

# Improving Users' Demographic Prediction via the Videos They Talk about

**Yuan Wang**, Yang Xiao, Chao Ma, and Zhen Xiao  
Peking University, China  
EMNLP2016, 3 November 2016

# Table of Contents

- Introduction
- Data
- Indirect Relationships between Users and Videos
- Evaluation
- Conclusion

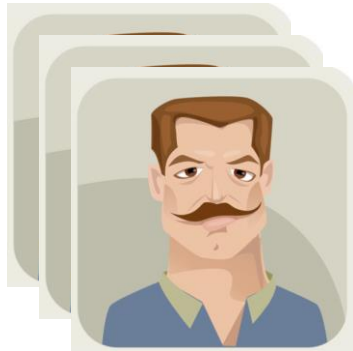
# Introduction



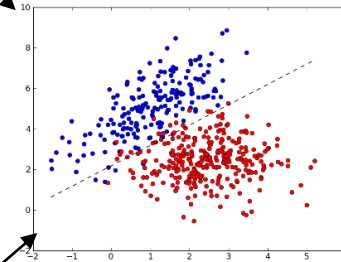
# Introduction



Love, pretty,  
work, ...



Sports, Financial,  
..., Car



Female

Male

# Introduction



TOP1 #我是歌手# 综艺

★ video

#我是歌手# 第四季 V 经纪招募，首唱淘汰赛制升级，到底谁去谁留呢？

阅读数：232.4亿 主持人：我是歌手



TOP2 #太阳的后裔# 电视剧

★ video

由宋仲基、宋慧乔主演的《太阳的后裔》主要讲述了特战部队海外派兵

阅读数：105.7亿 主持人：爱奇艺



TOP3 #二十四小时# 综艺

★ video

斯柯达速派【#二十四小时# 大结局】6位时空水手@陈坤 @徐峥 @韩

阅读数：14.9亿 主持人：浙江卫视二十四小时



4 #0408EXO出道四周年# 明星

EXO是韩国SM公司于2012年4月8日推出的韩国男子组合，分成EXO-

阅读数：2.4亿 主持人：EXO吧微博



5 #最好的我们# 电视剧

★ video

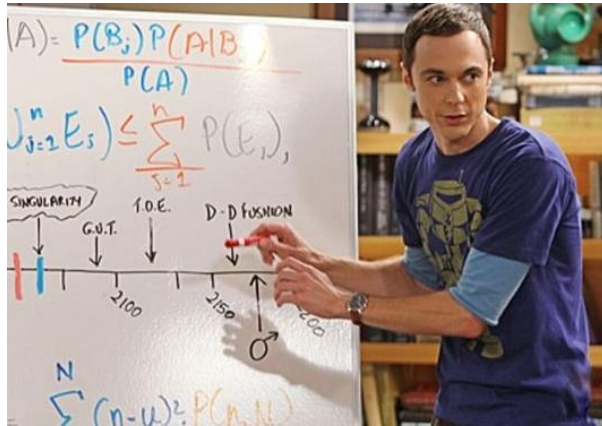
总有一人，是你耿耿于怀的青春。八月长安同名小说改编网络剧《最好

阅读数：2.1亿 主持人：爱奇艺

- Mean Girls, Pretty Woman, The Devil Wears Prada
- House of Cards, Mission Impossible, NBA



