

# Text Mining Applications

KhabarFarsi - ParsPack

پارس پک  
PARSPACK



سایت هوشمند خبری



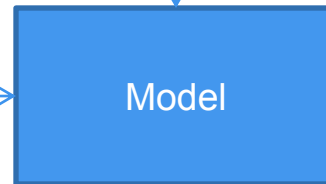
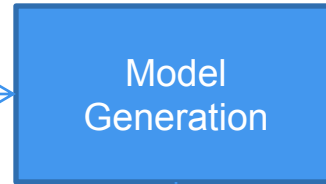
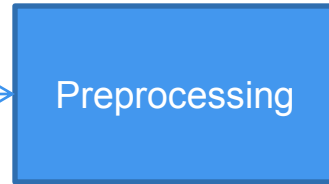
# What is Text Mining?

احساس	توییت
منفی	واقعا همراه چهارم گند زده با این خدماتش. یه روز در میون اینترنتش می‌پره
مثبت	همراه چهارم نسبت به قیمتش میرزه و خیلی هم بد نیست
منفی	این همراه چهارم هر هفته یه قطعی چیزی داره
خنثی	نمیدونم از همراه چهارم استفاده کنم یا همراه پنجم



# Whole Process

توییت	اح سا س
واقعا همراه چهارم گند زده با این خدماتش. یه روز در میون اینترنتش می‌پره	مذ فی
همراه چهارم نسبت به قیمتش میرزه و خیلی هم بد نیست	مذ بت
این همراه چهارم هر هفته به قطعی چیزی داره	مذ فی
نمیدونم از همراه چهارم استفاده کنم یا همراه پنجم	خذ تی

