



# JUST EAT

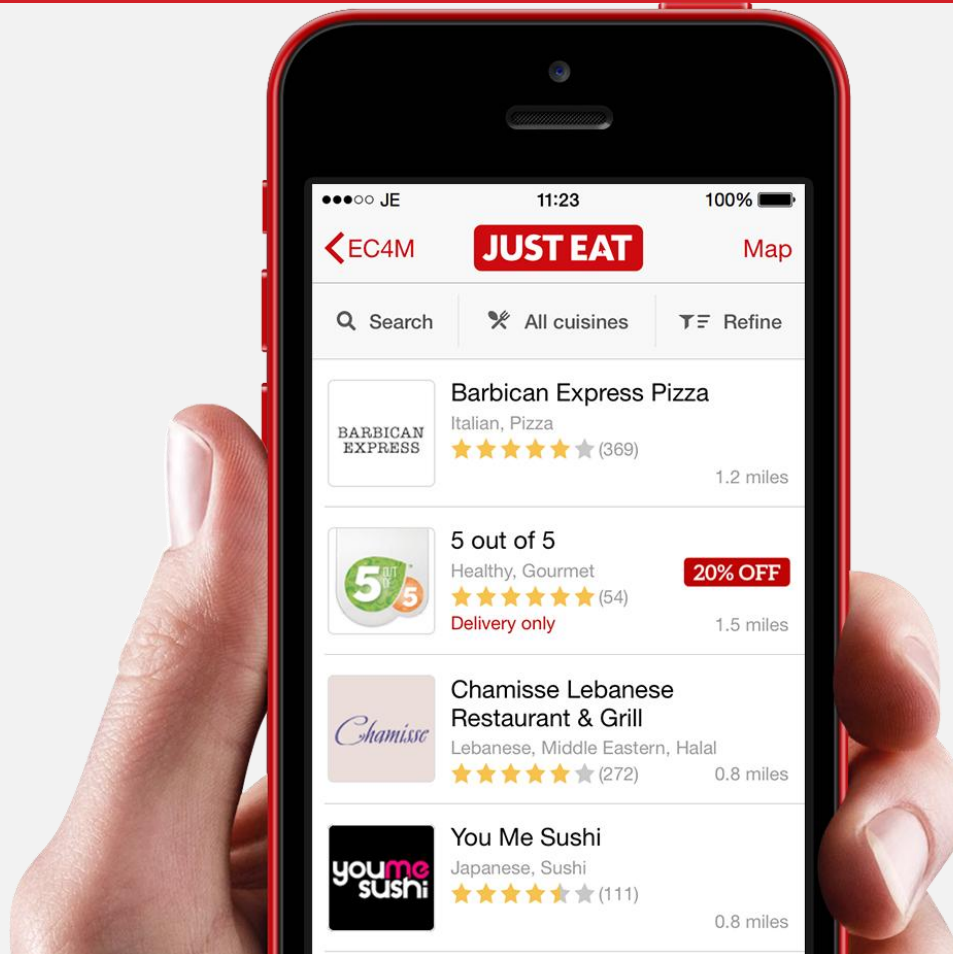
Natural language processing  
and  
Sentiment analysis



← **JUST EAT** →

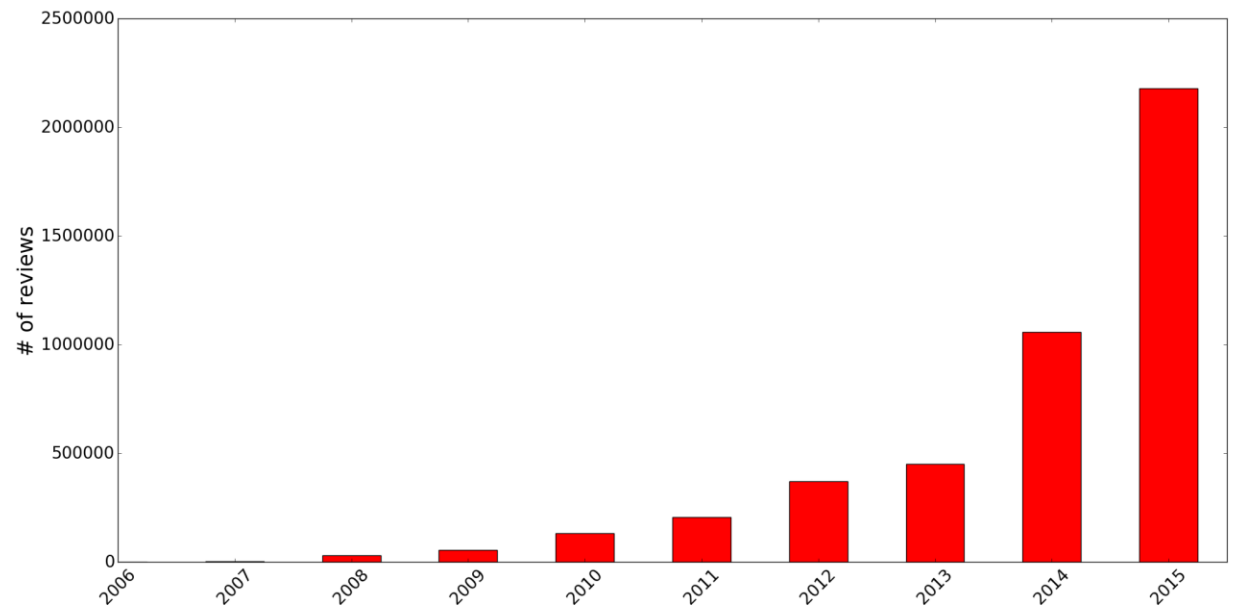
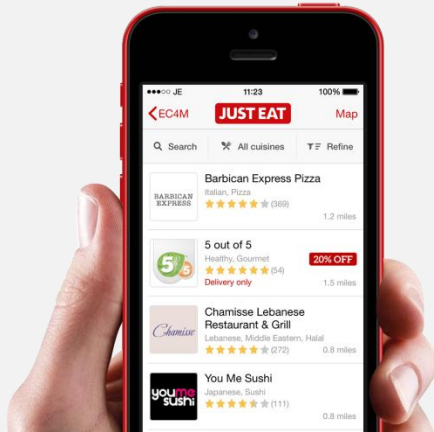