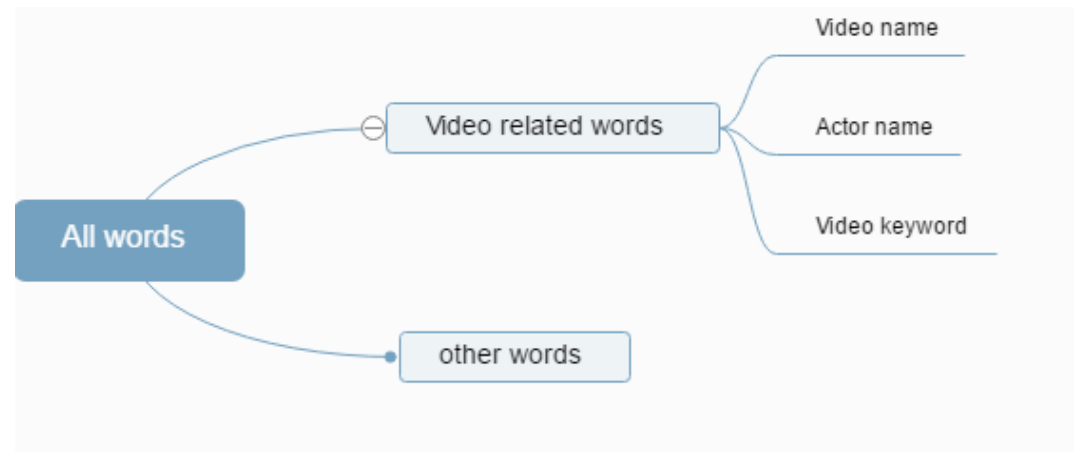
# Introduction



“Will the **Big Bang Theory** last into the next century?”

“**Sheldon** is so cool, I love him!”

“**Jim Parsons** was nominated for another Emmy Award”



# Introduction

YOUKU 优酷

IQIYI 爱奇艺

腾讯视频  
V.QQ.COM

搜狐 视频

电影: 星球大战7:原力觉醒 2016

1



评分: ★★★★★ 9.2 豆★ 7.1

别名: Star Wars:The Force Awa... 时长: 138分钟

上映: 2016-01-09 优酷上映: 2016-04-18

地区: 美国

导演: J.J.艾布拉姆斯

类型: 科幻/冒险/动作

主演: 哈里森·福特 / 马克·哈米尔

总播放: 35,848,818

评论: 22,560

顶: 32,210

指数: ▼

2

免费试看

★收藏 ↓下载 手机看 分享

概况

获奖

预告片

花絮

粉丝自制

MV

首映式

豆瓣

卢卡斯, 如是说

请看完电影后再来看影评。-----首先要声明的是, 这片子在制作和创作层面跟卢卡斯一点关系都没有了。卢卡斯在一次采访中的原话如下: "The ones that I sol..."

查看全文

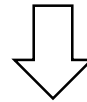
3

《星球大战: 原力觉醒》彩蛋和花絮总汇

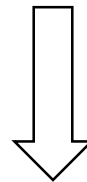
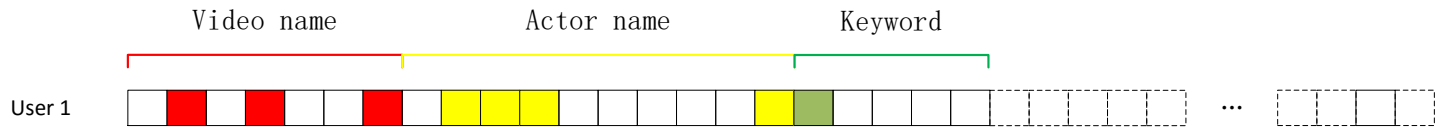
三刷影片后, 确认了不少彩蛋, 下面尽量按照时间顺序排列。2187 Finn的暴风兵编号为FN-2187, 2187是《新希望》中帝国军关押Leia公主的牢房号码。这个数字本来就是一个彩蛋, 《21-87》是由加拿大先锋导演Arthur...