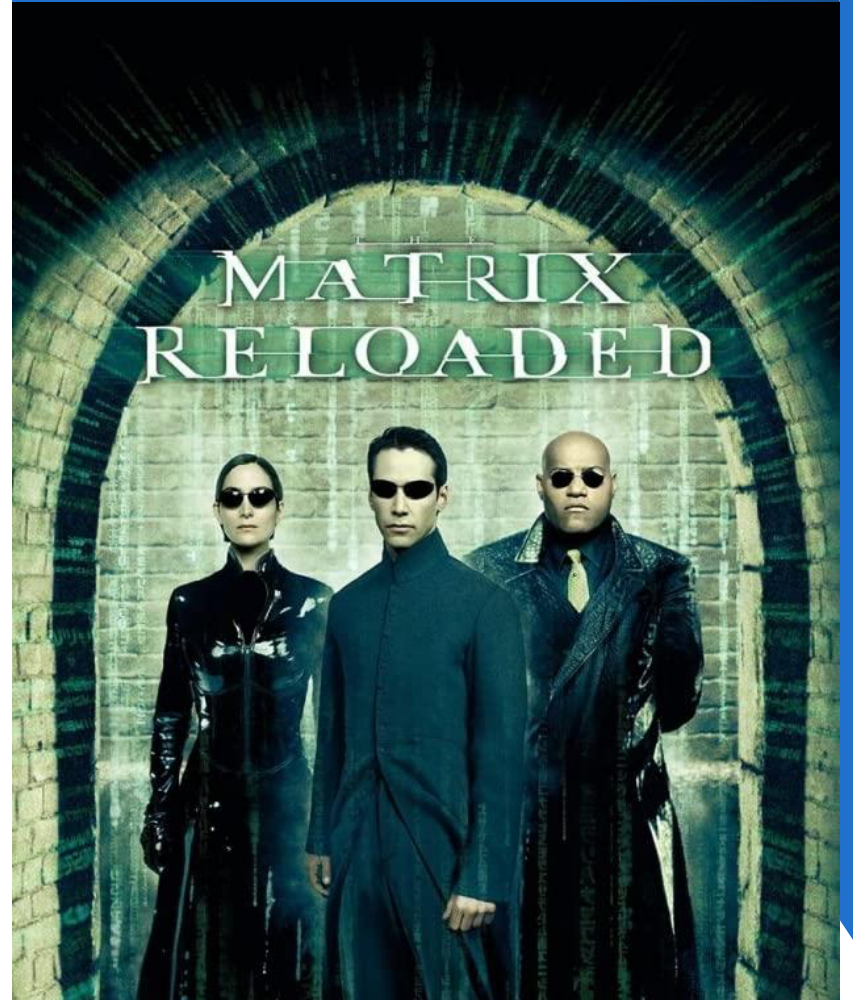


مثبت

همراه چهارم با همراه سوم قابل رقابته

# Matrix!

**Convert Data into  
Computer Understandable  
Format**





# Convert Text to Matrix (Feature Engineering)

- ▶ Bag of Words
- ▶ TF-IDF
- ▶ Word2Vec (Doc2Vec)
- ▶ LDA and LSI
- ▶ Glove, Bert, FastText and...



**KhabarFarsi**

استانها - پیش بینی آب و هوا - تبلیغات - صوت و ویدیو - سایت های خبری - جدیدترین خبرها - RSS

سه شنبه ۰۷ مرداد ساعت ۲۰:۴۴ به وقت تهران < مصادف با 28 Jul 2020



جستجو کنید ...



سایت هوشمند خبری

خبر فارسی یک موتور جستجوگر اینترنتی است و هیچ دخالت انسانی در دریافت و چینش اخبار وجود ندارد.

خانه سیاسی ورزشی اقتصادی فرهنگ و هنر علمی و پزشکی فناوری و ارتباطات اجتماعی بیر

برچسب‌های خبری

کرونا جمهوری اسلامی ایران خبر ایران سخنگوی وزارت بهداشت امام جمعه جمهوری اسلامی وزارت شنبه  
برگزار وزارت بهداشت تیم ملی فدراسیون فوتبال

## مهمترین اخبار

جدیدترین

۲ ساعت گذشته

۱۲ ساعت گذشته

۲۴ ساعت گذشته

## واکنش نیروی دریایی آمریکا به رزمایش سپاه

برترین ها ۲ ساعت پیش + ۱۰ سایت دیگر

خبرگزاری فارس: نیروی دریایی آمریکا به برگزاری مرحله نهایی رزمایش بزرگ دریایی و موشکی پیامبر اعظم (ص) واکنش نشان داد. ...



🔴 سرلشکر سلامی: توسعه تسلیحات ما منطبق بر شناخت نقاط ضعف و قوت دشمن است

خبرگزاری فارس / ۴ ساعت پیش

🔴 رزمایش پیامبر اعظم 14 | تصویربرداری ماهواره نور در صدمین روز ماموریت/ نسل ...

خبرگزاری فارس / ۴ ساعت پیش

## روحانی: با وثیقه سهام عدالت، کارت اعتباری بگیرید

برترین ها ۱ ساعت پیش + ۶۶ سایت دیگر

خبرگزاری فارس: حجت الاسلام حسن روحانی رئیس جمهور عصر امروز در جلسه شورای عالی فناوری اطلاعات گفت: وزرا و رؤسای دستگاه های اجرایی مسئول تحقق کامل دولت الکترونیکی هستند و شورای





## جدیدترین خبرها

## ورزشی

غیبشآوری: برای اسکوچیچ در اردوی نساجی اتاق بگیرند/دوره ...  
 خبرگزاری فارس ۷ دقیقه پیش

توسعه ورزش آپارتمانی برنامه ریزی و در همدان دنبال می شود  
 خبرگزاری مهر ۴ دقیقه پیش

پرسپولیس پس از دو روز استراحت تمرین کرد  
 ایرنا ۱۰ دقیقه پیش

تبریک سالگرد قهرمانی آسیا با طعنه به پرسپولیس/عکس  
 خبر آنلاین ۱۸ دقیقه پیش

## سیاسی

امام جمعه اصفهان: تبعیت از ولایت فقیه، وجه تمایز قوای ایرا...  
 ایرنا ۳ دقیقه پیش

از حربه آمریکا برای منع همکاری ایران و چین تا داستان وزیر...  
 باشگاه خبرنگاران ۷ دقیقه پیش

