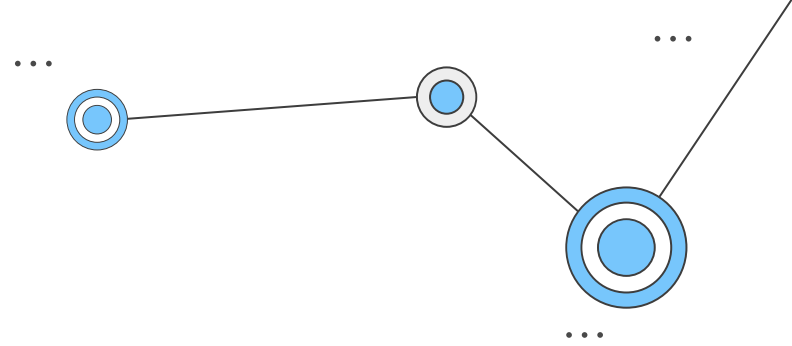


# Airline tweets sentiment analysis

Wirach (Joe)  
Leelakiatiwong



# Table of Contents

01

...

## The Problem

Why twitter? Why sentiment analysis?

02

...

## The Data

What data do we have?

03

...

## The approach

How can we flag negative sentiment?


04

...

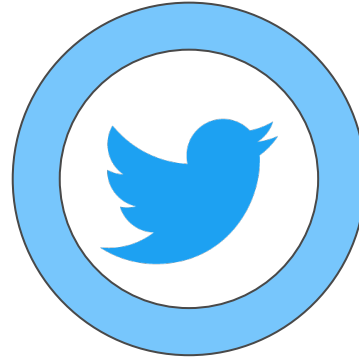
## The Conclusion

Model performance and further improvements





A little bit of  
background



# Twitter

**Twitter** is a '**microblogging**' system that allows you to **send** and **receive** short posts called **tweets**. Tweets can be up to **280 characters long** and can include links to relevant websites and resources.

...



**BBC Breaking News**

@BBCBreaking

France and Germany join the US in advising their nationals in Libya immediately [bbc.in/1rVmrDJ](https://bbc.in/1rVmrDJ)



RETWEETS  
596

FAVORITES  
223



**US Airways**

@USAirways

@ellerafter We welcome feedback, Elle. If your travel is complete, you can detail it here for review and follow-up:

[pic.twitter.com/vbeYgCuG25](https://pic.twitter.com/vbeYgCuG25)

Reply Retweet Favorite More



**Donald J. Trump**

@realDonaldTrump

Follow

Sorry losers and haters, but my I.Q. is one of the highest -and you all know it! Please don't feel so



**Jim B.**

@thetruefailure

Follow

Hey @Delta, you suck a lot and I hate having to fly your terrible airline.

7:37 AM - 13 Feb 2019

2 Likes



2



2



**Delta** @Delta · Feb 13

Replying to @thetruefailure

That's not good to hear, Jim. Pls follow/DM if I may be of assistance. \*AAB

[Send a private message](#)



1

...  
...  
...  
@AmericanAir **Love** the new planes for the JFK-LAX run. Maybe one day I will be on one where the amenities all function. #NoCharge #Ever



@AmericanAir leaving over 20 minutes **Late Flight**. No warnings or communication until we were 15 minutes Late Flight. That's called **shitty** customer svc

# Sentiment Analysis

Sentiment analysis is a natural language processing technique used to **determine whether data is positive, negative or neutral** to help businesses monitor brand and product sentiment in customer feedback, and **understand customer needs**.

...



01  
The

Problem



"Although airline industry — one of the world most competitive market — have been employing multiple feedback channels to collect customer feedback, they are all arduous. Thus making a social media analysis a more compelling approach. This project will identify **negative sentiment tweets** about the airline to categorize **them into subcategories and direct them to related departments to solve problems** and develop improvement strategies in the future."







02

The Data



# Data: Twitter US Airline Sentiment dataset (Kaggle)



## Coverage

Data was collected  
from 16th to 24th  
February of 2015

...



## Airlines

Covered 6 major airlines  
in the USA  
(United, US Airways, American  
Airlines, Southwest, Delta,  
Virgin)



## Size

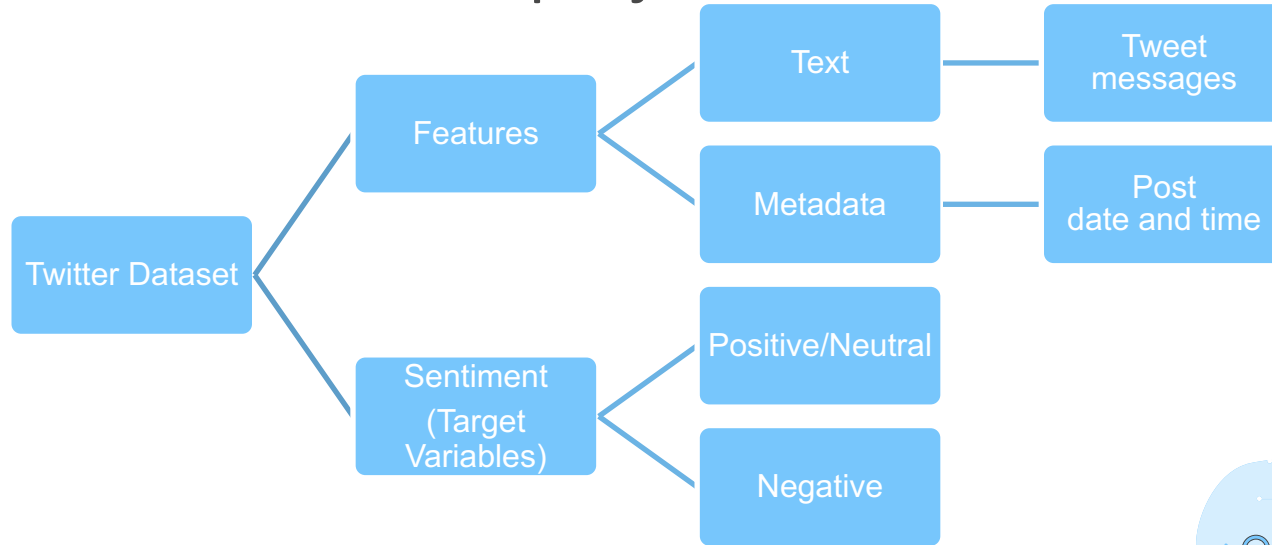
Rows : 14640  
Columns : 15

...



