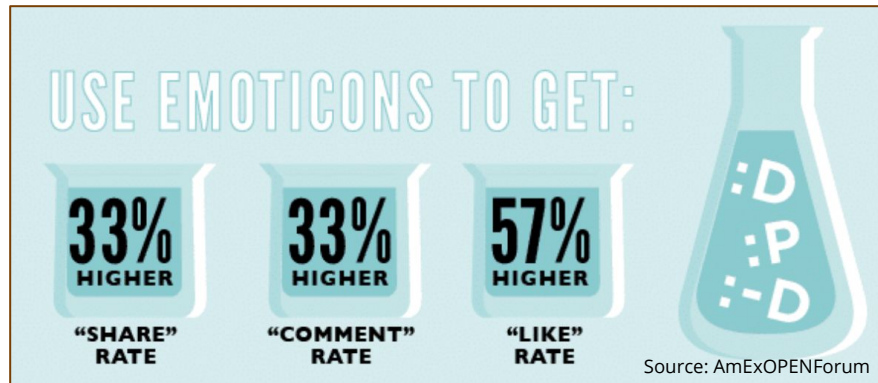
## Methods

2. Dataset & Preprocessing
  3. Multiclass vs Binary classification vs Personalization
  4. Models
  5. Evaluation
- 

6. Conclusion
7. Future work

# “ *Why should we use emojis?*”

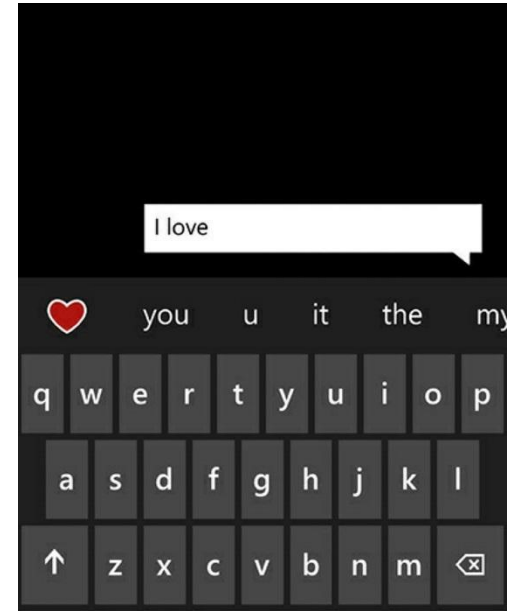
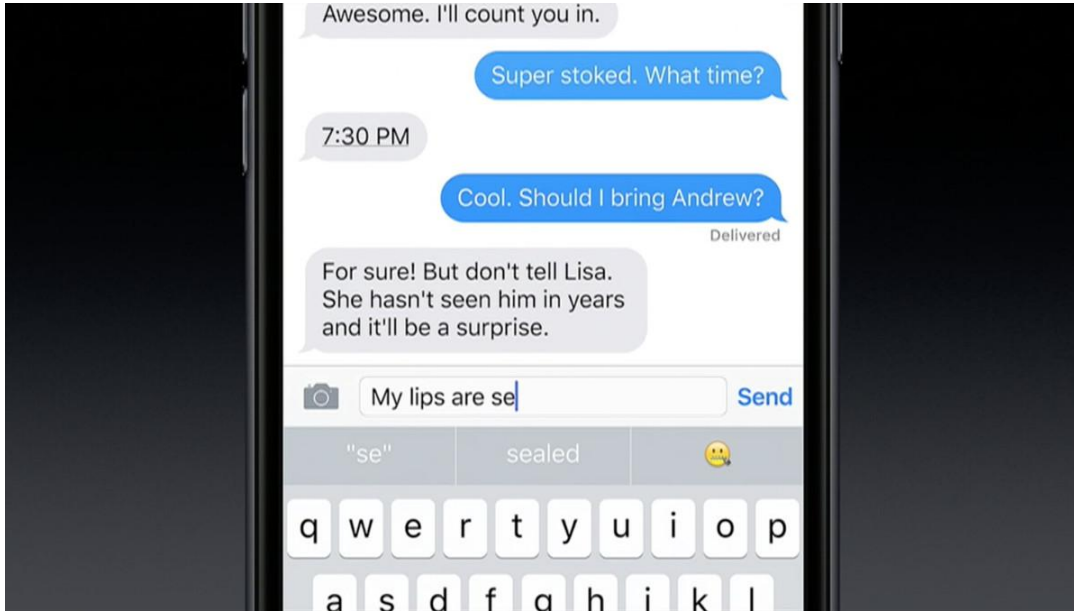
- *Makes it easier to express oneself*
- *Visual content is more likable*
- *Leads to better memorization*
- ...