گفتمان مقاومت می تواند در شکل دهی نظام آینده جهان جای...  
 ایرنا ۱۵ دقیقه پیش

همکاری آمریکا با عربستان برای جاسوسی از مخالفان ریاض  
 ایرنا ۲۵ دقیقه پیش

# 120,000

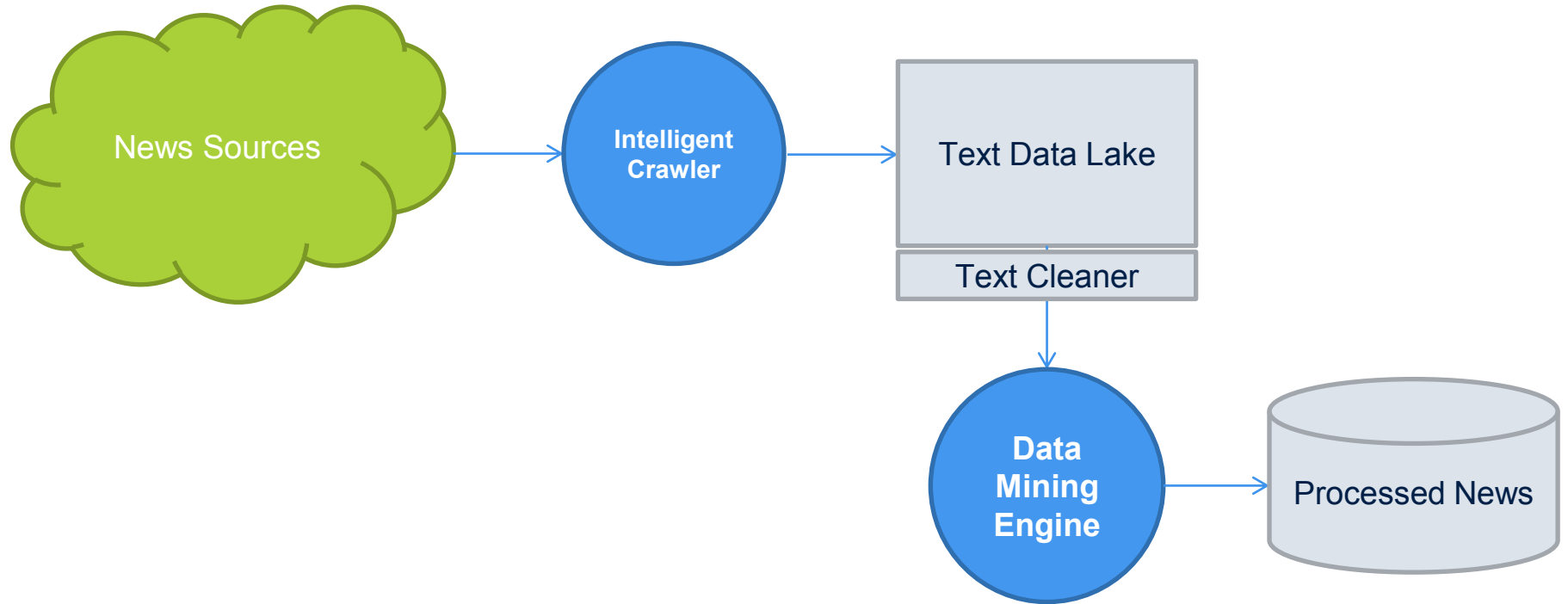