	tweet_id	airline_sentiment	airline	text	tweet_created
0	570306133677760513	neutral	Virgin America	@VirginAmerica What @dhepburn said.	2015-02-24 11:35:52-08:00
1	570301130888122368	positive	Virgin America	@VirginAmerica plus you've added commercials t...	2015-02-24 11:15:59-08:00
2	570301083672813571	neutral	Virgin America	@VirginAmerica I didn't today... Must mean I n...	2015-02-24 11:15:48-08:00
3	570301031407624196	negative	Virgin America	@VirginAmerica it's really aggressive to blast...	2015-02-24 11:15:36-08:00
4	570300817074462722	negative	Virgin America	@VirginAmerica and it's a really big bad thing...	2015-02-24 11:14:45-08:00
...	...	...	...	...	...
14635	569587686496825344	positive	American	@AmericanAir thank you we got on a different f...	2015-02-22 12:01:01-08:00
14636	569587371693355008	negative	American	@AmericanAir leaving over 20 minutes Late Flig...	2015-02-22 11:59:46-08:00
14637	569587242672398336	neutral	American	@AmericanAir Please bring American Airlines to...	2015-02-22 11:59:15-08:00
14638	569587188687634433	negative	American	@AmericanAir you have my money, you change my ...	2015-02-22 11:59:02-08:00
14639	569587140490866689	neutral	American	@AmericanAir we have 8 ppl so we need 2 know h...	2015-02-22 11:58:51-08:00

# Data Dictionary (in the scope of project)



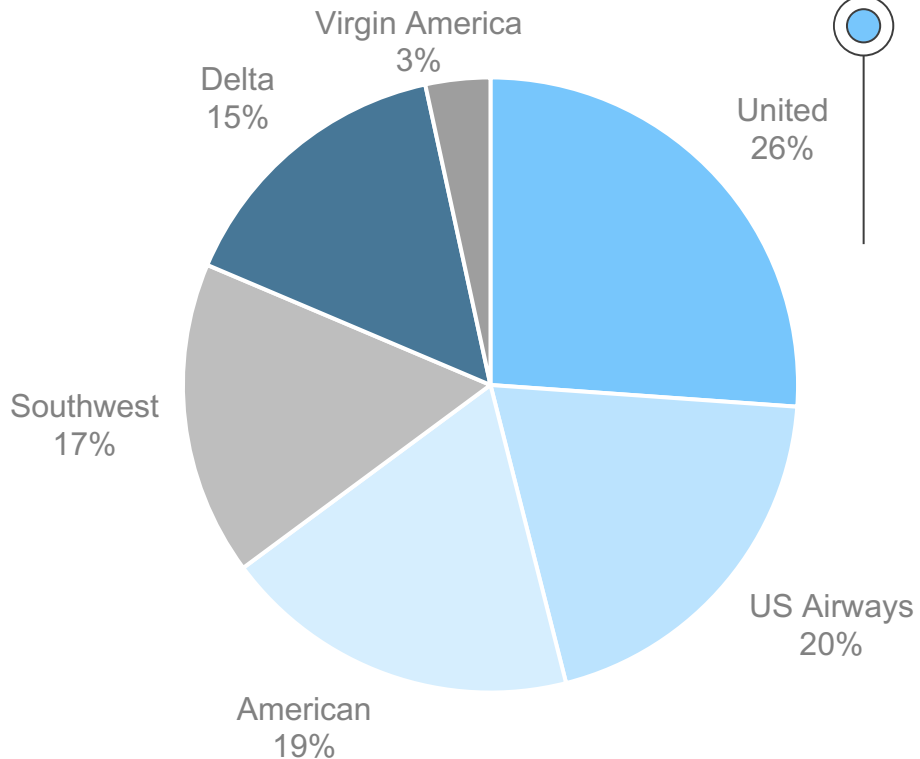
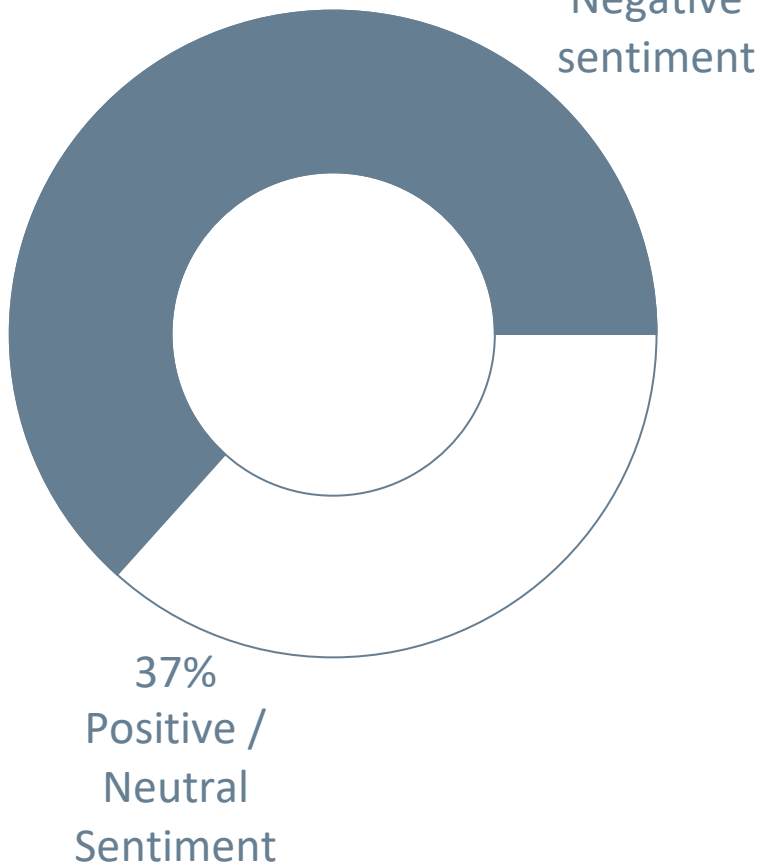