# JUST EAT

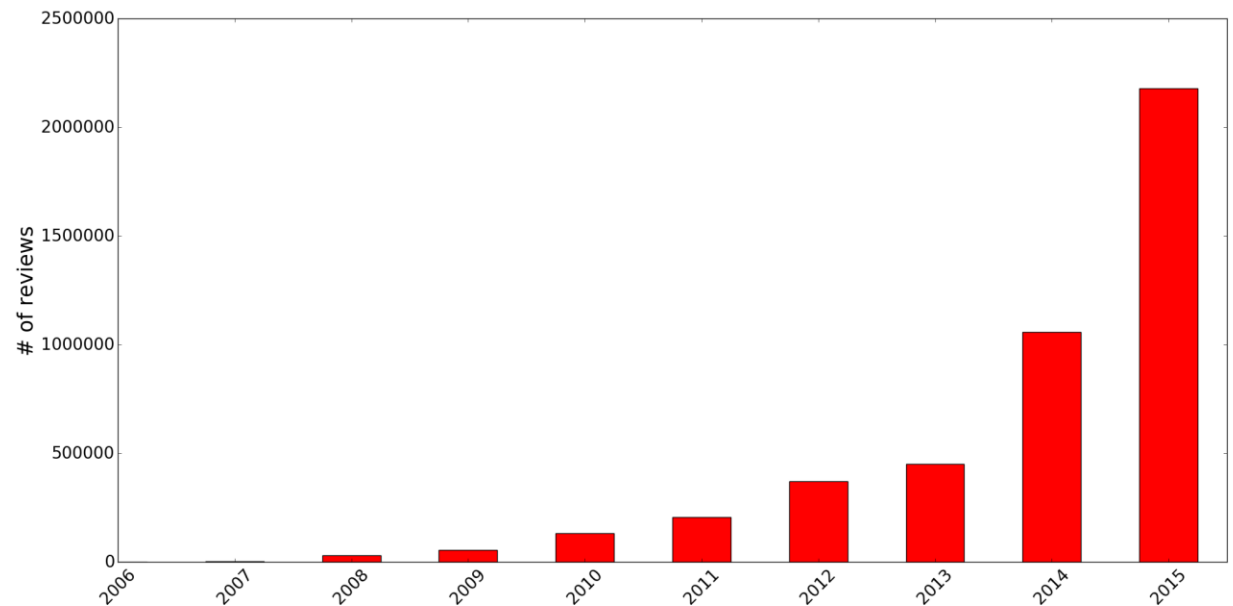


Rolando P. Hong Enriquez

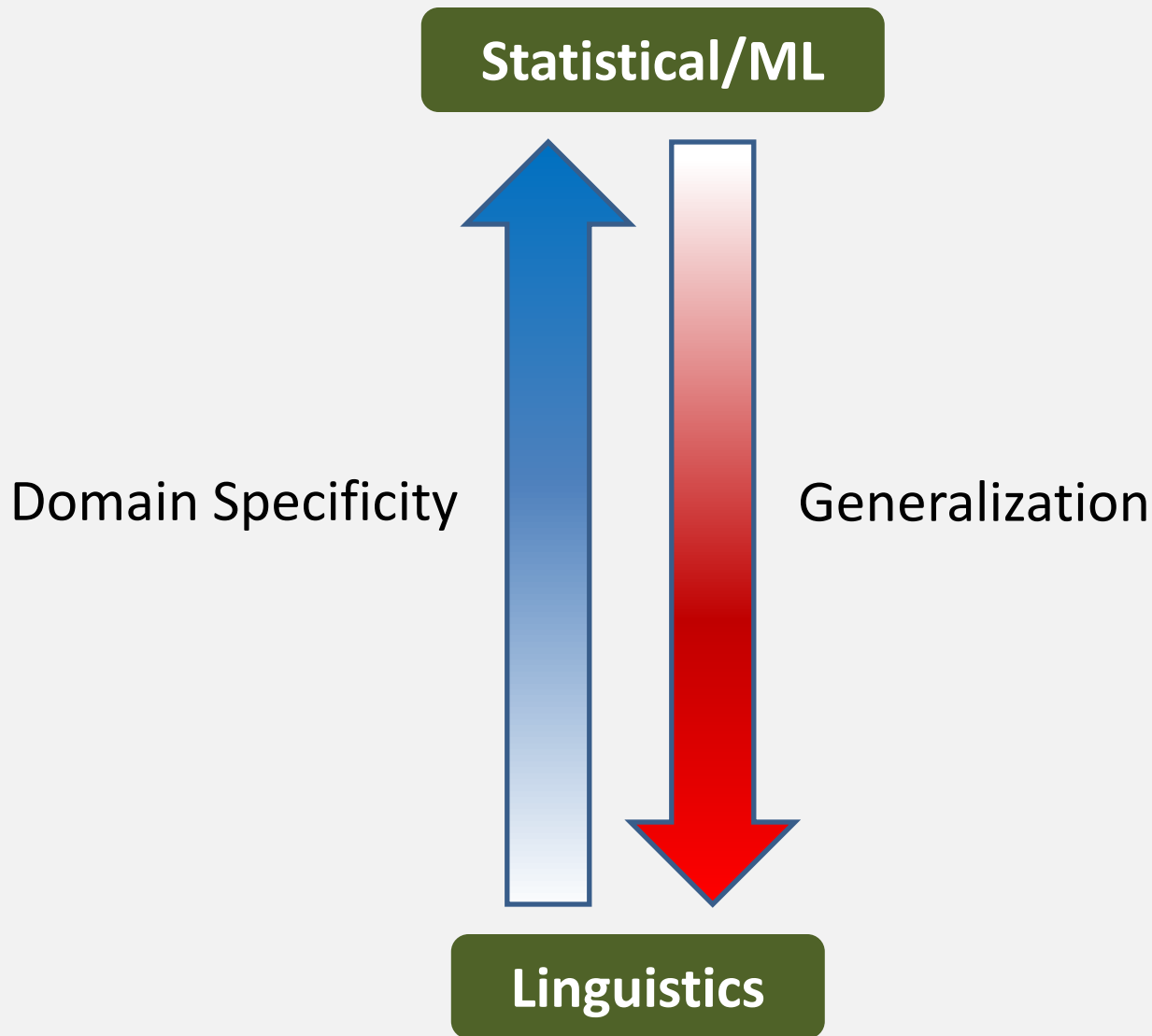
# JUST EAT



# How can we efficiently analyze the sentiment in these data?

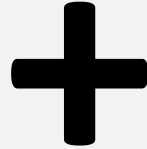


# Two Approaches



# Two Approaches

Statistical/ML

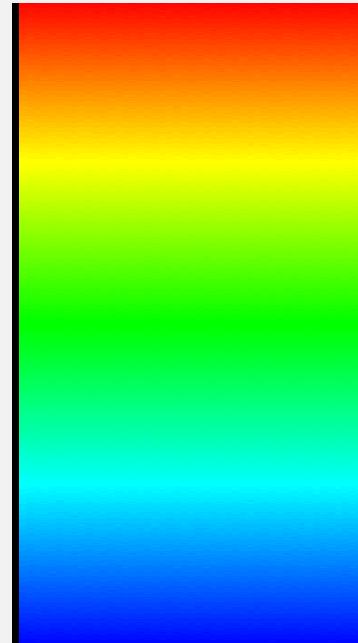


Linguistics



# Goal

High sentiment



Low sentiment

['word1' 'word2' 'word3']

# THE PIPELINE

[“Wheeeen I am alone, a NormallY enjoy a good pizza!! 😊”]

# THE PIPELINE

[“Wheeeen I am alone, a NormalY enjoy a good pizza!! 😊”]

# THE PIPELINE

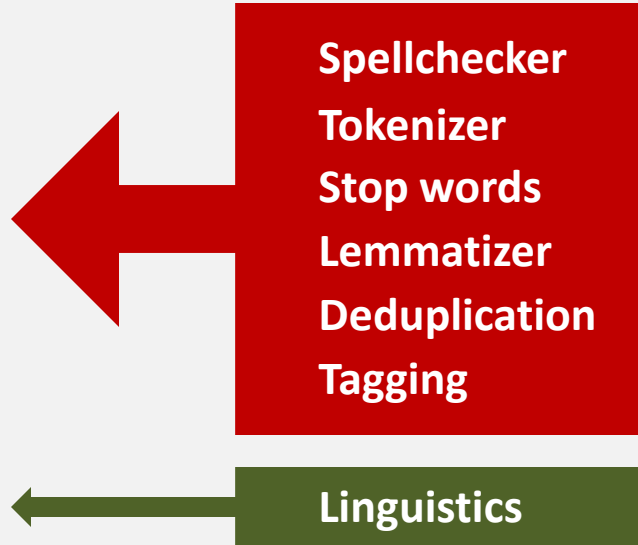
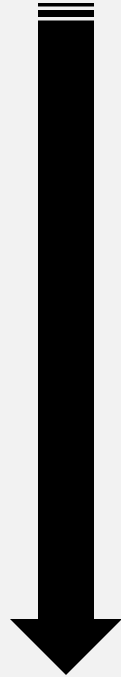
[“Wheeeen I am alone, a NormalY enjoy a good pizza!! 😊”]