News Every Day

# Glosseries

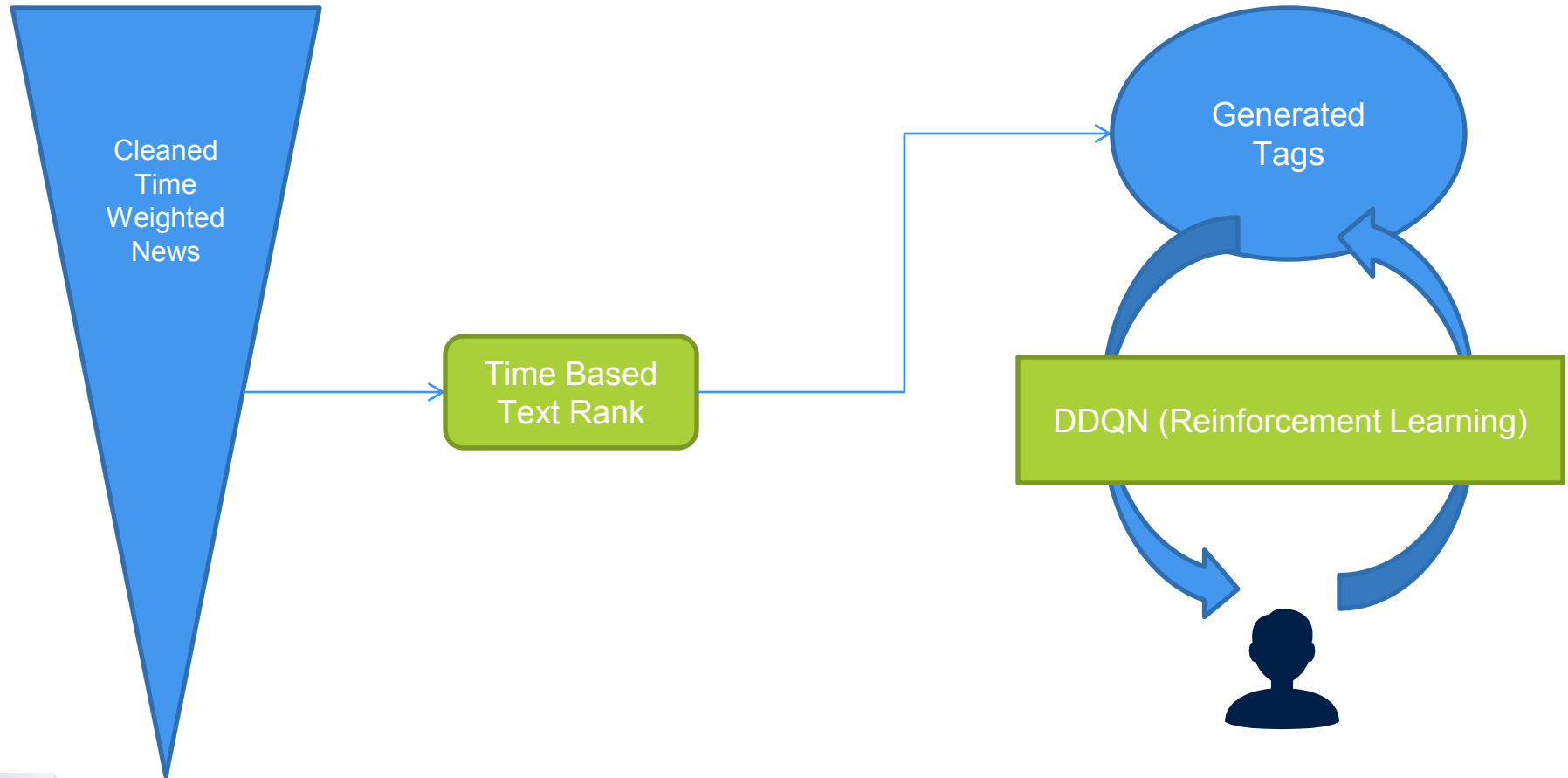


- ▶ Text Classification
- ▶ Text Clustering
- ▶ TextRank (Text Summarization)
- ▶ Topic Modeling
- ▶ Similarity Detection

.