# Exploratory Data Analysis





# Data Overview

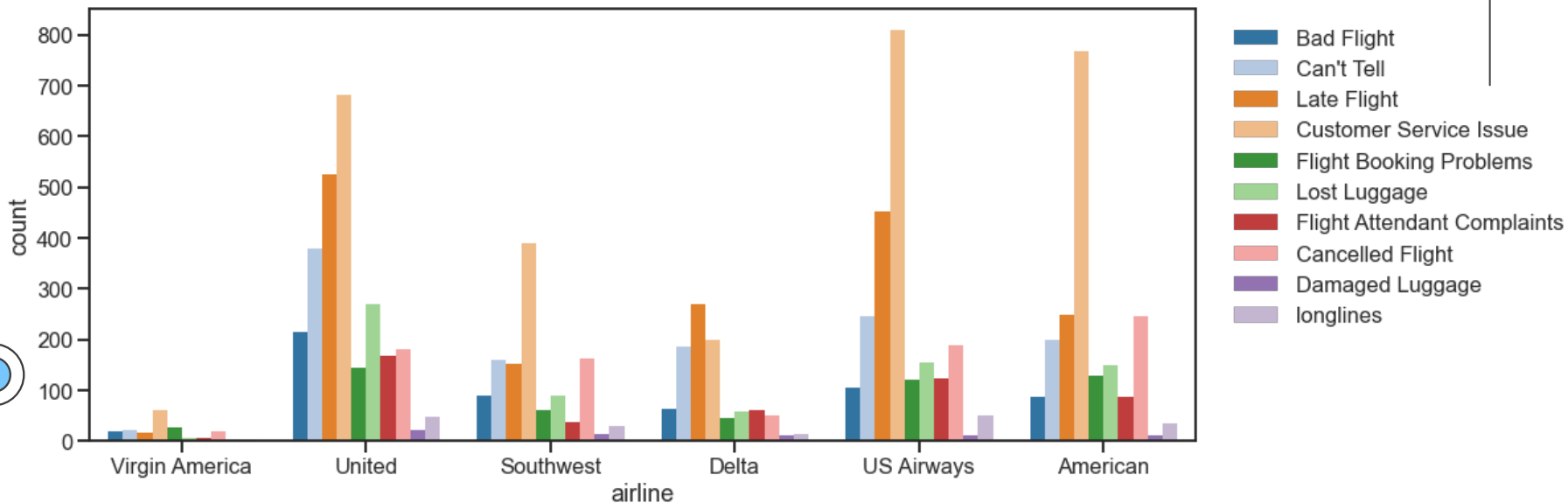




# Airlines Overview

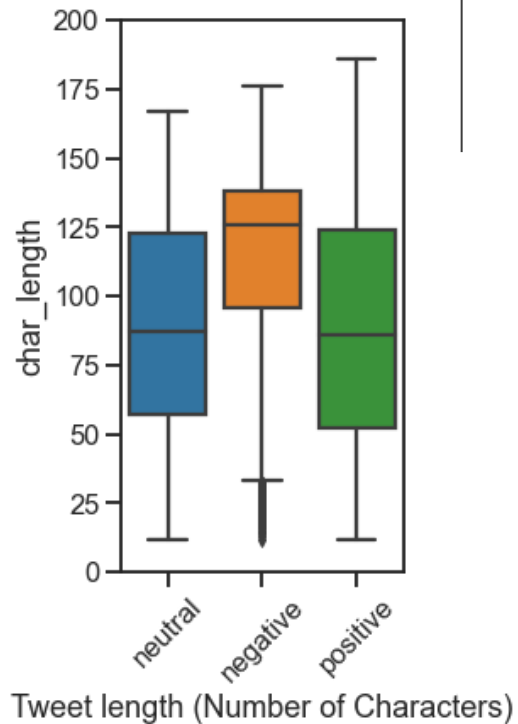
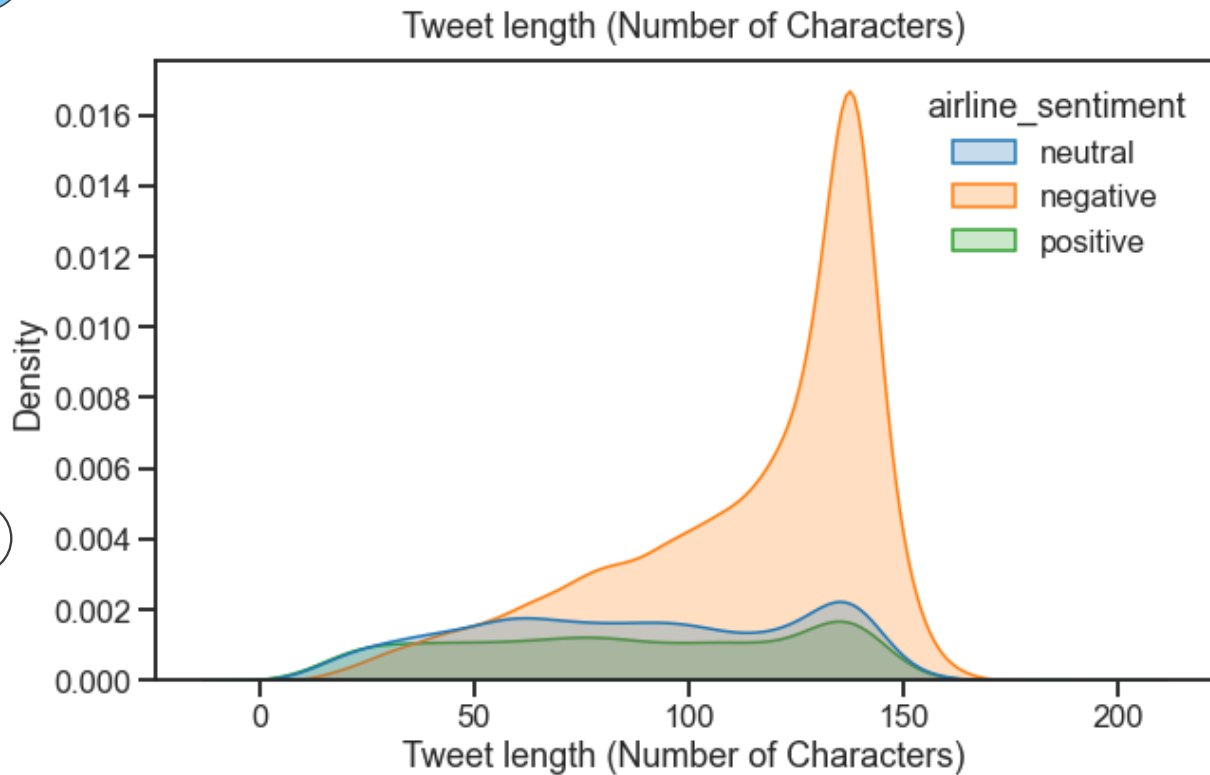