Spellchecker  
Tokenizer  
Stop words  
Lemmatizer  
Deduplication  
Tagging

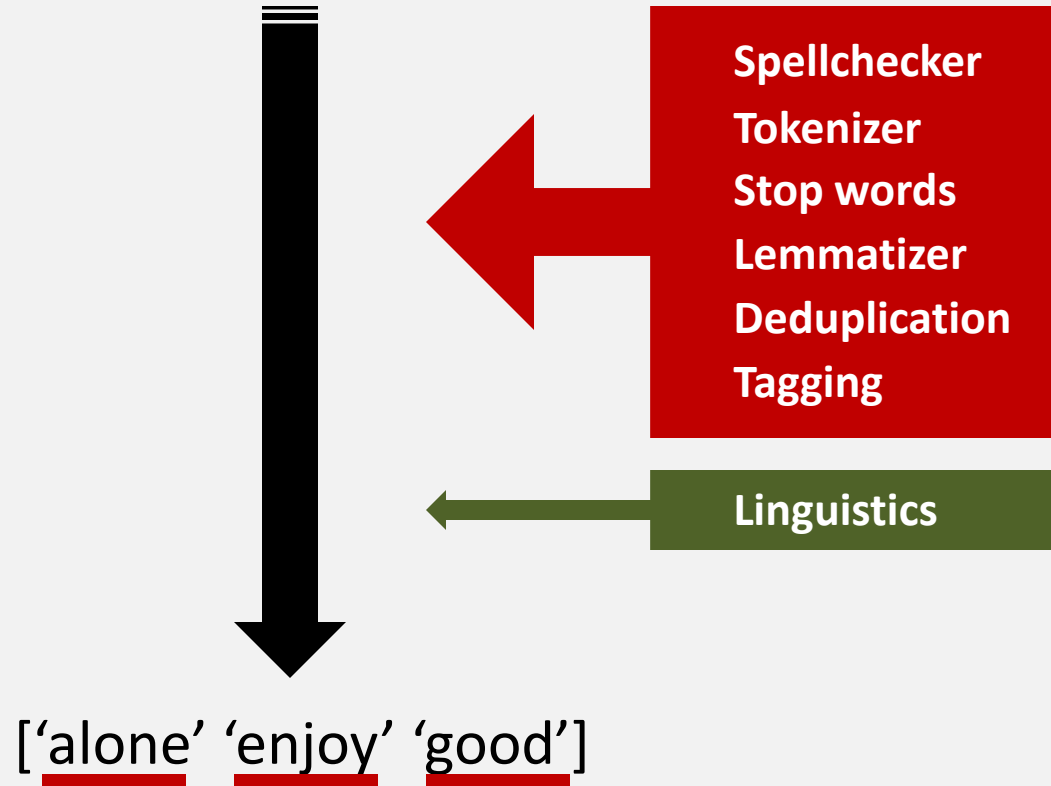
# THE PIPELINE

[“Wheeeen I am alone, a Normally enjoy a good pizza!! 😊”]



# THE PIPELINE

[“Wheeeen I am alone, a Normally enjoy a good pizza!! 😊”]



# THE PIPELINE

["She runs a lot!!  
When I run that much  
I normally need to enjoy a good pizza alone in the evening."]

# THE PIPELINE

["She runs a lot!!  
When I run that much  
I normally need to enjoy a good pizza alone in the evening."]



# Tokenize

["She runs a lot!!  
When I run that much  
I normally need to enjoy a good pizza alone in the evening."]

# Tokenize

['She' 'runs' 'a' 'lot' '!' '!' 'When' 'I' 'run' 'that' 'much'  
'I' 'normally' 'need' 'to' 'enjoy' 'a' 'good' 'pizza'  
'alone' 'in' 'the' 'evening' '.']

# Remove stop words

['She' 'runs' 'a' lot' '!' '!' 'When' 'I' 'run' 'that' 'much'  
'I' 'normally' 'need' 'to' 'enjoy' 'a' 'good' 'pizza'  
'alone' 'in' 'the' 'evening' '.']

# Remove stop words

['She' 'runs' 'lot' 'I' 'run' 'much' 'I' 'normally'  
'need' 'enjoy' 'good' 'pizza' 'alone' 'evening']

# Lemmatize

['She' 'runs' 'lot' 'I' 'run' 'much' 'I' 'normally'  
'need' 'enjoy' 'good' 'pizza' 'alone' 'evening']

# Lemmatize

['She' 'run' 'lot' 'I' 'run' 'much' 'I' 'normally'  
'need' 'enjoy' 'good' 'pizza' 'alone' 'evening']

# Eliminate duplicates

['She' 'run' 'lot' 'I' 'run' 'much' 'I' 'normally'  
'need' 'enjoy' 'good' 'pizza' 'alone' 'evening']

# Eliminate duplicates

['She' 'run' 'lot' 'much' 'I' 'normally' 'need' 'enjoy'  
'good' 'pizza' 'alone' 'evening']



# Tagging

['She' 'run' 'lot' 'much' 'I' 'normally' 'need' 'enjoy'  
'good' 'pizza' 'alone' 'evening']

- Verbs & Adverbs
- Adjectives
- Nouns & Personal Pronouns

# Eliminate objective parts

['She' 'run' 'lot' 'much' 'I' 'normally' 'need' 'enjoy'  
'good' 'pizza' 'alone' 'evening']

- Verbs & Adverbs
- Adjectives
- Nouns & Personal Pronouns

# Eliminate objective parts

['run' 'lot' 'much' 'normally' 'need' 'enjoy'  
'good' 'alone']

— Verbs & Adverbs

— Adjectives

# Valence Aware Dictionary for sEntiment Reasoning [1]

['run' 'lot' 'much' 'normally' 'need' 'enjoy'  
'good' 'alone']

- Verbs & Adverbs
- Adjectives

# VADER



['run' 'lot' 'much' 'normally' 'need' 'enjoy'  
'good' 'alone']

— Verbs & Adverbs  
— Adjectives

## VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text

**C.J. Hutto**

Georgia Institute of Technology, Atlanta, GA 30032  
cjhutto@gatech.edu

**Eric Gilbert**

gilbert@cc.gatech.edu

# VADER

['run' 'lot' 'much' 'normally' 'need' 'enjoy'  
'good' 'alone']

