# Predicting the Popularity of User Generated Content

Jussara M. Almeida

Department of Computer Science Federal University of Minas Gerais - Brazil

jussara@dcc.ufmg.br

#### User Generated Content (UGC)

- Unprecedented scale and growing rate
  - YouTube: over 72 hs of video uploaded per min more than 4 billion videos watched daily

– Foursquare: over 40 million people worldwide over 4.5 billion check-ins millions more on a daily basis

#### User Generated Content (UGC)

Popularity distribution is highly skewed





#### Most objects attract little attention

#### User Generated Content (UGC)

Popularity distribution is highly skewed

_	ou 🜩		10.15		زگر								
<	Back	Gra	and Central Ma	rket	Share								
<ul> <li> the most beautiful train station in the city. (2 tips)</li> <li>The freshest Swiss Appenzeller cheese I found in NY (2 tips)</li> <li>Don't order the embryos in a glass (trust us, your date would rather you not), but do stroll casually and pick up some snacks throug</li> <li>BRAVO</li> </ul>													
							Don't order the embryos in a glass (trust us, your date would rather you not), but do stroll casually and pick up some snacks throug						
							View All 36 Tips						
	$\sim$		62 67		Tin								
	Friend	ds	Q Explore	M	arisa								



# Ongoing Research

Which factors impact UGC popularity evolution and how can we exploit them to build simple yet reasonably accurate popularity prediction models?

# Popularity Prediction: Why?

- Content distribution services (CDNs, caching)
- Searching services
- Advertising and marketing strategies
- Content filtering, ranking, and recommendation
- Customer feedback (Foursquare tips)
- Understand human dynamics of information consumption processes

# Popularity Prediction: Challenges

- Multitude of factors with potential influence
  - Content itself
  - Social neighborhood or influence zone of user
  - Mechanisms available to drive users to content (search, recommendation, top lists)
  - Specific characteristics of application (ranking by creation time)
  - External factors

Popularity Prediction: Case Studies

• YouTube videos

• Foursquare tips

Popularity Prediction: Case Studies

YouTube videos

• Foursquare tips

#### YouTube Datasets

- Top: 18k videos that appeared on top lists
- YouTomb: 103k videos with copyright violation
- Random: 22k videos selected based on random queries

#### YouTube Datasets



#### YouTube Datasets



# Popularity Evolution: Analysis

- How fast does a video become popular?
- How concentrated is popularity?
- Are there clear popularity trends?
- How content/referrer features correlate with trends?

#### How Fast?

Fraction of time until X% of popularity reached



#### How Fast?

Fraction of time until X% of popularity reached

![](_page_14_Figure_2.jpeg)

50% of videos take at most 65% of lifetime to reach 90% of views

#### How Fast?

![](_page_15_Figure_1.jpeg)

#### For 50% of the videos:

- YouTomb: ≤ 21% of lifetime to reach 90% of views
- Top: ≤ 65% of lifetime for same 90%
- Random: ≤ 87% for same 90%

#### How Concentrated? Fraction of views on peak week 1.0Prob. (Fraction of views $\leq$ f) 70 $0.0 \times 10^{-10}$ f) 70 $0.0 \times 10^{-10}$ f) 3<sup>rd</sup> peak week ...... $2^{nd}$ peak week Peak week 0.0 0.20.80.40.61.0Fraction of views on peak week - f(a) Top

![](_page_17_Figure_0.jpeg)

At least 50% of views on peak week for 60% of videos

#### How Concentrated?

![](_page_18_Figure_1.jpeg)

For 60% of videos, most popular week consists of:

- At least 50% of views for Top
- At least 40% of views for YouTomb
- At least 5% of views for Random

## Are There Popularity Trends?

![](_page_19_Figure_1.jpeg)

KSC Clustering: time shift and scale invariants [Yang2011]

• 4 Clusters in all datasets

# Types of Content per Cluster

#### Fraction of videos per YouTube category

![](_page_20_Figure_2.jpeg)

Distribution within each cluster differs from whole dataset (chi-squared test) 21

## How Do Users Find This Content?

Fraction of views per type of referrer

![](_page_21_Figure_2.jpeg)

Search is very important, but also internal browsing Featured is important for videos in CO and C1 (user retention) Different distributions depending on cluster