Negative Sentiment Categories of different airlines



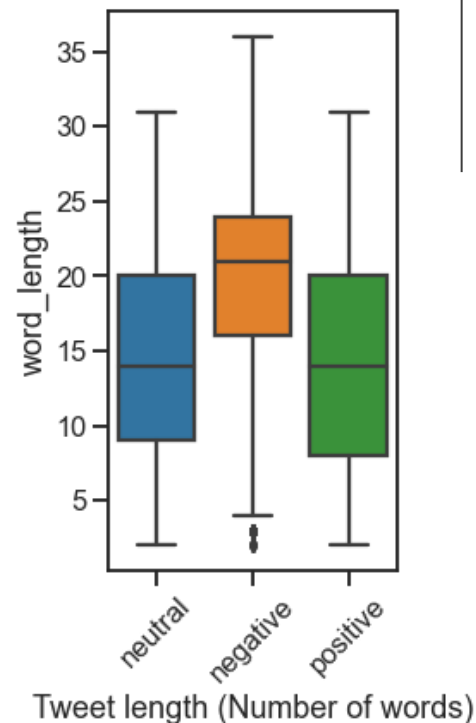
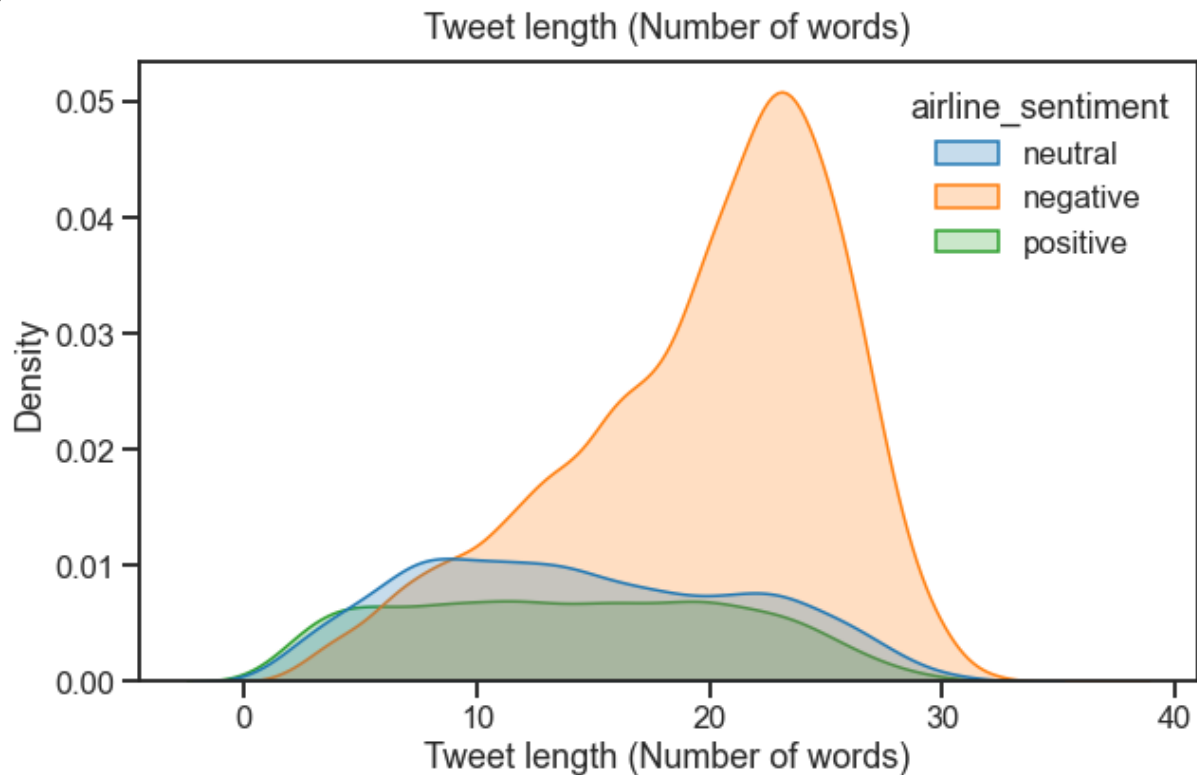