1. Emojis are used as an emotion or a word in a sentence.

“

## *Use Cases*



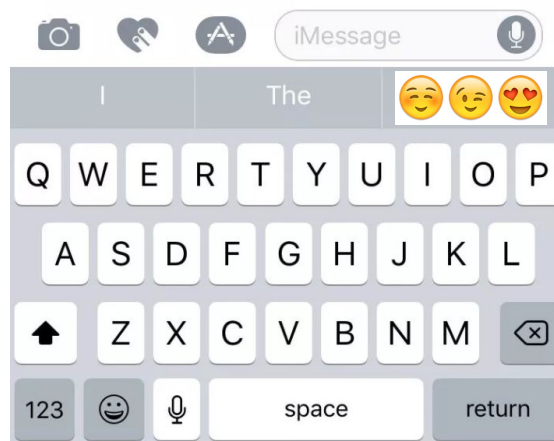




Now watch this.  
Read 8:40 AM

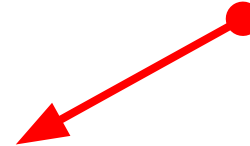


Now watch this.  
Read 8:40 AM



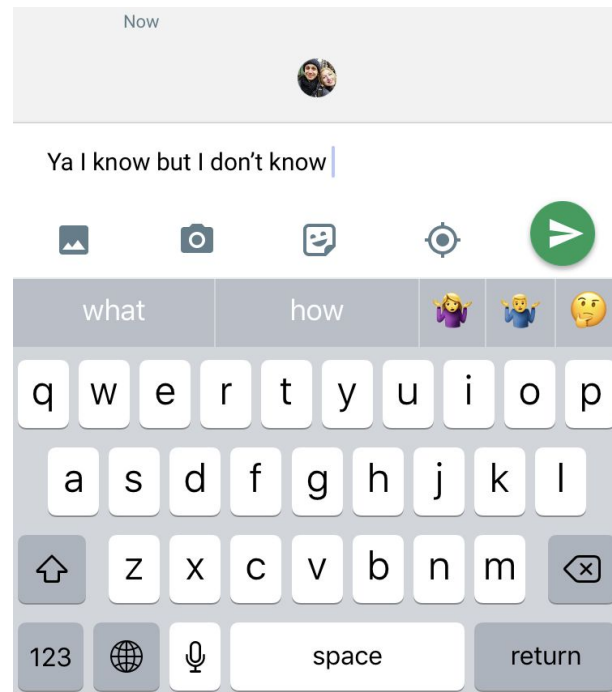
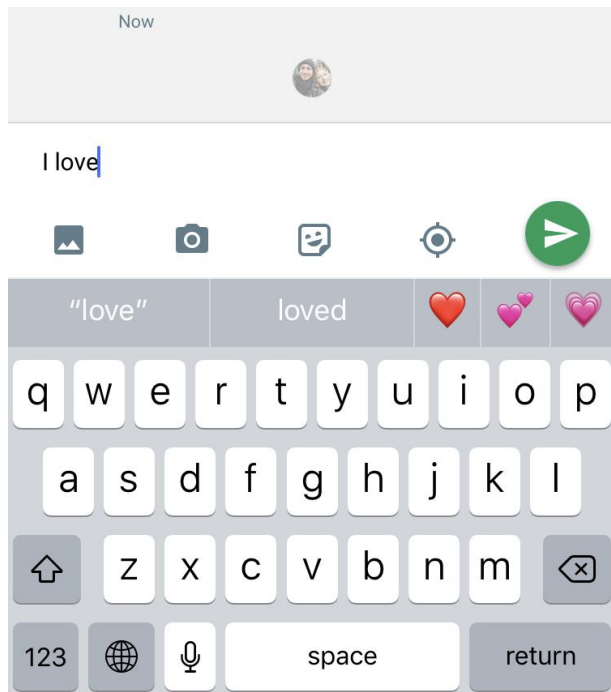
**The idea:**  
Given the text  
that an user  
wrote, return top  
k emojis relevant  
to the text.

**Goal**



Upgrade from  
recommending up to 3  
emojis and using n-grams.