— VADER score  $\neq 0$



# VADER

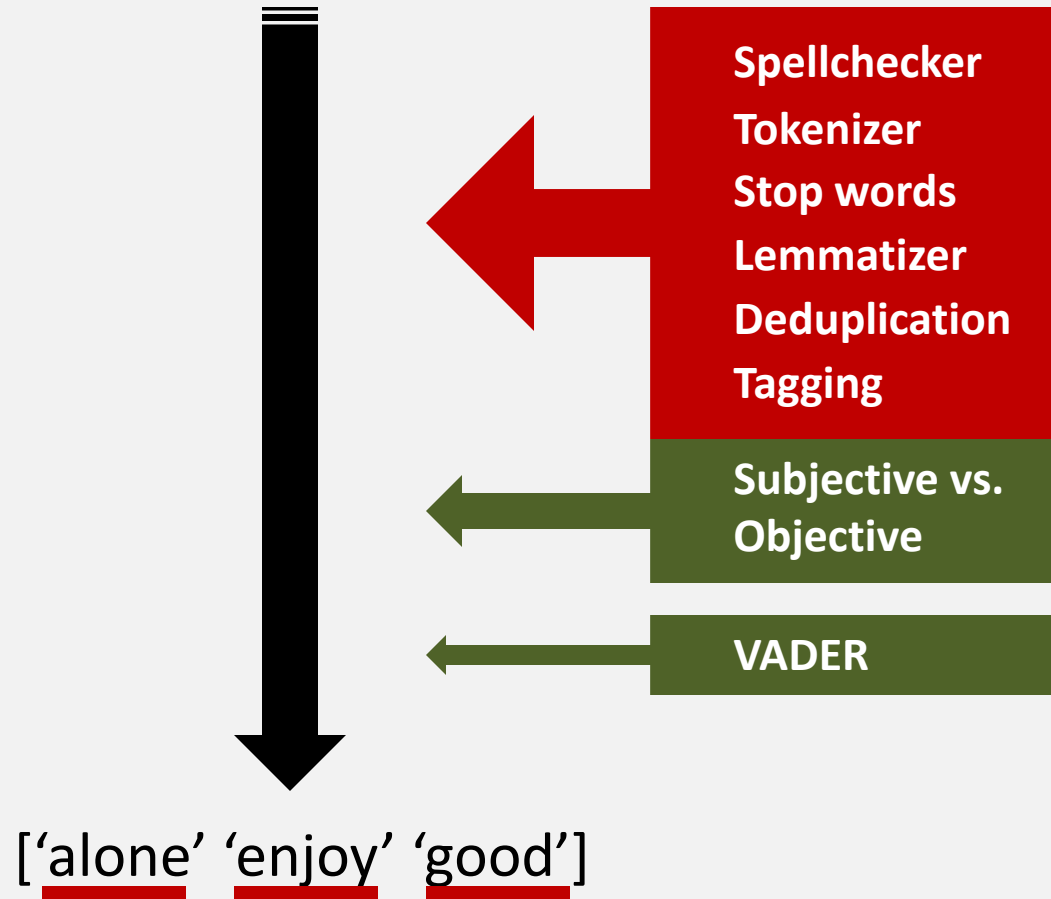
['enjoy' 'good' 'alone']

 VADER score  $\neq 0$



# THE PIPELINE

[“Wheeeen I am alone, a Normally enjoy a good pizza!! 😊”]