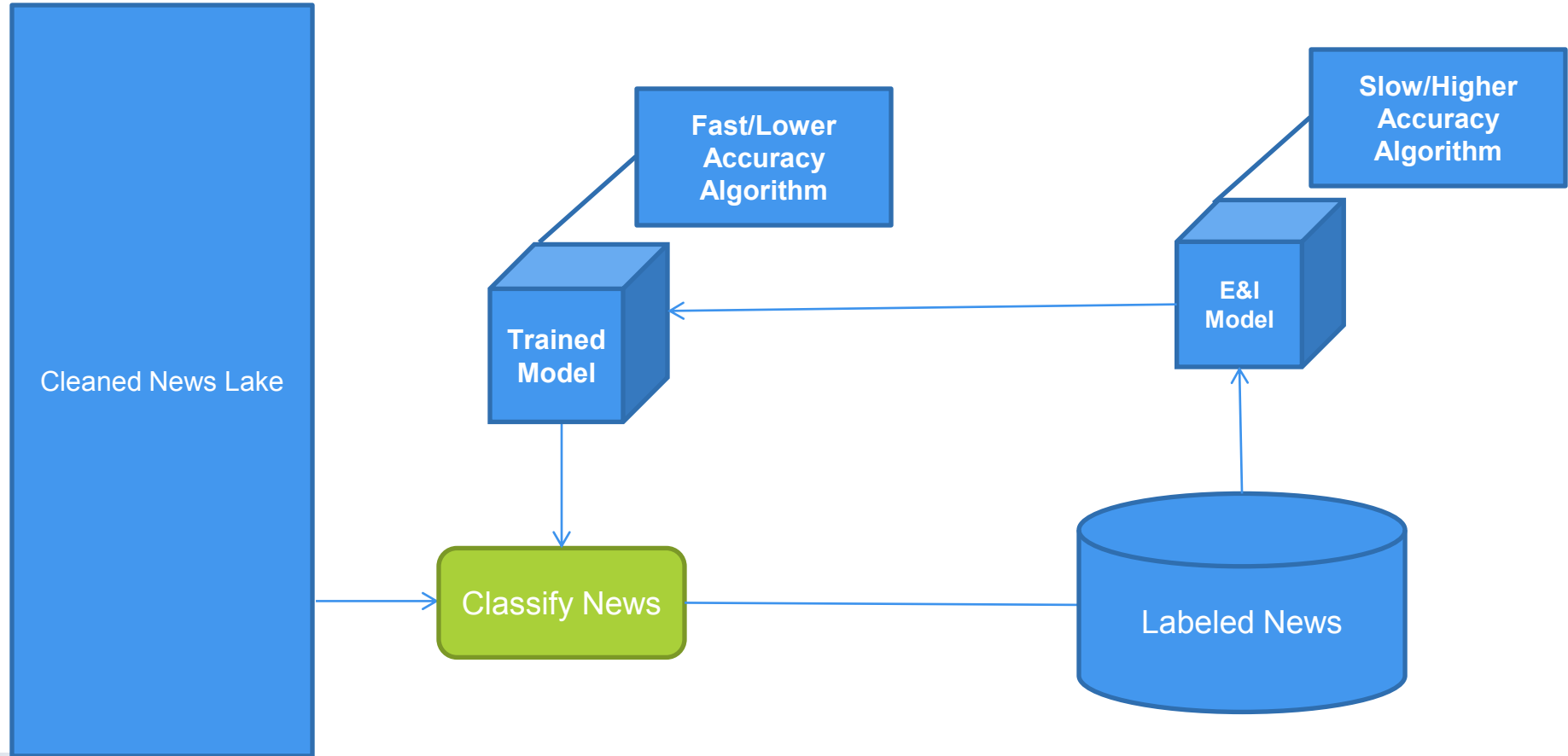
# Model 02 – Popular Tags



// >> Weighted Text Rank  
>> Double Deep Q  
Reinforcement Learning



# Model 03 – News Classification



“

>> DNN

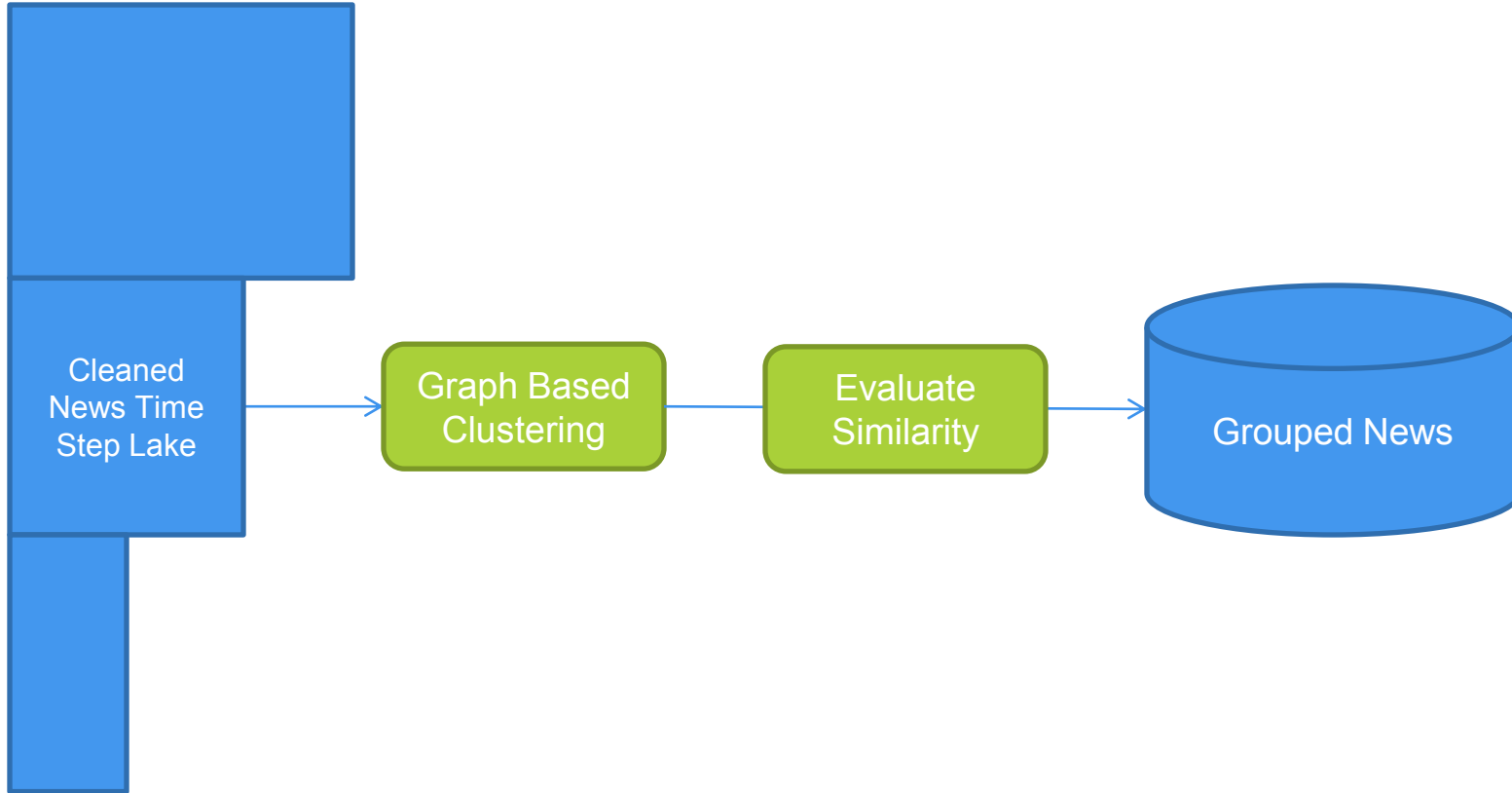
>> LSTM Text Classifier

>> SVM, Random Forest  
(With Meta Classification  
and Boosting Features)





## Model 04 – News Clustering



//

>> HDBSCAN

>> Cosine & Jaccard  
Combined Similarity  
Detection

>> Spectral Clustering





# Some Useful Python Libraries

- ▶ Tensorflow, Scikit-learn,
- ▶ Gensim (Word2Vec, Doc2Vec)
- ▶ NLTK, HAZM
- ▶ CoreNLP
- ▶ spaCy

.