“  
*Existing system iOS*



# Problem Description

*The original idea: Emoji Recommendations - Given the message/text that an user wrote, return top  $k$  emojis relevant to the text.*

*Problem:*

*What are the use cases?*

*Are emojis predictable?*

*Look at the tweets (with no emojis)*

*Guess what emoji describes the emotion of the text*



## Problem Description

*Try guessing?*



That's Alex!! I love her



- ?

## Problem Description

*Try guessing?*



That's Alex!! I love her



-



## Problem Description

*Try guessing?*



You're right man, I just  
can't get hurt like that  
again



- ?

## Problem Description

*Try guessing?*



You're right man, I just  
can't get hurt like that  
again



—



## Problem Description

*Try guessing?*



You're right man, I just  
can't get hurt like that  
again



—



*Laughing to lighten the sentence.*

Reverse Problem:

*What does this emoji mean?*



Reverse Problem:

*What does this emoji mean?*



*Neutral face*



*expressionless face*



Reverse Problem:

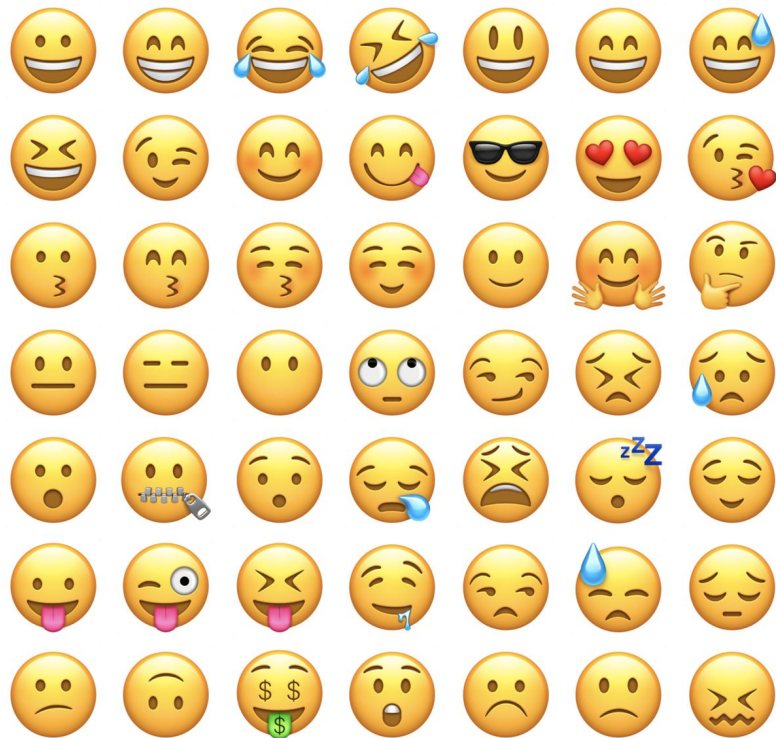
*What does this emoji mean?*

*Or these ones?*

- *We educate ourselves in emojis*
- *They are more ambiguous than words*

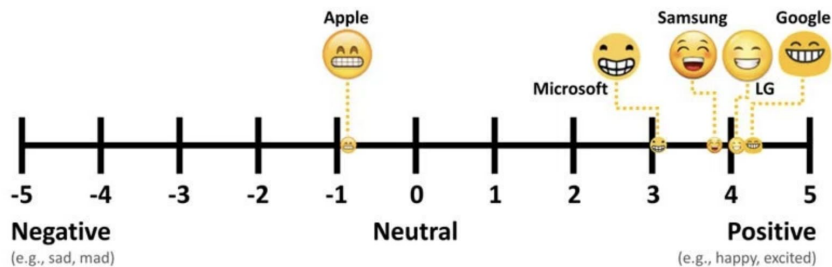


# 1 Emoji = 1 Emotion ?

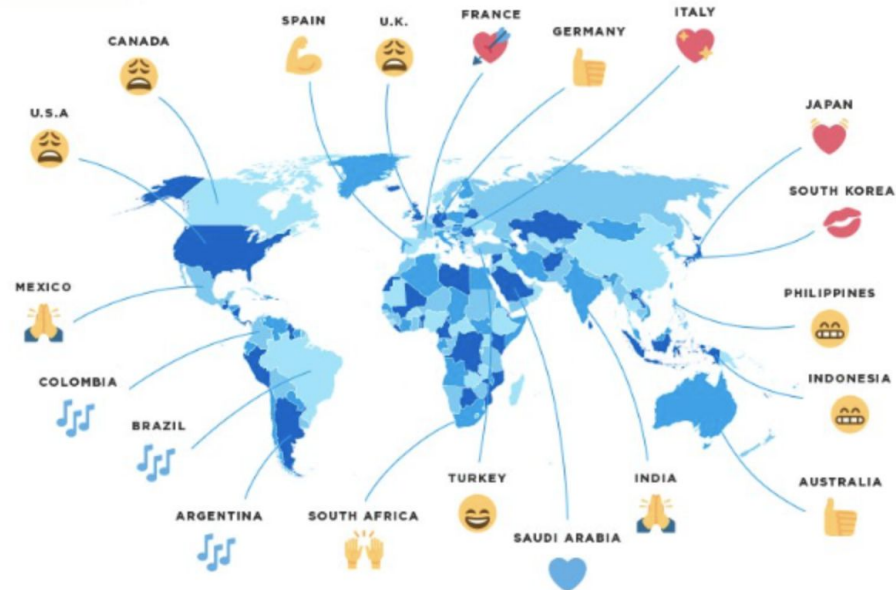




- Emoji M:N Emotion
- Emoji usage can be can depend on culture, nationality, gender, age, social circle ...