# Tweet length (#characters)



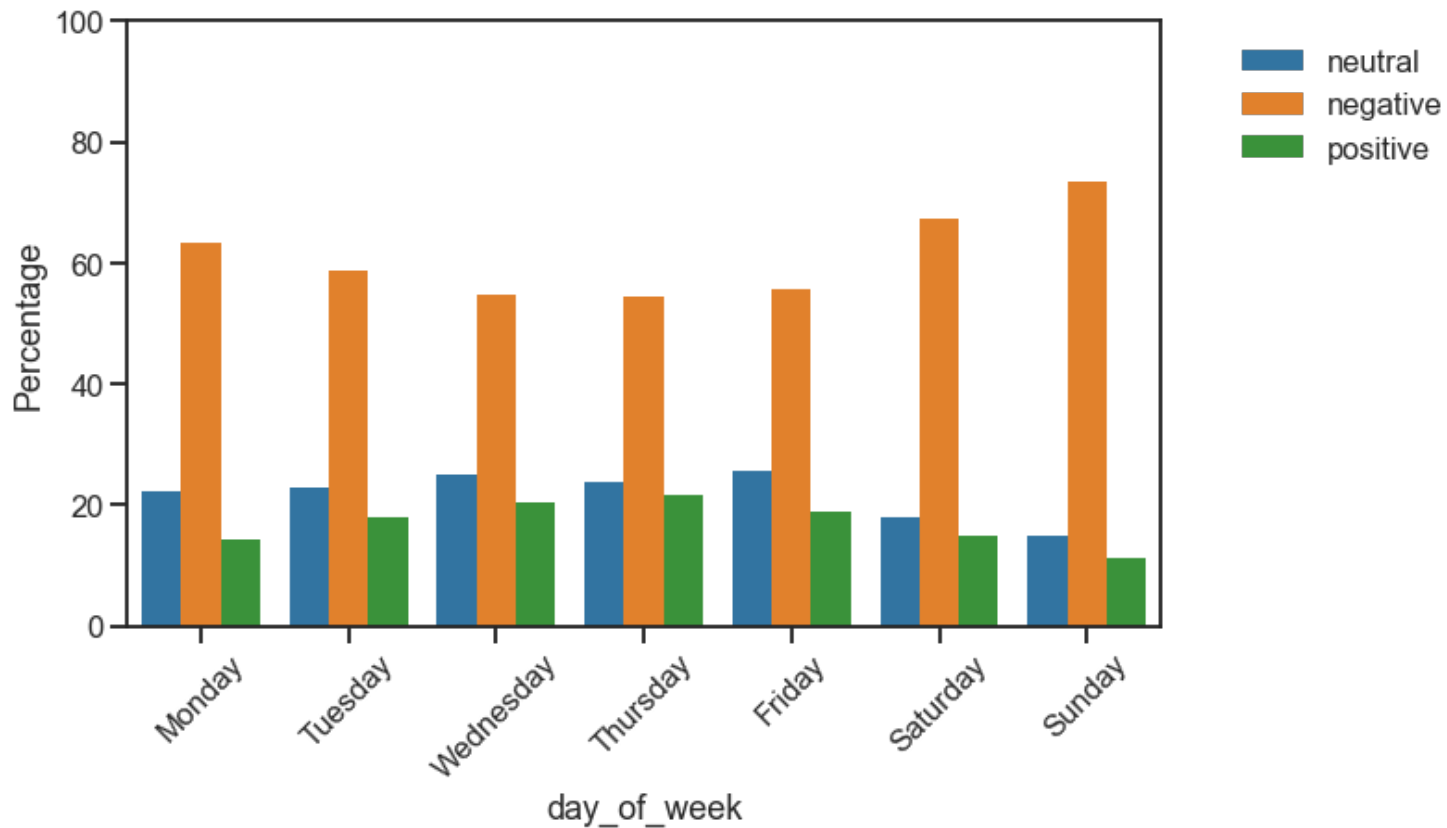


# Tweet length (#words)

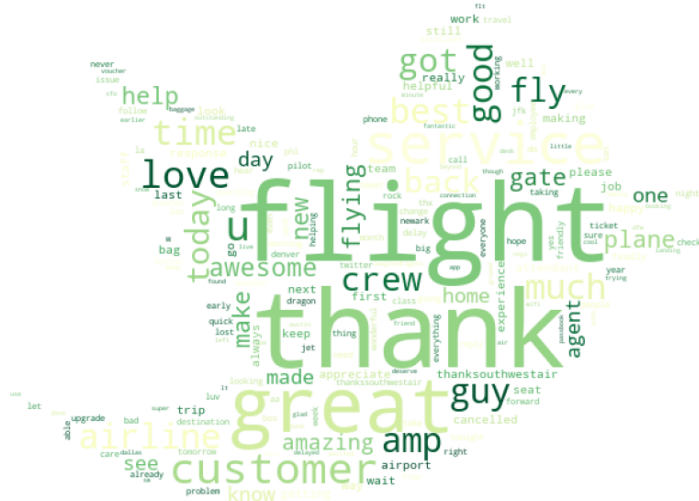




# Days of week



## Positive Sentiment



## Positive Sentiment



## Negative Sentiment

# 03 The

# Approach



# Algorithm outline

**Metamodel**  
“Take metadata  
such as day of  
week, tweet length  
as input”



**Textual model**  
“Take processed  
tweet as input”



**Combined  
model**  
“Combine prediction  
from both model and  
predict final  
probability”



\*Baseline Accuracy =  
62.7%



# Text Data Preprocessing



Hello, this is a TEST message for aviation Lover "forums". please visit at (<https://americanairlines.com>) #awesomeairline

hello, this is a test message for aviation lover "forums". please visit at (<https://americanairlines.com>) #awesomeairline

hello this is a test message for aviation lover forums please visit at #awesomeairline

convert to lower

Remove HTML / special char.

Segment hashtag

Remove stopwords

Remove airline entity

hello this is a test message for aviation lover forums please visit at awesome airline

hello this is a test message for aviation lover forums please visit at airline

hello test message aviation lover forums please visit airline

Lemmatize

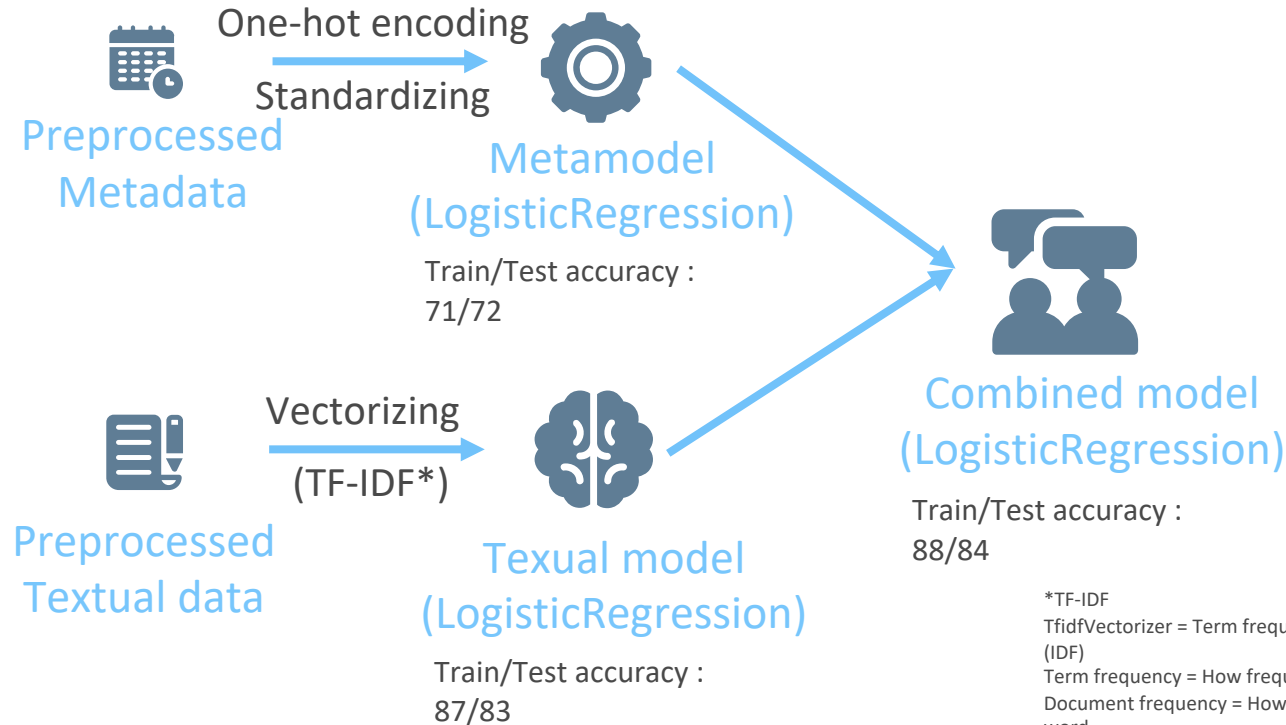
hello test message aviation love forum please visit airline

0.26 0.47 0.56 0.34 0.23 0.52 ...