查看全文

# Introduction

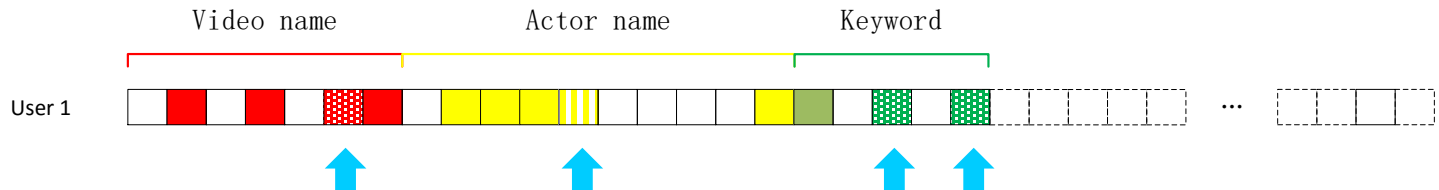


With the help of



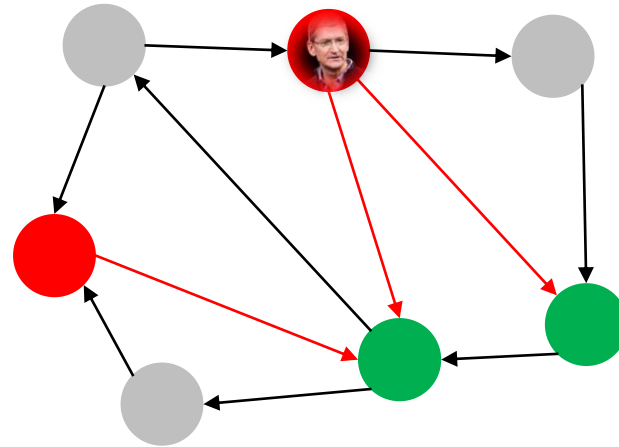
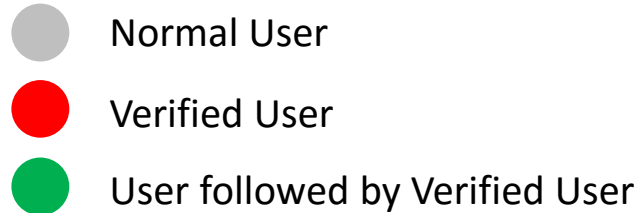
With the help of

???





# Data



Attribute	Completion Rate	Categories
Gender	95.019%	Male, Female
Age	18.604%	Teenage (<18), Youngster (18-24), Young (25-34), Mid-age(>34)
Education BG	17.443%	University, Non-University
Marital Status	2.203%	Single, Non-Single

**Table 1:** Demographic attributes and corresponding categories

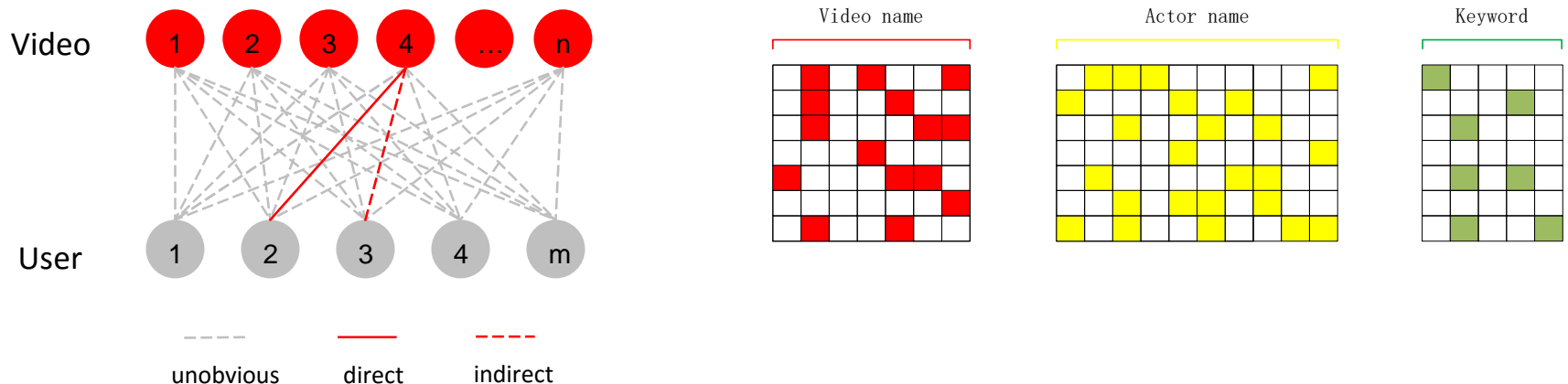
# Data

	<b>Video</b>	<b>Actor</b>	<b>Keyword</b>
Variety show	344	1007	2925
Movie	306	741	2049
TV	197	515	1302
Total	<b>847</b>	1422	4094

**Table 2:** Statics of video relevant information (There is an overlap between the three collections of actors and keywords.)