- Emoji representations vary for different platforms



## TOP-TWEETED EMOJIS BY COUNTRY



# Complexity

*There are 2753 in v11.0β*

- *It is hard to find an emoji that you need*
- *There can be 2753+ labels for classification problem*

*Emojis in a chat application can be used as:*

- An answer to the previous text if it was a question,
- A reply for the previous text,
- A next word in the sentence.

# Complexity

## *Emojis in any text can be used for:*

- Expressing an emotion or a word.
- As a combination to describe:
  - a) **a phrase** - “☕🕒?” - what time should we get coffee?  
“👏👏👏👏👏” - go warriors.
  - b) **the strength of emotion** -  
“😭😭😭” - very sad. “😂😂😂” - very funny.
  - c) **any visual that has no separate emoji in the unicode yet.**  
“Every damn corner 🌮🚚 #TacoTuesday”.  
“🐔🐣🐤🐥” - life cycle.  
🏢🏃 - city parkour.  
😭🔫 - kill me now



# 1.

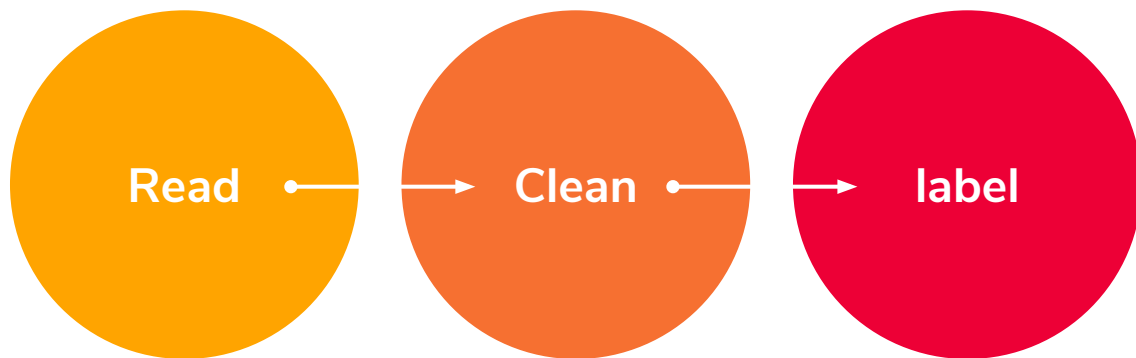
## Dataset & Preprocessing

*Getting Tweets that  
include emojis*



# Dataset

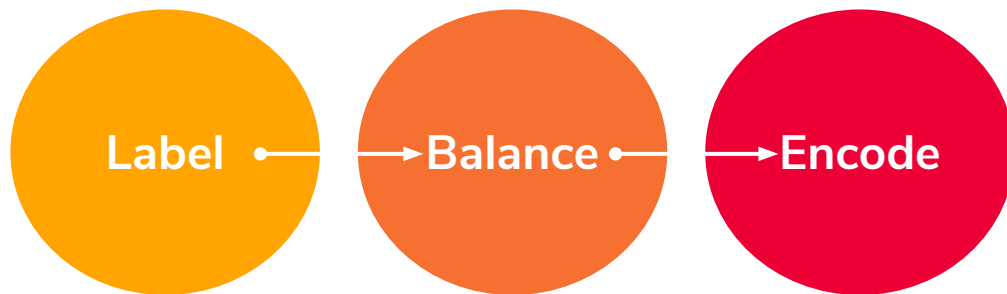
- Twitter Streaming API - 1,560,000 tweets
  - Filters applied: Language = EN,
  - Emoji is in “**My emoji set**” (74)
- Personalized dataset 120,000 tweets(2) + 20,000 tweets(10)
- Cleaning: URL-s, spelling, #-s, words that contain numbers...
- Handling tweets with several emojis and combojis.
  - Add combojis as a label