- Most previous work: linear regression models [Szabo2008,Pinto2013,Radinsky2012]
  - -Fixed target dates
  - -Fixed monitoring periods

$$\widehat{N}(t_f) = \alpha(t_f, t_r) N(t_r)$$

- Most previous work: linear regression models [Szabo2008,Pinto2013,Radinsky2012]
  - -Fixed target dates
  - -Fixed monitoring periods

$$\widehat{N}(t_{f}) = \alpha(t_{f}, t_{r})N(t_{r})$$
  
Future pularity

- Most previous work: linear regression models [Szabo2008,Pinto2013,Radinsky2012]
  - -Fixed target dates
  - -Fixed monitoring periods

$$\widehat{N}(t_f) = \alpha(t_f, t_r) N(t_r)$$
Current
Popularity

- Most previous work: linear regression models [Szabo2008,Pinto2013,Radinsky2012]
  - -Fixed target dates
  - -Fixed monitoring periods

$$\widehat{N}(t_f) = \alpha(t_f, t_r) N(t_r)$$
Scaling

- Most previous work: linear regression models [Szabo2008,Pinto2013,Radinsky2012]
  - -Fixed target dates
  - -Fixed monitoring periods

$$\widehat{N}(t_f) = \alpha(t_f, t_r) N(t_r)$$

Specialized models: accuracy improvements
 [Pinto2013]

-Predict popularity trend (our goal)

### Popularity Prediction Fixed Monitoring Periods?

![](_page_27_Figure_1.jpeg)

Shortest monitoring period required for accurate predictions varies across videos

How Early Can We Predict the Popularity of a Video?

- Time is not the metric
  - Remaining interest: fraction of remaining views after prediction
- Our solution:
  - Predict popularity trend: classification task
    - Trend = class = cluster
    - Clustering and classification algorithms
  - Prediction accuracy x remaining interest
    - Solves tradeoff on a per-video basis

# Our Popularity Trend Prediction Strategy

- Given a newly uploaded video:
  - Iterate over possible monitoring periods  $t_r$ :
    - Compute probability of video following each trend/class based on early popularity measures (up to  ${\sf t}_{\rm r})$
    - Take largest probability p and associated class  $C_i$ 
      - If p exceeds minimum confidence of  $C_i$ , stop
- Parameters (learned from training set):
  - Minimum and maximum monitoring periods per trend/class
  - Minimum confidence per trend/class.

# Our Popularity Trend Prediction Strategy

![](_page_30_Figure_1.jpeg)

![](_page_31_Figure_0.jpeg)

![](_page_32_Figure_0.jpeg)

![](_page_33_Figure_0.jpeg)

Prob. features: shortest monitoring time for prediction probabilities of video belonging to each class

# Our Popularity Trend Prediction Strategy

![](_page_34_Figure_1.jpeg)

### Trend Learner: Probability Features

- Probability of time series belonging to a trend/class is proportional to inverse distance to centroid of class
  - Distance metric: scale and shift invariants [Yang2011]

![](_page_35_Figure_3.jpeg)

### Trend Learner: Probability Features

- Probability of time series belonging to a trend/class is proportional to inverse distance to centroid of class
  - Distance metric: scale and shift invariants [Yang2011]

![](_page_36_Figure_3.jpeg)

### Trend Learner: Probability Features

- Probability of time series belonging to a trend/class is proportional to inverse distance to centroid of class
  - Distance metric: scale and shift invariants [Yang2011]

![](_page_37_Figure_3.jpeg)

# General Algorithm

1: function TRENDEXTRACTION $(D^{train})$ 

2:  $k \leftarrow 1$ 

- 3: while  $\beta_{CV}$  is not stable do
- 4:  $k \leftarrow k+1$
- 5:  $\mathbf{C}_D \leftarrow KSC(D^{train}, k)$
- 6: end while
- 7: Store centroids in  $C_D$
- 8: end function
- 9: function TRENDLEARNER( $\mathbf{C}_D, D^{train}, D^{test}$ )
- 10:  $\boldsymbol{\theta}, \boldsymbol{\gamma}, \mathbf{P}^{train} \leftarrow LearnParams(D^{train}, \mathbf{C}_D)$ 
  - $\triangleright$  learn classifier from  $\mathbf{P}^{train}$  and object features of  $D^{train}$
- 11:  $TrainERTrees(D^{train}, \mathbf{P}^{train} \bigcup \text{ object features})$
- 12:  $\mathbf{t}, \mathbf{P} \leftarrow MultiClassProbs(D^{test}, \mathbf{C}_D, \boldsymbol{\theta}, \boldsymbol{\gamma})$
- 13: **return t**,  $PredictERTrees(D^{test}, \mathbf{P} \bigcup object features)$

14: end function

# General Algorithm

![](_page_39_Figure_1.jpeg)

## General Algorithm

![](_page_40_Figure_1.jpeg)