# Discover Indirect Relationships

- Unobvious relationship
  - $N * M$  pairs
- Direct relationship
  - User 2, “Will the Big Bang Theory last into the next century?”
- Indirect relationship
  - User 3 posts, “Sheldon is so cool, I love him!”



# Discover Indirect Relationships

## Step 1

$$P(v_n) = \frac{\text{num}(\text{users watched the } n_{th} \text{ video})}{\text{num}(\text{users})}$$

$$P(w_{ni}|v_n) = \frac{\text{num}(\text{users watched the } n_{th} \text{ video and mentioned the } ni_{th} \text{ keyword})}{\text{num}(\text{users watched the } n_{th} \text{ video})}$$

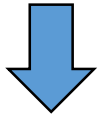
$$P(a_{nj}|v_n) = \frac{\text{num}(\text{users watched the } n_{th} \text{ video and mentioned the } nj_{th} \text{ actor})}{\text{num}(\text{users watched the } n_{th} \text{ video})}$$

## Step 2

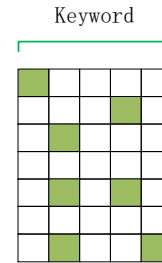
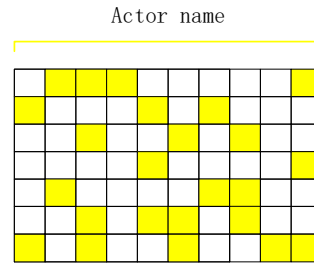
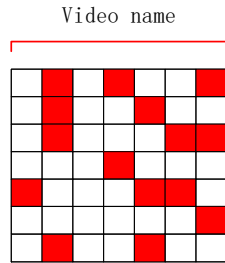
$$\begin{aligned} P(v_n|W_m, A_k) &= \frac{P(W_m, A_k|v_n) * P(v_n)}{P(W_m, A_k)} \\ &= \frac{\prod_{w_{ni} \in W_m} P(w_{ni}|v_n) * \prod_{a_{nj} \in A_k} P(a_{nj}|v_n) * P(v_n)}{P(W_m, A_k)} \end{aligned}$$

# Discover Indirect Relationships

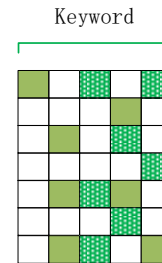
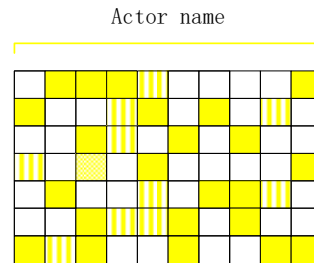
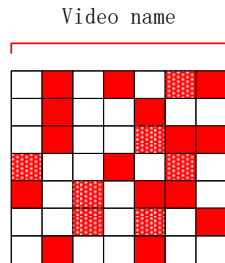
Direct



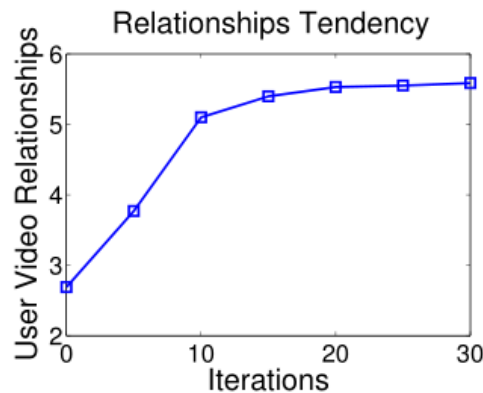
Direct + Indirect  
(more denser)