Vectorize



# The Result

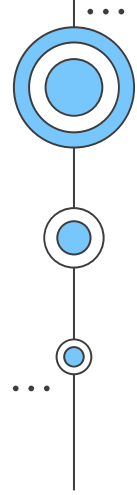


\*TF-IDF

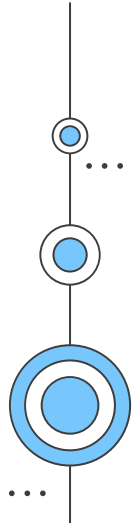
TfidfVectorizer = Term frequency (TF) x Inverse of Document Frequency (IDF)

Term frequency = How frequency of the word occurred in that document  
Document frequency = How frequency of the document containing that word

Words that occur often in one document but don't occur in many documents contain more predictive power (will get high vectorized value)



# 04 The Conclusion





# Performance Summary



## Textual model

Learned from tweet messages

83.5 % Accuracy

## Combined model

Combine prediction from both models

84.2 % Accuracy



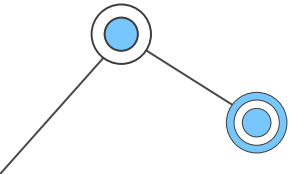
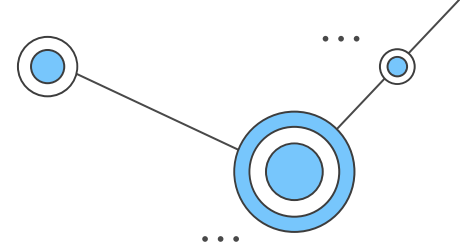
## Metamodel