# Binarize

1 2 3 4 5  
['play' 'funny' 'sad' 'good' 'bad'] Minimal Dictionary

# Binarize

1 2 3 4 5  
['play' 'funny' 'sad' 'good' 'bad'] **Minimal Dictionary**

'This food is very good' **Review**

# Binarize

1 2 3 4 5  
['play' 'funny' 'sad' 'good' 'bad'] **Minimal Dictionary**

'This food is very good' **Review**

[0 0 0 1 0] **Binarized  
Review**

# Small Intro to ML models

# ML models

**Discriminative**



e.g., Linear/Logistic Regression

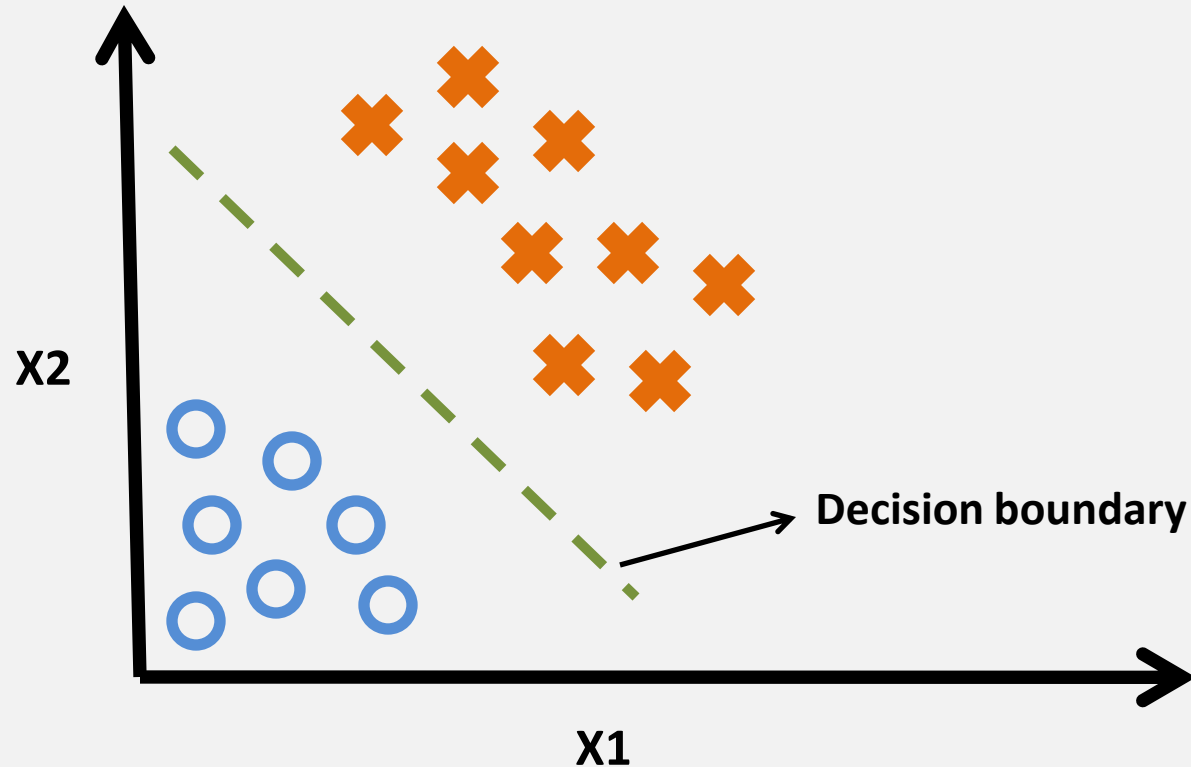