### Experimental Evaluation: F1 Metric

Use probabilities directly: take largest probability					
			Cla	ssifier	
		Р	$\mathbf{P} + \mathbf{ERTree}$	ERTree	TrendLearner
Тор	Micro F1 Macro F1	$.48 \pm .06$ $.44 \pm .06$	$.48 \pm .06$ $.44 \pm .06$	$.58 \pm .01$ $.57 \pm .01$	$.62 \pm .01$ $.61 \pm .01$
Random	Micro F1 Macro F1	$.67 \pm .02 \\ .69 \pm .02$	$.62 \pm .01$ $.63 \pm .01$	$.65 \pm .01 \\ .63 \pm .01$	$.71 \pm .01$ $.70 \pm .01$

Experimental Evaluation: F1 Metric					
Use probabilities as features of ERTree					ilities as ERTree
	Classifier				
		Р	$\mathbf{P} + \mathbf{ERTree}$	ERTree	TrendLearner
Тор	Micro F1 Macro F1	$.48 \pm .06$ $.44 \pm .06$	$.48 \pm .06$ $.44 \pm .06$	$.58 \pm .01$ $.57 \pm .01$	$.62 \pm .01$ $.61 \pm .01$
Random	Micro F1 Macro F1	$.67 \pm .02 \\ .69 \pm .02$	$.62 \pm .01$ $.63 \pm .01$	$.65 \pm .01 \\ .63 \pm .01$	$.71 \pm .01$ $.70 \pm .01$

Experimental Evaluation: F1 Metric					
Use only object features					
		Classifier			
		Р	$\mathbf{P}$ +ERTree	ERTree	TrendLearner
Top	Micro F1 Macro F1	$.48 \pm .06 \\ .44 \pm .06$	$.48 \pm .06 \\ .44 \pm .06$	$.58 \pm .01 \\ .57 \pm .01$	$.62 \pm .01$ $.61 \pm .01$
Random	Micro F1 Macro F1	$.67 \pm .02 \\ .69 \pm .02$	$.62 \pm .01 \\ .63 \pm .01$	$.65 \pm .01 \\ .63 \pm .01$	$.71 \pm .01$ $.70 \pm .01$

Experimental Evaluation: F1 Metric					
Our solut				lution	
			Cla	ssifier	
		Р	$\mathbf{P}+\mathbf{ERTree}$	ERTree	TrendLearner
Тор	Micro F1 Macro F1	$.48 \pm .06$ $.44 \pm .06$	$.48 \pm .06$ $.44 \pm .06$	$.58 \pm .01$ $.57 \pm .01$	$.62 \pm .01$ $.61 \pm .01$
Random	Micro F1 Macro F1	$.67 \pm .02 \\ .69 \pm .02$	$.62 \pm .01$ $.63 \pm .01$	$.65 \pm .01 \\ .63 \pm .01$	$.71 \pm .01$ $.70 \pm .01$

### Results: Summary

- Promising results: accuracy and remaining interest
  - median of 68% of views remaining after prediction (Top)
  - median of 32% of views remaining after prediction (Random)
- Specialized models to predict popularity at future date
  - 1. Predict popularity trend
  - 2. Apply specialized regression model for predicted trend
  - $\geq$  improve prediction accuracy by 30% (median)

Popularity Prediction: Case Studies

YouTube videos

• Foursquare tips

# Foursquare: Tips and Likes

![](_page_47_Picture_1.jpeg)

Tip = Micro-review

Popularity = total number of likes Popularity ≈ helpfulness ≈ quality<sup>8</sup>

## Reviews v.s. Micro-Reviews

#### Epinions 😜 😑 😮

![](_page_48_Picture_2.jpeg)

#### Canon EOS 1100D / Rebel T3 Digital Camera ★ ★ ★ ★ 11 consumer reviews Average Rating: Excellent

![](_page_48_Figure_4.jpeg)

#### 2 stars 1 star

Absolutely amazing camera 

User Rating:	Excellent	Pro
Ease of Use:		Co
Durability:		Th
Battery Life:		
Photo Quality:		
Shutter Lag		

os: great photo quality ons: The LCD feels like it's blurry sometimes e Bottom Line: A camera of high guality

Brought on Amazon.com,this camera is absolutely amazing! My boyfriend has the T4i and although that has some pretty interesting features. I can't help but to love this baby more. I went to a hockey game the other night and even though we sat pretty far away the pictures that I took with this came out excellent! I find this camera to be more lightweight than the Nikon 3100 (I used to rent out my University's camera) and easier to use as well. But that's just my opinion! I'm a Canon girl through and through.