Learned from metadata

72.1 % Accuracy

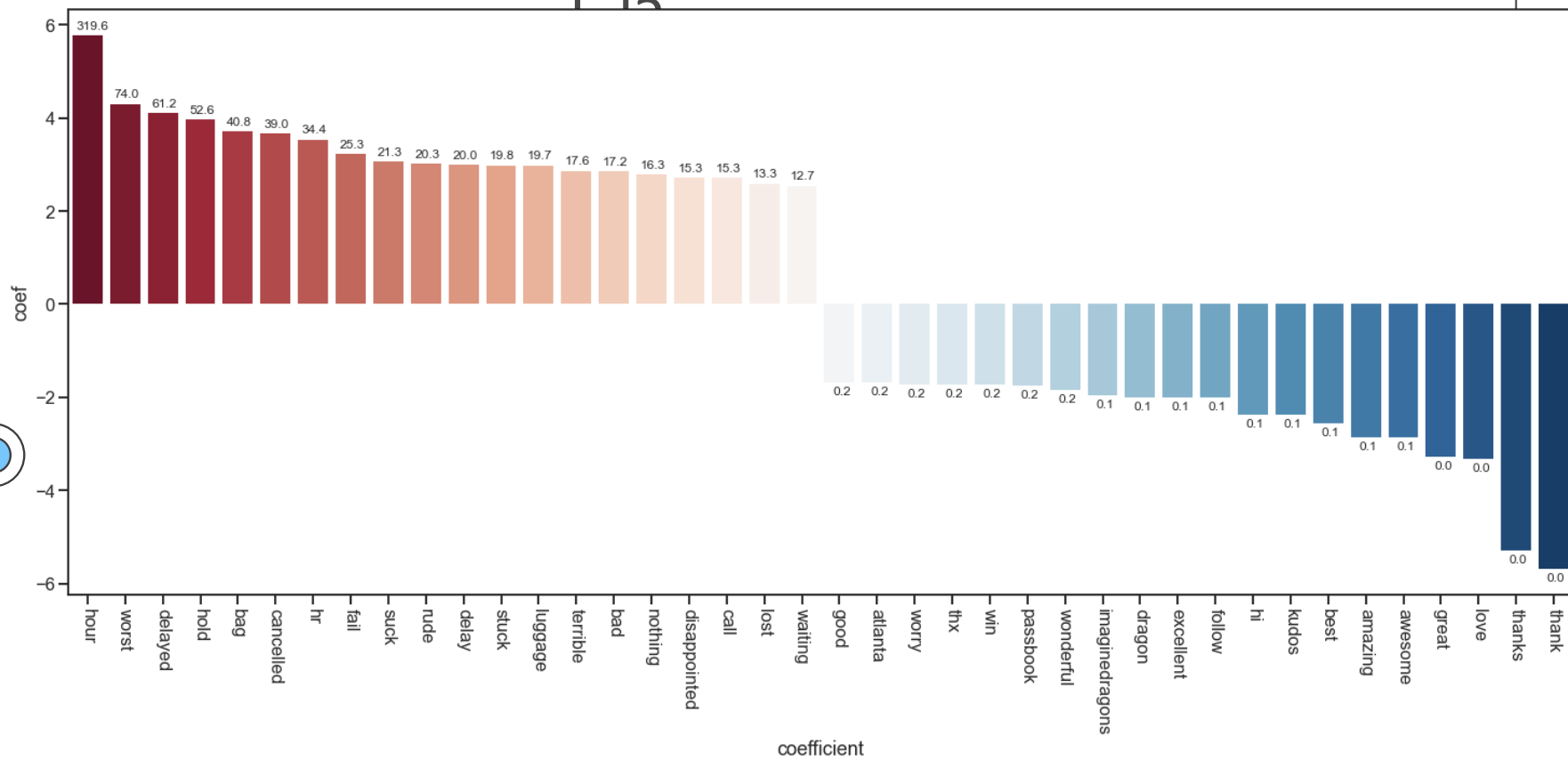
Baseline accuracy

62.7%



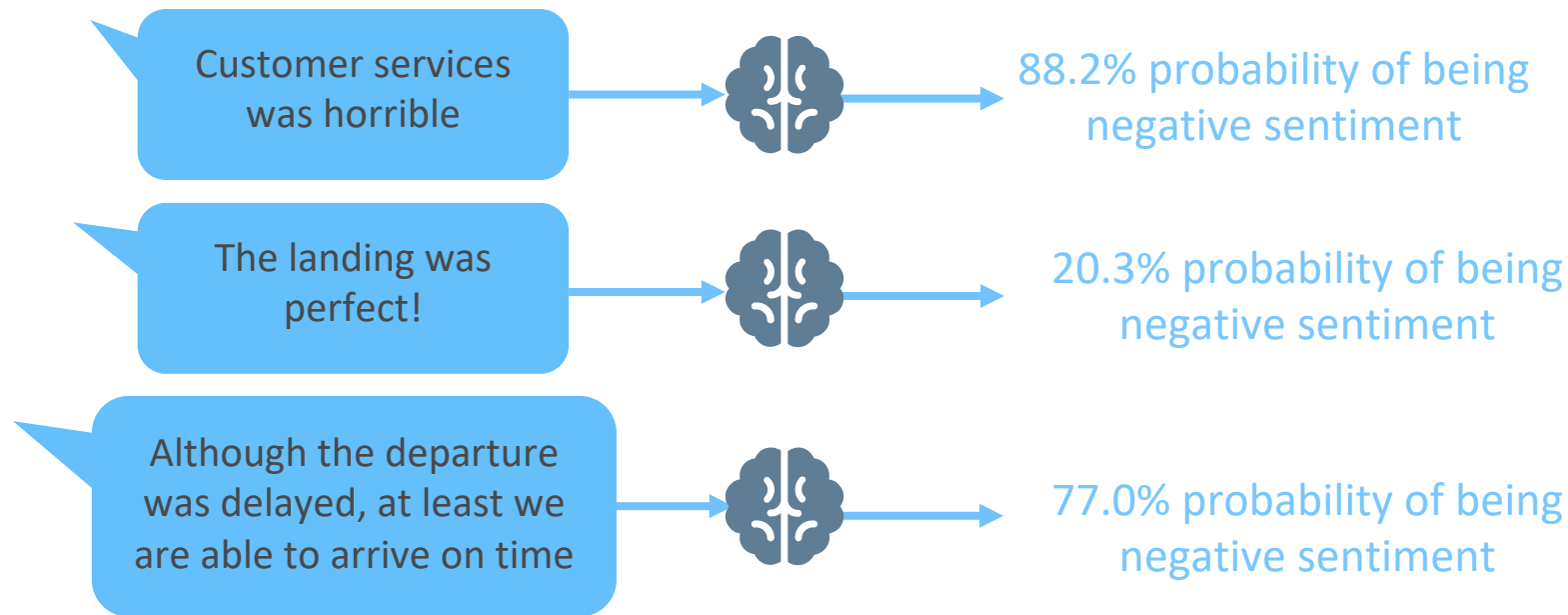
# What can be inferred from the

Top 20 keywords distinguishing negative sentiment from positive sentiment





# Prediction and Limitation





# Conclusion



01

## How can it be used?

Hooked with the data pipeline to filter the negative sentiment

02

## Limitation

A bag of words model may not effectively predict complex sentences

03

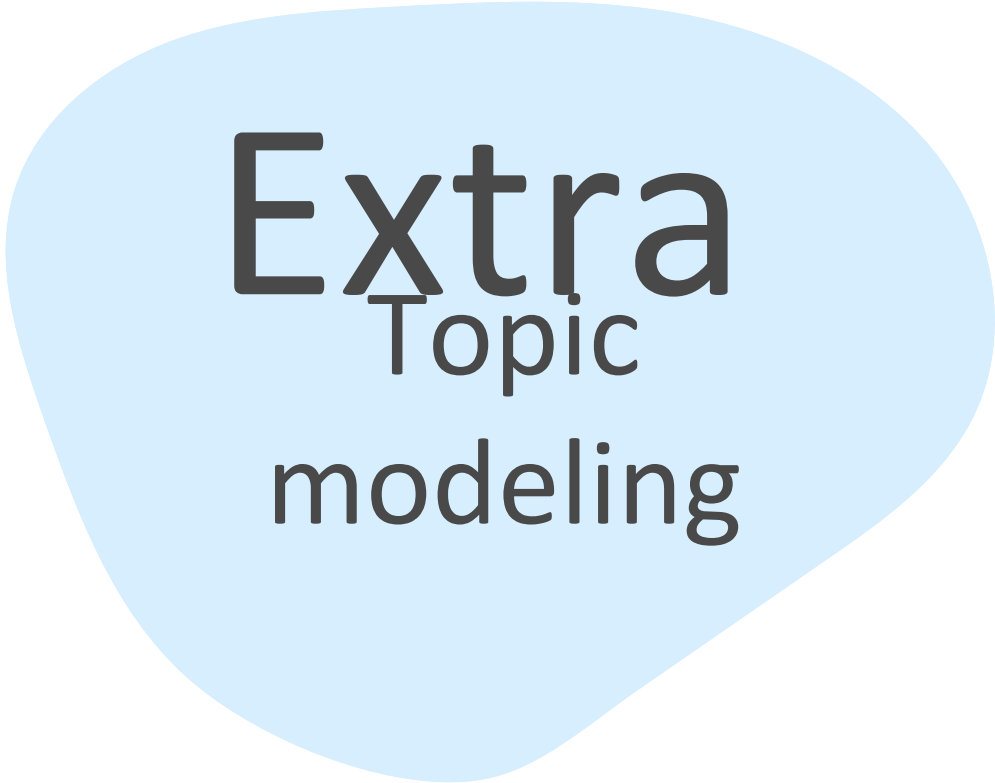

## Further improvement (The data)

- Gathering data from extended time span
- Acquiring more data to train deep learning model


04

## Further improvement (The model)

- Train more complex model (require massive dataset)
- Build classifier to subcategorize negative feedbacks