**Generative**



e.g., Naïve Bayes

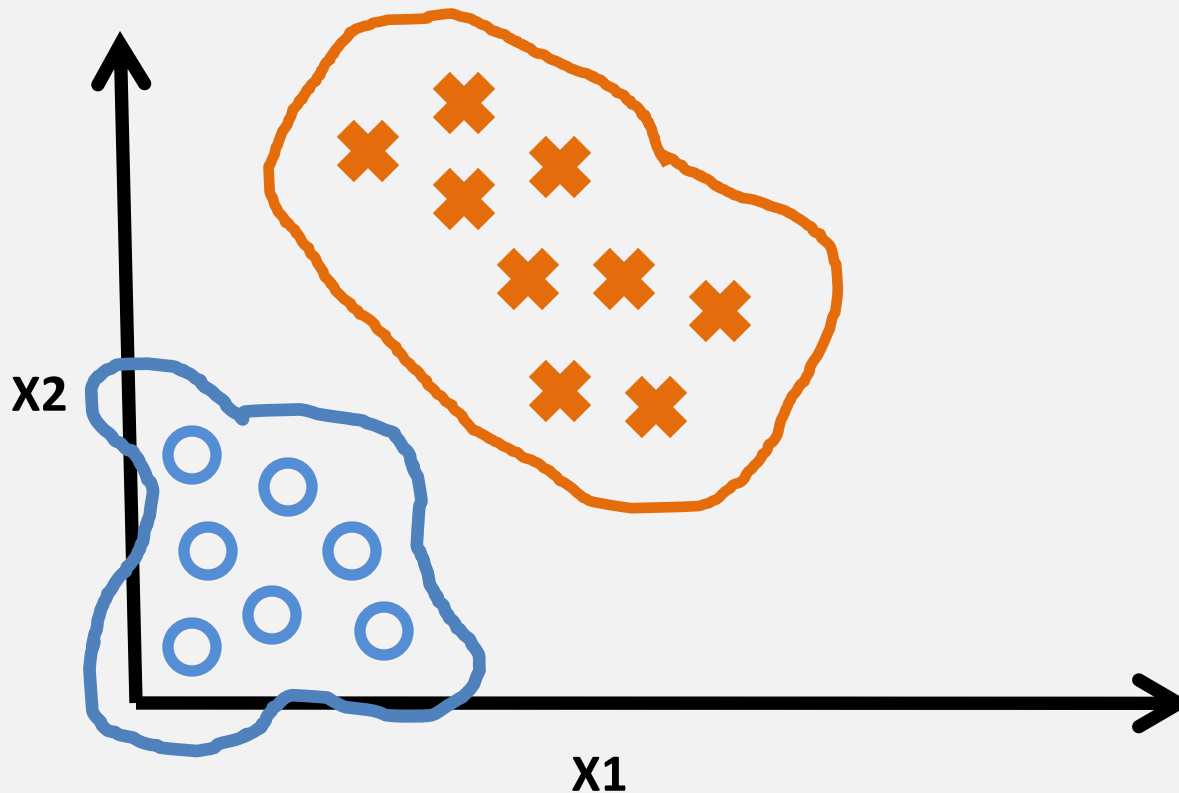
# ML models

## Discriminative



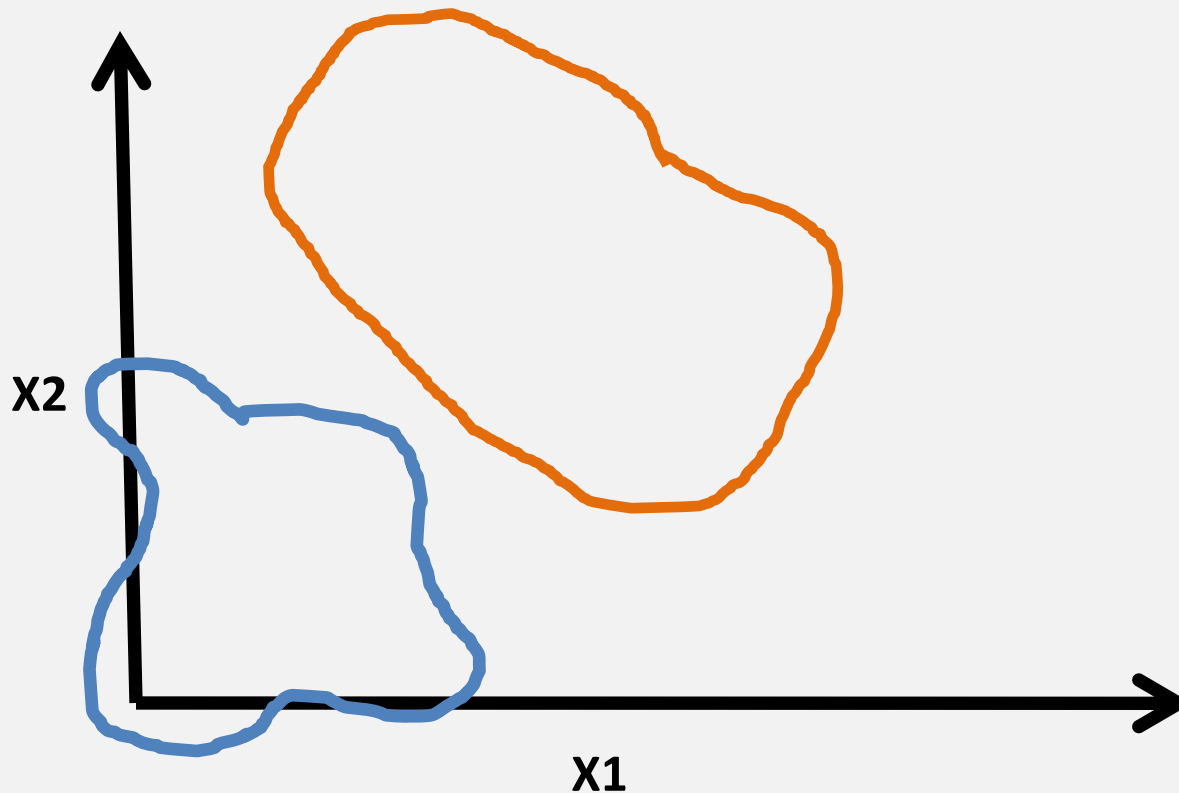
# ML models

Generative



# ML models

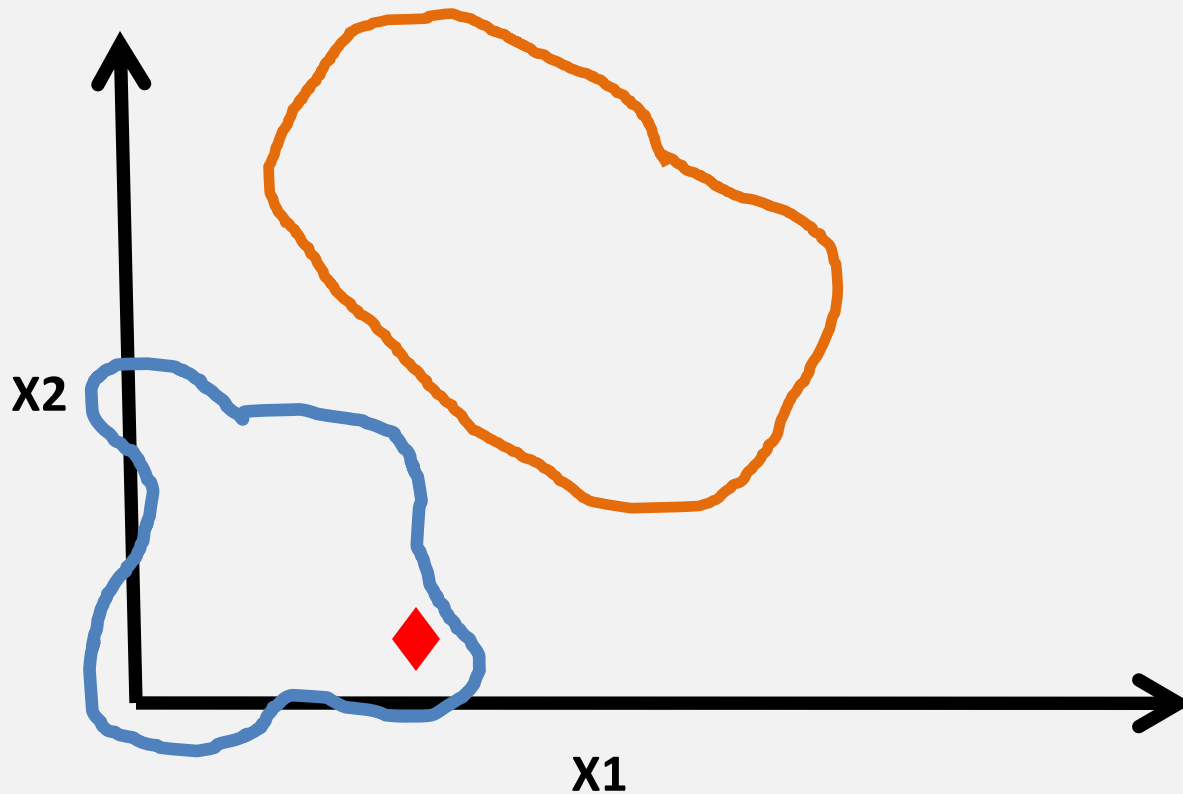
**Generative**





# ML models

Generative



# Training data

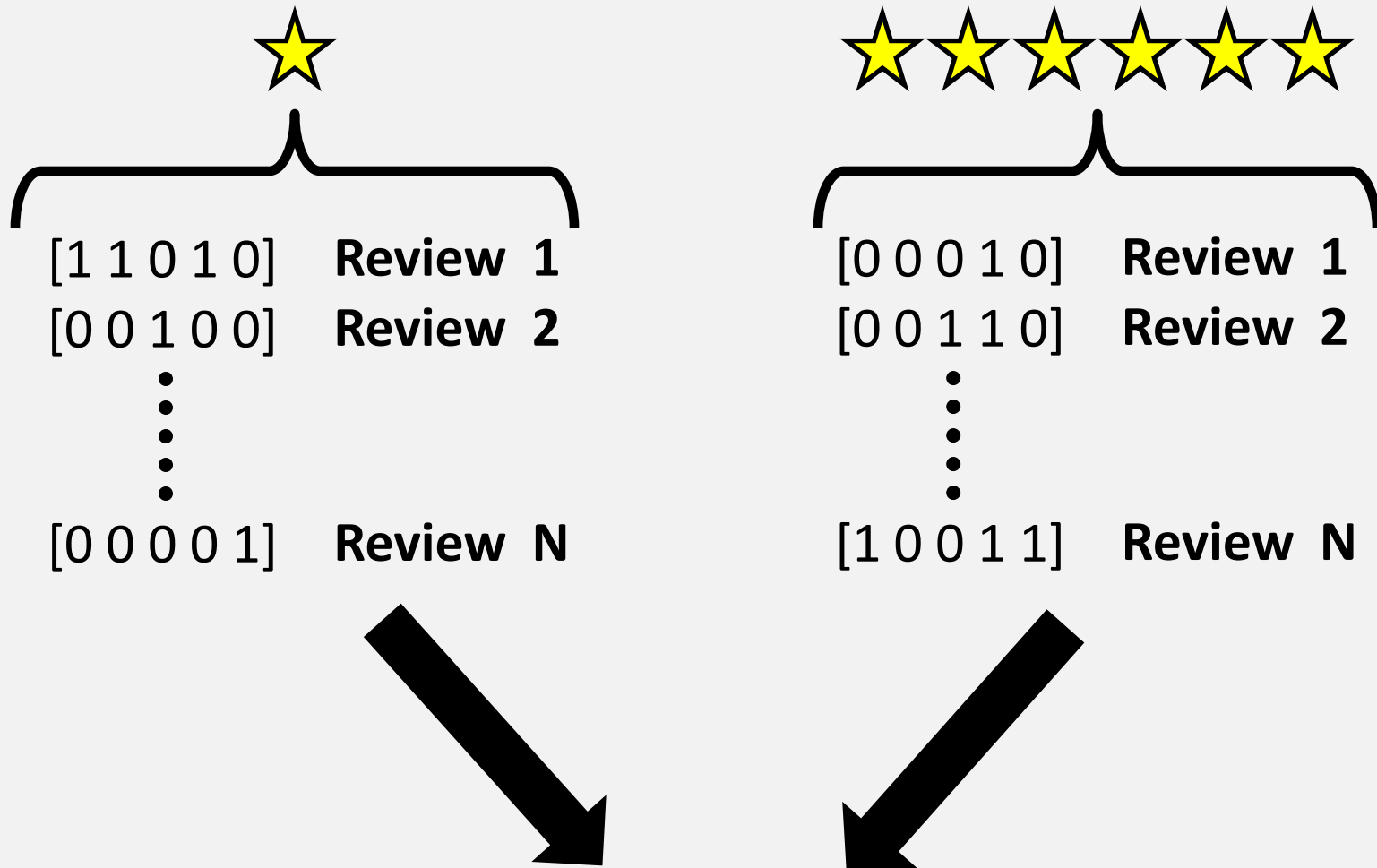


[1 1 0 1 0]	Review 1
[0 0 1 0 0]	Review 2
⋮	
[0 0 0 0 1]	Review N



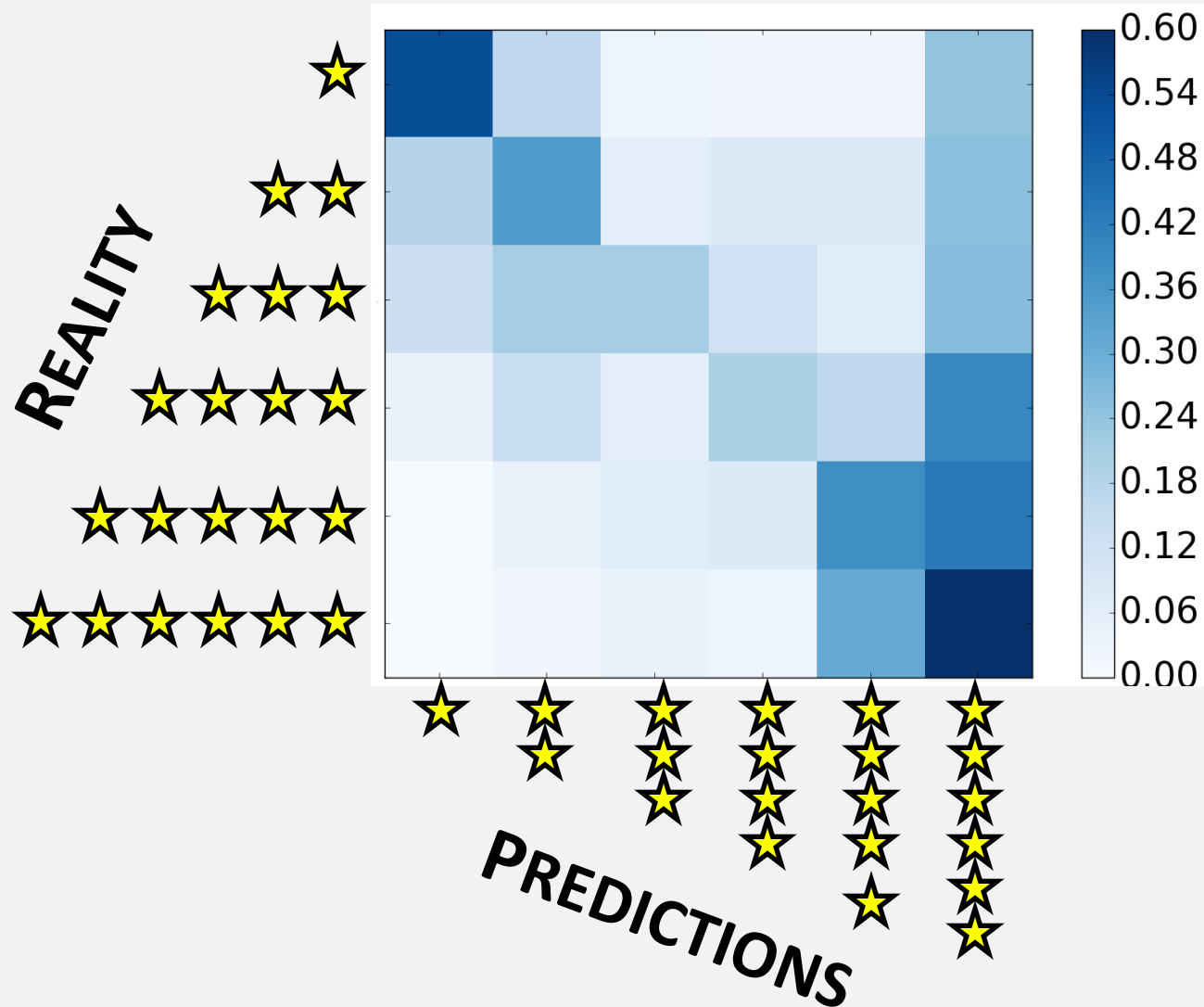
[0 0 0 1 0]	Review 1
[0 0 1 1 0]	Review 2
⋮	
[1 0 0 1 1]	Review N