Also, comparing it still to the 3100. I find the T3 has sharper, brighter colors especially when you toy around with the "Creative Auto" setting and change up some of the settings there. I believe the grip for the 3100 was rubberized and I will admit that it does take a little adjusting to for the non-rubberized grip, but now after using it for almost two months it's just something natural to me and whenever I use my boyfriend's T4i, the rubberized grip just feels weird to me now! The battery life, for me and how I use it, is excellent! I even left the switch "on" for a whole weekend by accident after taking pictures all day, and picked it right back up and continued to shoot more pictures. I've only had to charge it two or three times so far, and I tend to take at least 20 to 30 pictures a day at the very least. Some days I'll take up to 300 or more!

This camera was the best purchase I've made, and I am not disappointed one bit by the performance, the body, and the picture quality!

#### foursquare

![](_page_48_Figure_13.jpeg)

the bagged chocolate chip cookies are addictive! olive rolls are delicious.

Save 🖤 Like

# Tip Popularity Prediction Tasks

- Predict the popularity of a tip at a given future date
- Predict the *popularity level* of a tip at a given future date
- Rank tips based on predicted popularity at future date

# Tip Popularity Prediction Tasks

- Predict the popularity of a tip at a given future date
- Predict the *popularity level* of a tip at a given future date
  - Classification task
- Rank tips based on predicted popularity at future date

# Popularity Levels

Popularity level	# of likes/tip	# tips
Low	< 5	703,827
High	≥ 5	3,427

- High class imbalance: severe impact on prediction!!!
- Similar for other definitions of popularity levels

# Problem Statement

- A tip is a tuple (p,u,v) where
  - p: features extracted from the tip's content
  - -u: features associated with the tip's author
  - -v: features of venue where tip was posted

![](_page_52_Figure_5.jpeg)

# Prediction Algorithms

- Support Vector Machine (SVM) classifier
  - Linear e RBF kernels
- Regression
  - Support Vector Regression (SVR)
    - Linear e RBF kernels
  - Simpler linear regression (OLS)
- Median number of likes of tips previously posted by the user (baseline)

# Features: Tip's Author

- Number of tips
- Number of likes received
- Number of likes given
- Number of distinct venues
- Number of mayorships
- Is user mayor of venue?
- Number of friends/followers
- Number of likes from SN
- Number of tips posted by SN
- Number of likes given by SN
- Visibility of user in venue
- Type of user

Activities in the system

Social network

## Features - Venue

- Number of tips
- Number of likes
- Number of checkins
- Number of unique visitors
- Is venue verified?
- Venue category
- Position of tip in ranking by # likes
- Position in ranking by posting time

Activities

Characteristics

#### Others

# Features - Tip

- Number of characters
- Number of words
- Number of urls or emails
- % nouns
- % adjectives
- % adverbs
- % verbs
- % punctuation marks
- Positive scores (average)
- Negative scores (average)
- Neutral scores (average)

Amount of content

#### Part-of-speech tags

#### Polarity/Sentiment

#### Experimental Evaluation: Macro-Recall

![](_page_57_Figure_1.jpeg)

Best approach: combine user + venue features OLS, SVM and SVR: similar results (but OLS is simpler) 58

#### Experimental Evaluation: Recall of Low Popularity

![](_page_58_Figure_1.jpeg)

Median of likes: best results (class imbalance)

#### Experimental Evaluation: Recall of High Popularity

![](_page_59_Figure_1.jpeg)

SVR (RBF): slightly better than OLS and SVM

# Results: Summary

- Best results:
  - OLS: similar to SVM and SVR, but simpler
  - combination of user and venue features as inputs
- Ranking of features by importance (Information Gain)
  - Top-3 related to user: number of likes in previous tips
  - $-4^{\text{th}}$ : size of social network of user
  - 6<sup>th</sup>-7<sup>th</sup>: popularity of venue (# visitors, # checkins)

### **Important User Features**

![](_page_61_Figure_1.jpeg)

Average # likes

#### # friends/followers

# Summary

- Popularity prediction of UGC: challenging task
  - Multitude of external and internal factors
  - Inherent characteristics of application
  - Highly skewed popularity distribution: severe imbalance impacts efficacy of regression/classification
- Current work:
  - YouTube: popularity trends and measures
  - Foursquare: popularity measures, levels and ranking
- Next steps: new features
  - YouTube: user and social network, content
  - Foursquare: geographical aspects

# Foursquare Dataset

- Almost 7 million tips
- 5,7 million likes
- 1,8 million users
- 3,2 million venues