# Extra Topic modeling



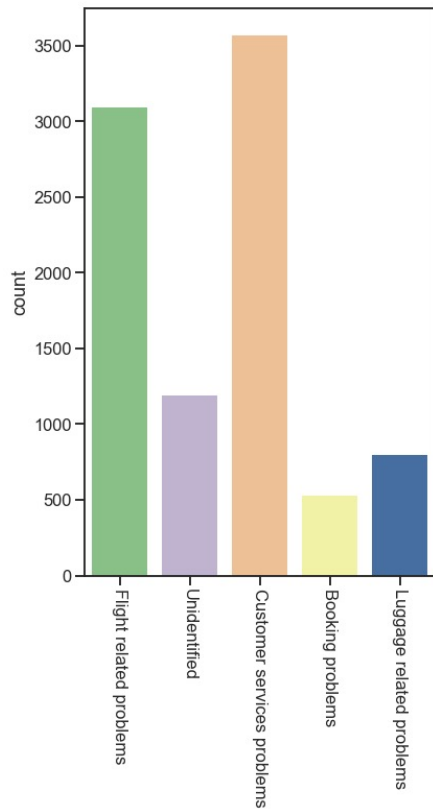
# Without topic modeling



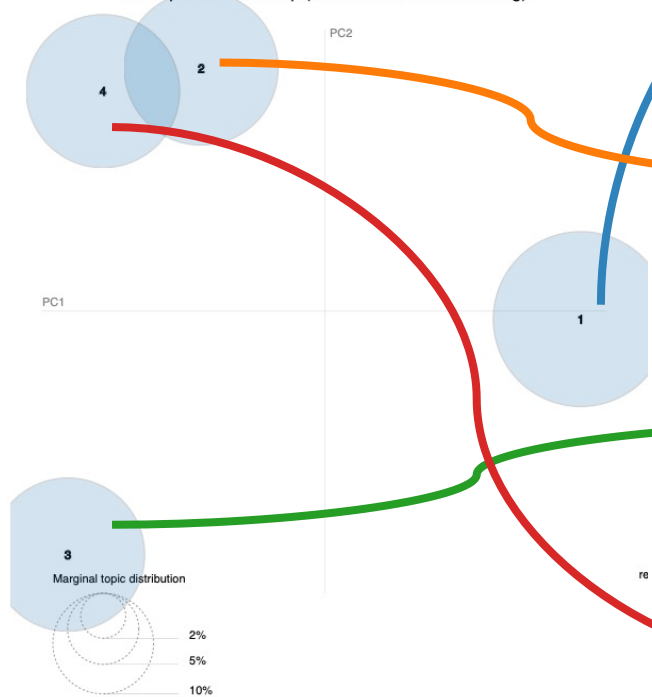


# Ground truth and topic modeling result

Ground truth



Intertopic Distance Map (via multidimensional scaling)



flight

cancel  
amp tomorrow bad gate  
flightle delay dfw  
flightled

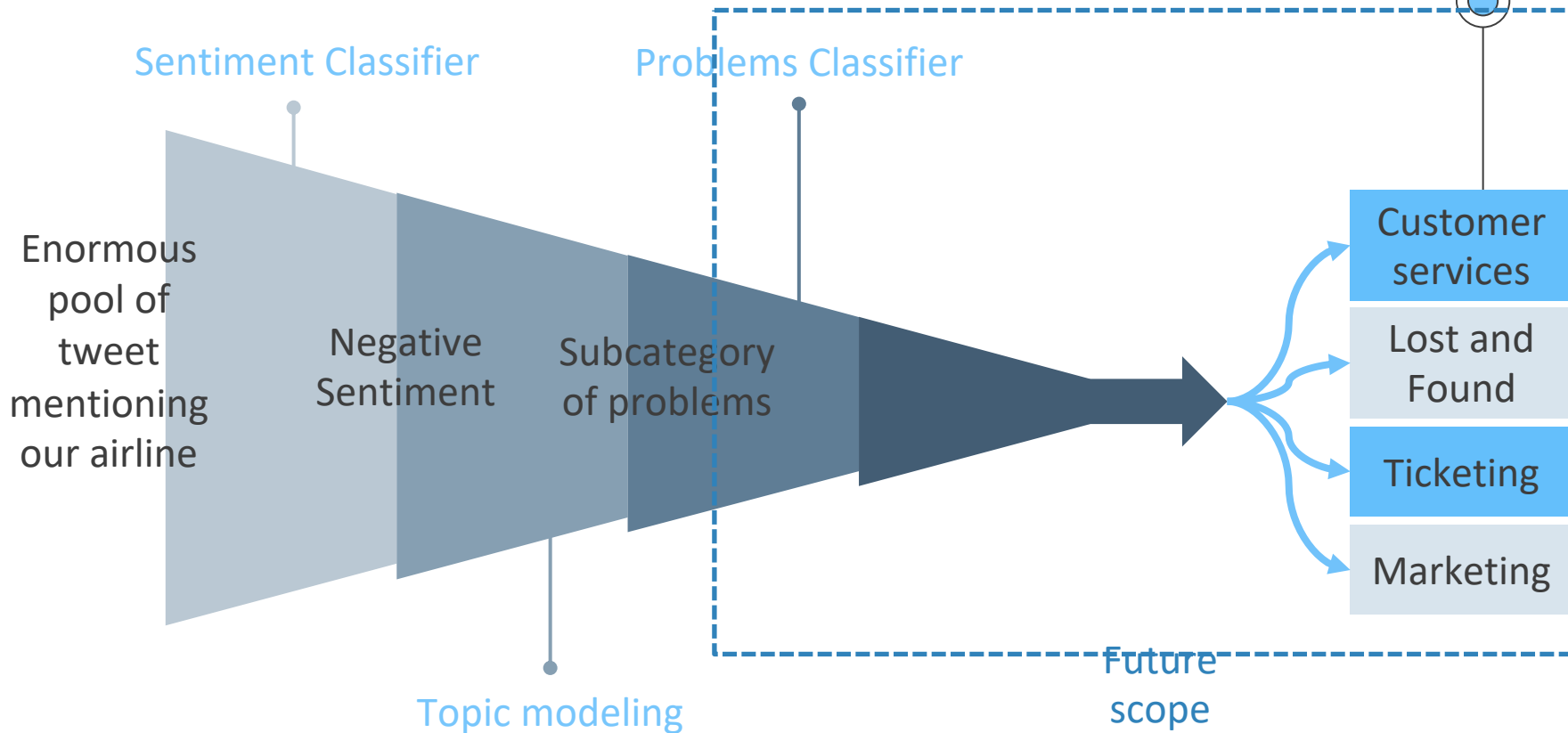
aa tell  
plane work  
guy weather  
agent seat  
time luggage

leave phone  
help reservation  
hour hold bag  
rebook minute  
ticket

call back wait  
day  
service customer  
still airport issue  
response



# What can topic modeling do?





# Thanks!

Do you have any questions?

Wirach.lee@gmail.com

+66 81 453 6326

Or visit the project repo at

[https://github.com/Joeycook/DSI\\_CapstoneProject](https://github.com/Joeycook/DSI_CapstoneProject)

**CREDITS:** This presentation template was created by [Slidesgo](#),  
including icons by [Flaticon](#), infographics & images by [Freepik](#) and  
illustrations by [Stories](#)

Please keep this slide for attribution