# Training data



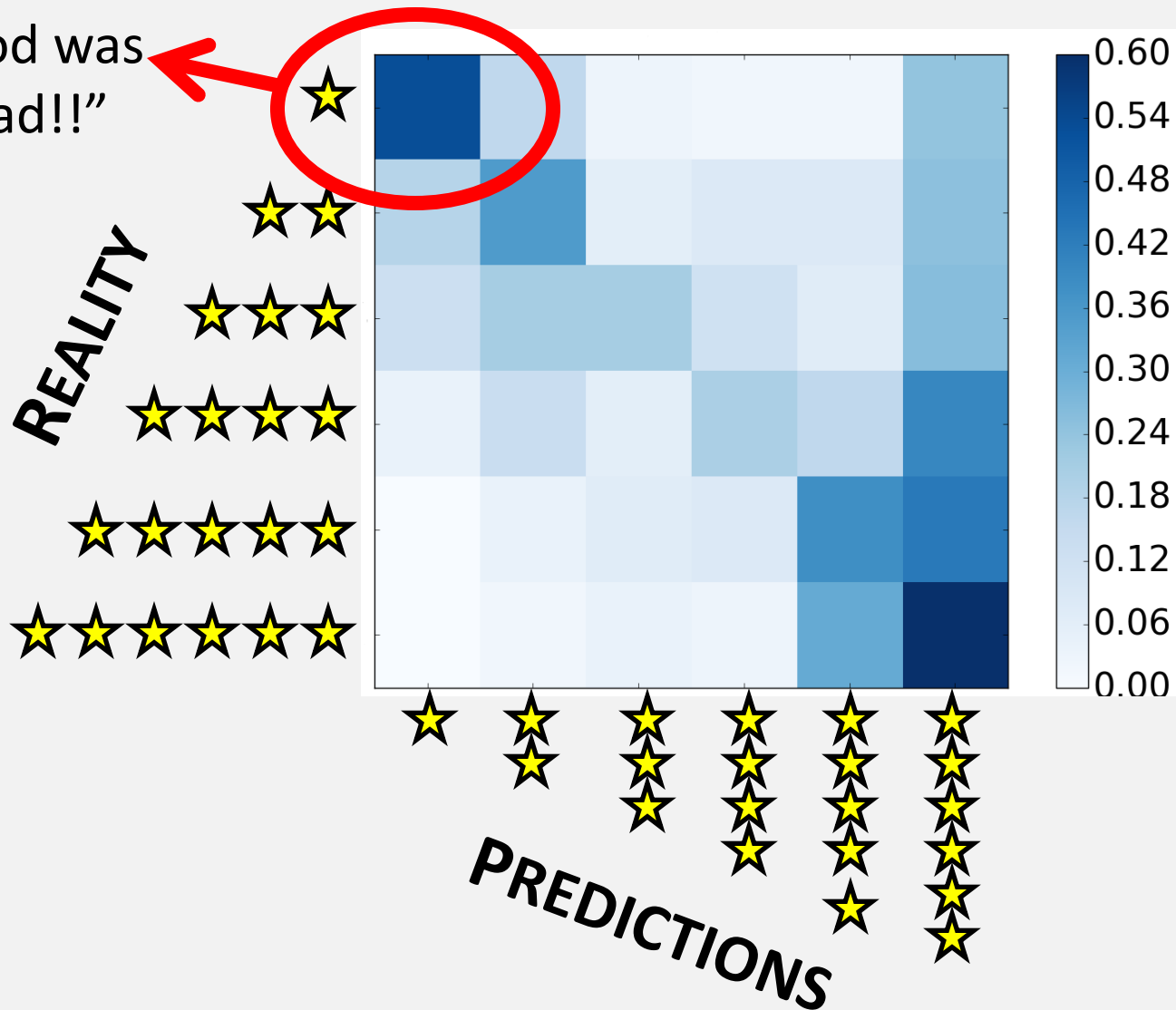
**Bernoulli NB  
Classifier**

# Testing the Classifier



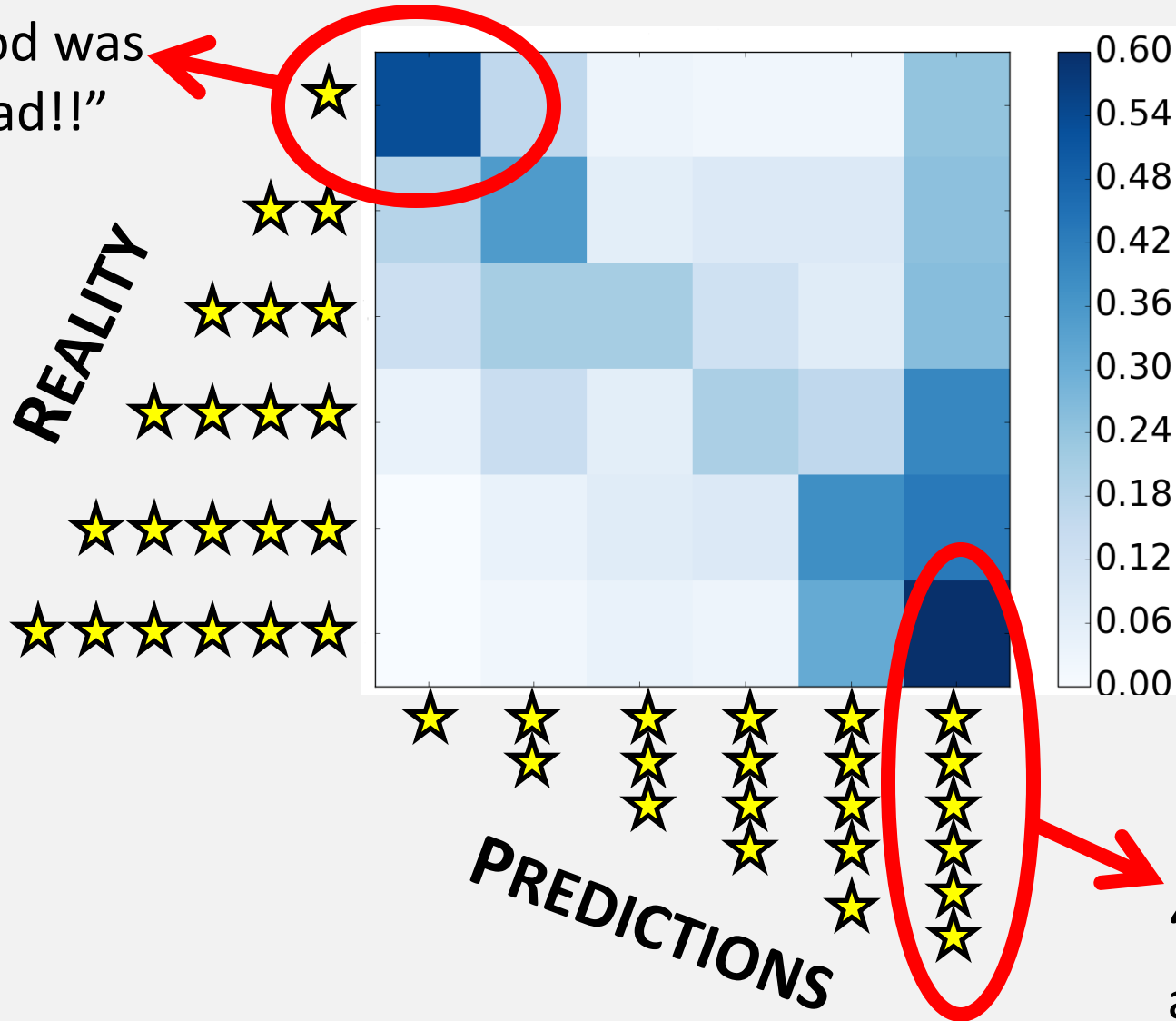
# Testing the Classifier

“the food was  
really bad!!”



# Testing the Classifier

“the food was  
really bad!!”



# Some comparisons

## METHODOLOGY

## Open Access

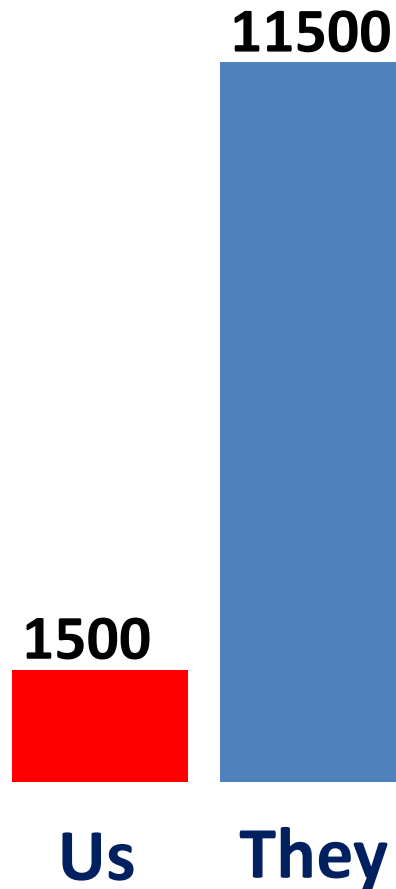
### Sentiment analysis using product review data



Xing Fang\* and Justin Zhan

# Some comparisons

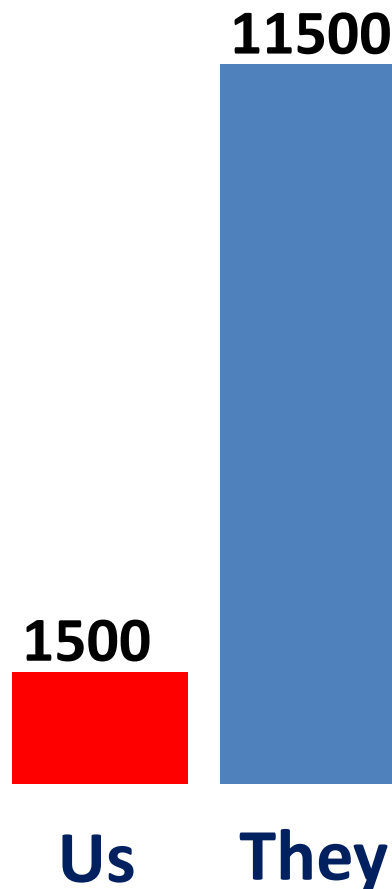
## Number of features



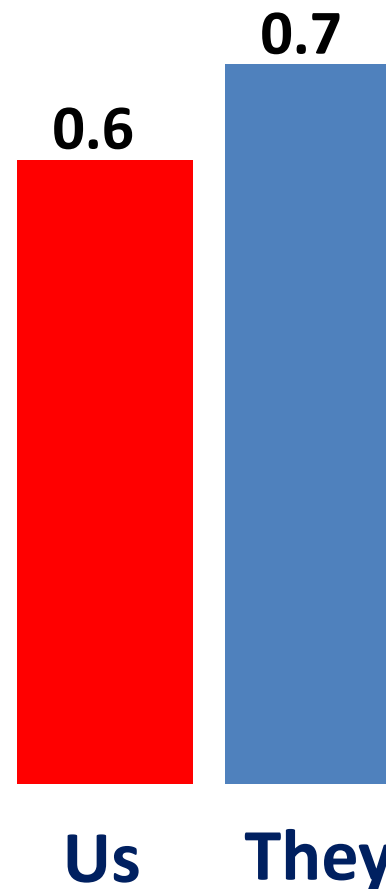