Two Baseline Model

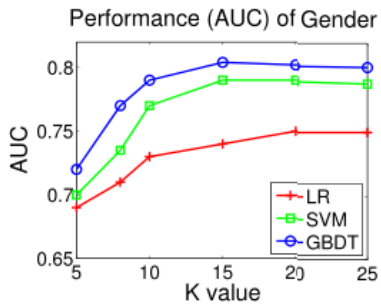


Two Indirect Relationship  
Based Model

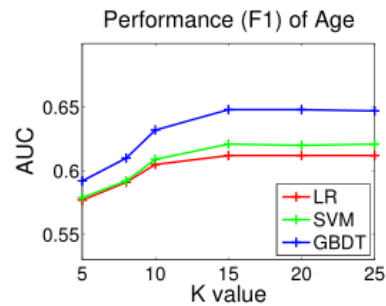


# Discriminant Model

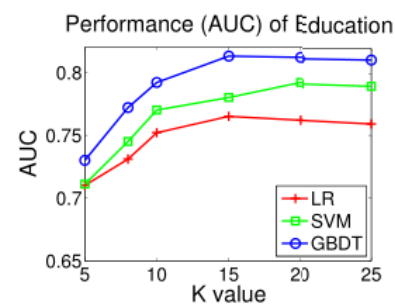
- Matrix Factorization<sup>1</sup>, **K=20**
- LR<sup>2</sup>, SVM<sup>2</sup>, **GBDT<sup>3</sup>**



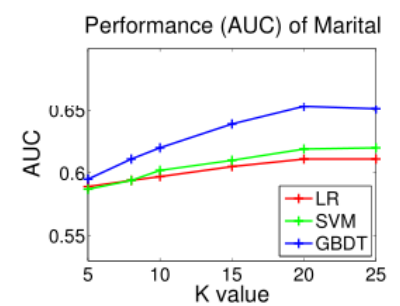
(a) AUC of Gender



(b) F1 of Age



(c) AUC of Education BG



(d) AUC of Marital Status

- 1 libFFM
- 2 liblinear
- 3 XGBoost

# Generative Model

- Calculate video demographic tendency
- Calculate user demographic attribute
- Smooth the result

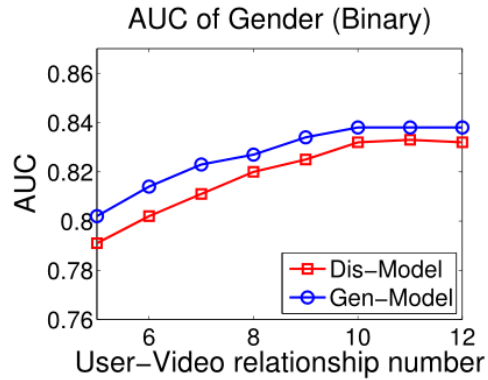
# Evaluation

		<b>Precision</b>	<b>Recall</b>	<b>F1</b>	<b>AUC</b>
Gender 1	Dis-Baseline	0.720	0.714	0.717	0.730
	<b>Dis-Model</b>	0.786	0.779	0.783	<b>0.812</b> ↑ 11.2%
	Gen-Baseline	0.701	0.687	0.694	0.707
	<b>Gen-Model</b>	0.799	0.802	0.801	<b>0.825</b> ↑ 16.7%
Age 2	Dis-Baseline	0.569	0.541	0.554	*
	<b>Dis-Model</b>	0.642	0.653	<b>0.648</b> ↑ 16.8%	*
	Gen-Baseline	0.529	0.504	0.516	*
	<b>Gen-Model</b>	0.663	0.645	<b>0.654</b> ↑ 26.7%	*
Education BG 3	Dis-Baseline	0.707	0.716	0.711	0.730
	<b>Dis-Model</b>	0.788	0.801	0.795	<b>0.809</b> ↑ 11.1%
	Gen-Baseline	0.680	0.659	0.669	0.690
	<b>Gen-Model</b>	0.790	0.808	0.799	<b>0.812</b> ↑ 17.7%
Marital Status 4	Dis-Baseline	0.565	0.549	0.557	0.571
	<b>Dis-Model</b>	0.657	0.640	0.648	<b>0.659</b> ↑ 15.4%
	Gen-Baseline	0.572	0.550	0.560	0.581
	<b>Gen-Model</b>	0.682	0.691	0.687	<b>0.696</b> ↑ 19.8%

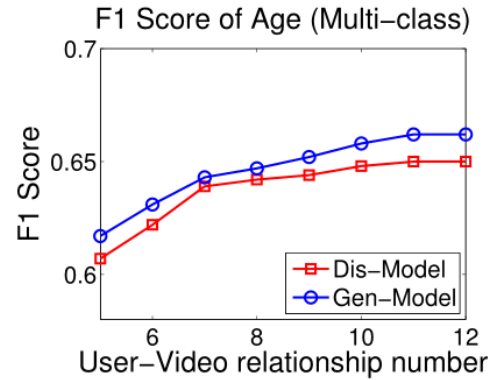
**Table 3:** Prediction accuracy based on users' video describing words. Classes have been balanced.