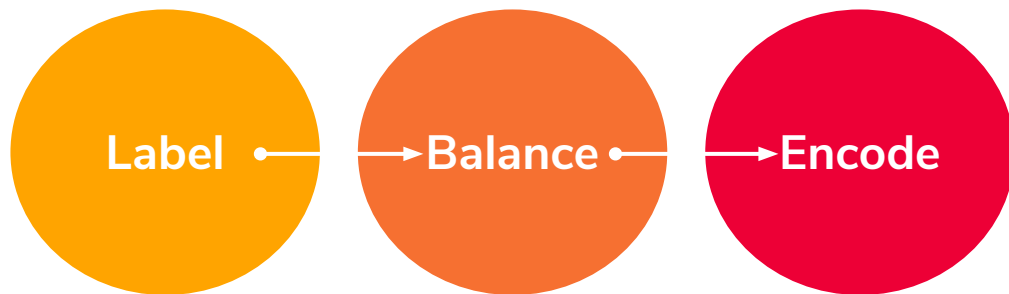
# Preprocessing

- Labeling:
  - The first emoji in the tweet
  - Yes / No
  - Separate entry for each emoji (not for combos)
- Balancing the data
- Feature extraction - frequency of emojis, device, daytime...
- Encoding:
  - One Hot Encoding
  - Word Embeddings



# Models

- MultiClass model\*: General dataset classification for up to 60 labels
- Binary model - LSTM:
  - Create a binary classifier for each emoji on balanced data (5)
- Personalized model - LSTM:
  - Feature extraction - frequency of emojis, device, daytime...
  - Labels are only previously used emojis and combojis
  - Next word prediction using n-grams and stupid backoff
  - Mapping words to emojis
  - Scoring based on confidence (emoji frequency, model accuracy, predicted probability)





# MultiClass vs Binary Classification vs Personalization

Change in Dataset / Change in Algorithm

# Multi Class vs Binary

## *Complexity vs. Usability*



### **Multi Class**

1. More complex problem.
2. More data but...
3. How many emojis to use?
4. Which ones to use?
5. Does not work good with not frequent emojis.



### **Binary**

1. Easier problem.
2. Introduces Bias & needs balancing.
3. Harder to detect with many choices of emojis with the same emotion.
4. Emoji choice is still limited

# vs Personalized

## *Complexity vs. Usability*



### **Personalized**

1. More complex problem.
2. Less data
3. More features
4. Less ambiguity
5. More emoji options
6. Suggesting new emojis (next word prediction)
7. Comboji support
8. ...

# Implementation



## NLTK

### **The preprocessing:**

Word tokenizer

### **Training:**

Naive Bayes Classifier



## Scikit

### **Training:**

Scikit classifiers:

Logistic Regression

Stochastic Gradient  
Descent



## Keras on TensorFlow

### **Preprocessing:**

Embedding Vector

Sequence

### **Training:**

LSTM neural network

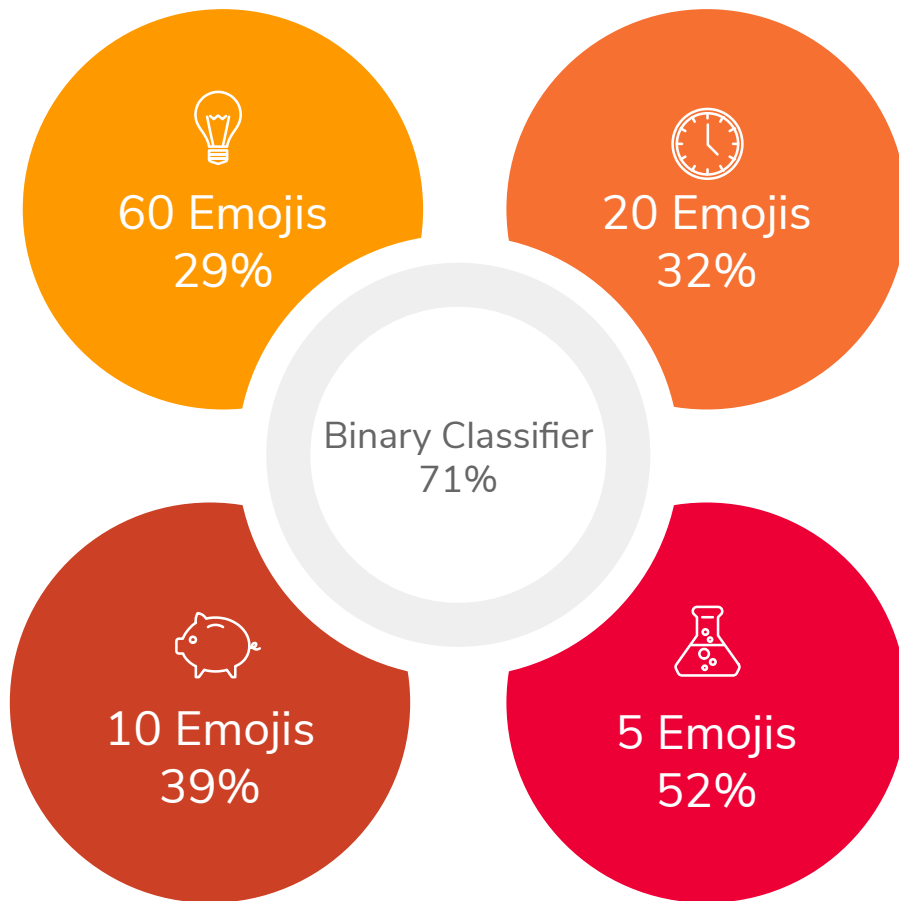
# Evaluation





# Evaluation Accuracy

Binary : 🤔  
365 000 tweets



# Evaluation Precision

Combined binary classifiers for 5 emoji P@1 - 73%

Barbieri et al. on Twitter dataset

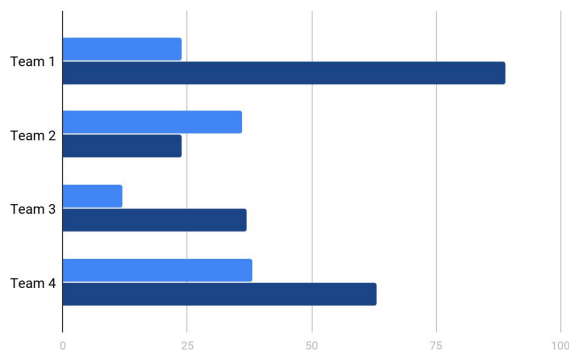
General B-LSTM classifier with top 5 emojis - p@1 - 65%

Human evaluation for the same was - 50%

Xie et al. on Weibo dataset

Hierarchical LSTM for understanding dialogue top 10 emojis

Points scored



Period 1 35%  
Period 2 65%

# Different methods

Results, results...

Accuracy	5	10	20+
Naive Bayes	48%	39%	32%
Logistic Regression	52%	38%	32%
Stochastic Gradient Descent	13%	9%	7%

# Evaluation

-

## Precision

Next word prediction - 13.5%

LSTM model with additional features on personalized data  
(generated labels on average - 125):

P@1 - 61%

P@3 - 74%

LSTM with Next word prediction and scoring:

P@1 - 40%

P@3 - 34%

Maybe because there was not a recommender system...

“Alternative” way of evaluating the recommender system



## Conclusion

Recommending emojis is a very complicated problem

It is fairly possible to break it down to subtasks -  
predict usage of one emoji.

Factors that help are:

Recommending several emojis from the same classifier

Using additional features about user

Using personalized data





# Future Work



## Future Work

### Dataset



Collect Facebook data for using more features

### Recruit people



Real world experiment

### Human Evaluation



Volunteers? ;)

# Resources

## *For the idea...*

- **Analyzing Twitter Sentiment of the 2016 Presidential Candidates** by Delenn Chin, Anna Zappone, Jessica Zhao
- **Emotional Chatting Machine: Emotional Conversation Generation with Internal and External Memory** by Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, Bing Liu
- **Are Emojis Predictable?** By Francesco Barbieri, Miguel Ballesteros, Horacio Saggion
- ... etc

Thanks to:



Dr. Lipyeow Lim and Dr. David Chin  
*for guidance*

Ed White and Gene Park  
*for donating their tweets*

Nicolas Golubovic, Muzamil Yahia, Mark Nelson and  
Takumi Aoki  
*for friendly support and reviewing my work*

That was all



It's freezing here



Thanks for listening



Do you



have any Questions?