# Some comparisons

## Number of features



## F1 scores



# Improving

Bigger Dictionary



Pipeline



Try different ML model (e.g. SVM)

**Thank you**

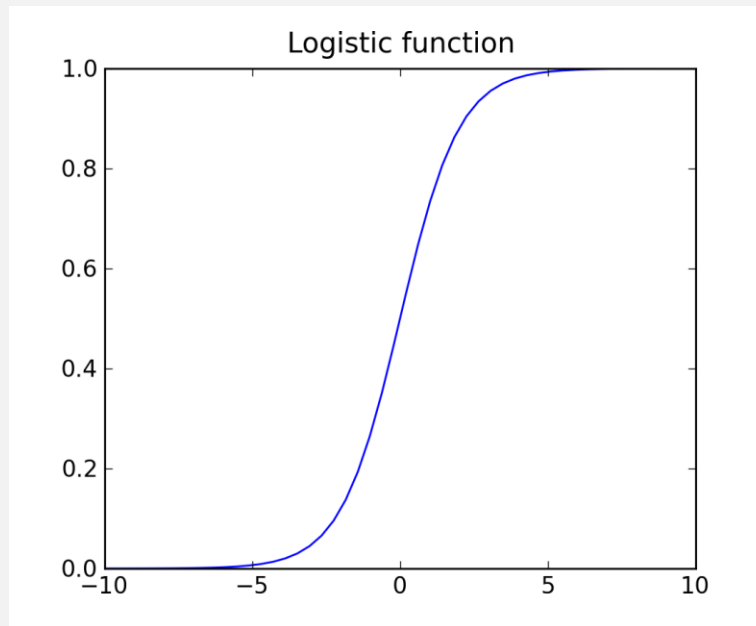
# ML models

**What does all that means from the perspective of probabilities?**

# ML models

## Discriminative

$p(y|x)$  Estimates “directly” the probability of  $y$  given  $x$



$$y = \frac{1}{1 + e^x}$$

$$y = \begin{cases} 0 \\ 1 \end{cases} \quad \text{Classifying in classes '0' and '1'}$$

e.g.,  $y = 1$  if  $x > 0$

# ML models

## Generative

$p(x|y)$       Estimates “directly” the probability of  $y$  given  $x$