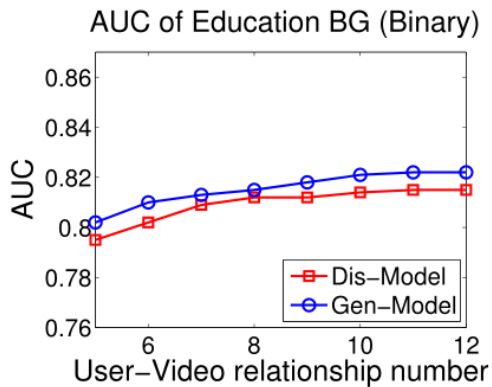
# Evaluation



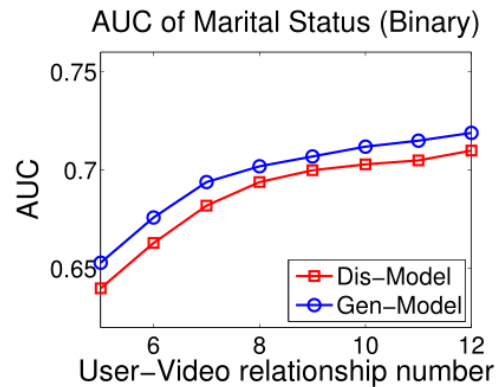
(a) AUC of Gender



(b) F1 Score of Age

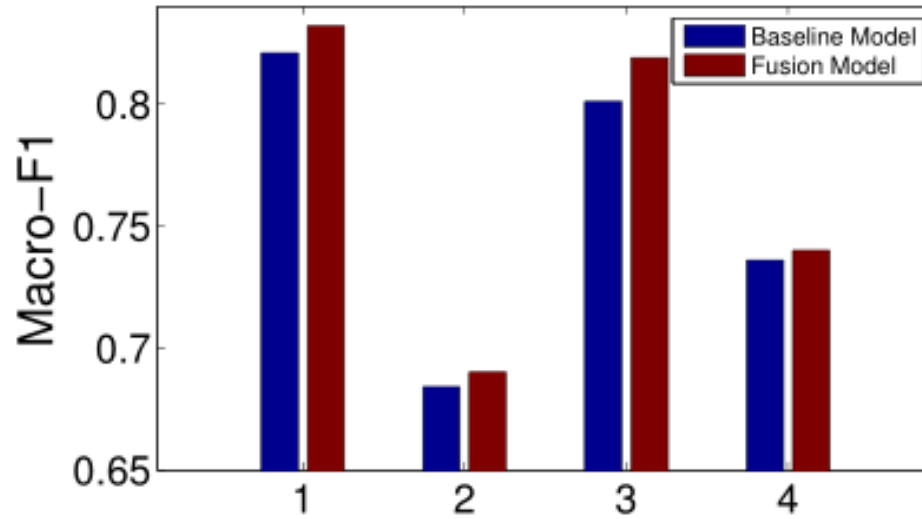


(c) AUC of Education BG



(d) AUC of Marital Status

# Evaluation



**Figure 5:** Results of Fusion Model evaluation (Macro-F1).

# Conclusion

- Our motivation is that user's video related behavior is usually under-utilized on demographic prediction tasks.
- With the help of third-party video sites, we detect the direct and indirect relationships between users and video describing words, and demonstrate this effort can improve the accuracy of users' demographic predictions.
- To our knowledge, this is the first work which explores demographic prediction by fully using users' video describing words.
- This framework has good scalability and can be applied on other concrete features, such as user's book reading behaviors and music listening behaviors.

Thanks!