# ParsPack

Ticketing, Log Analyter, Customer Care

# Matrix Again! But...

Convert Text Data Into Matrix  
and Concatenate with Other  
Features





نام	مسعود کاویانی	آدرس ایمیل	kaviani.masoud@gmail.com
موضوع		اولویت	متوسط
بخش	پشتیبانی فنی	سرویس مربوطه	هیچکدام

پیام

با سلام  
به نظر میرسه یکی از سرورها مون دچار کندی سرعت شده  
تعدادی از کاربرها میگن سایت خیلی دیر لود میشه  
ممنون از رسیدگیتون

## فناوری ابری مقیاس پذیر

از یک استارت آپ کوچک تا پشتیبان تان هستیم.



شروع کنید ↓



# Multi Label Classification/Regression

شدت خطا	مشکل	دپارتمان	تیکت
۹۰	ISP	سرور دانلود	کلا هیچی دانلود نمیشه، حتی نمی‌تونم به سرورم وصل بشم، پرداختی رو هم انجام دادم و فاکتورم پرداخت شده، ممنون میشم رسیدگی کنین
۵۷	Load	اشتراکی	از صبح سایت بالا میاد ولی حس می‌کنم به سرعت گذشته نیست
۴۱	Load, FTP	اشتراکی	الان با اینترنت ایرانسل دارم آپلود می‌کنم ولی حدودا یک ساعت طول کشیده
۸۴	DNS, Ping	CDN	من تصاویر اپلیکیشن رو بردم روی سی دی ان ولی دبیایگ که می‌کنم ۴۰۴ میده



# What is Text Mining?

تیکت	میانگین تیکت در ماه	تعداد سایت‌ها	تعداد سرورهای اختصاص ی	تعداد دامنه‌ها
کلا هیچی دانلود نمیشه، حتی نمی‌تونم به سرورم وصل بشم، پرداختی رو هم انجام دادم و فاکتورم پرداخت شده، ممنون میشم رسیدگی کنین	۳	۳	۱	۱۰
از صبح سایت بالا میاد ولی حس می‌کنم به سرعت گذشته نیست	۱	۱۳	۸	۱۲
الان با اینترنت ایرانسل دارم آپلود می‌کنم ولی حدودا یک ساعت طول کشیده	۱۲	۱۷	۰	۱
من تصاویر اپلیکیشن رو بردم روی سی دی ان ولی دیباگ که می‌کنم ۴۰۴ می‌ده	۱۴۰	۸	۲	۱





# Servers and Tickets



# Customer Care





# Prediction Customer Churn

ریزش	نرخ رشد سرویس ها در ماه	تعداد سرویس ها	شیب رشد پرداختی در ماه	میانگین مبلغ پرداختی در ماه	شیب رشد تیکت	تعداد تیکت در هفته	کلمات مهم تیکت ها
○	۰/۵	۱	۱/۱	۱۲۰۰۰۰	۰/۴۲	۴	سرور - هاست - سایت - خراب - کند - فکر - پول - فاکتور
۱	۱	۵	۱	۴۹۰۰۰	۲/۴۱	۱	خراب - کمک - همکاری - هلند - کنار
○	۰/۴	۳	۰/۹	۸۹۰۰۰	۰/۴۱	۱	هاست - عکس - تصویر - پینگ - ۴۰۴ - اررور
○	۱	۱	۱	۱۲۰۰۰	۱۲/۱	۰	سرور - کمک - ایران - ۴۰۴

# Log Analyzing

Logs  
Generate  
Very Fast and  
it is Hard to  
Mine these  
Logs at Real  
Time

```
root@jessie:~# grep addr /var/www/logs/tecmintlovesnginx.error.log --color=auto
2015/12/17 18:14:33 [error] 2234#0: *2 limiting connections by zone "addr", cli
nt: 192.168.0.25, server: tecmintlovesnginx.com, request: "GET /index.html HTTP
.0", host: "192.168.0.25"
2015/12/17 18:14:33 [error] 2234#0: *4 limiting connections by zone "addr", cli
nt: 192.168.0.25, server: tecmintlovesnginx.com, request: "GET /index.html HTTP
.0", host: "192.168.0.25"
2015/12/17 18:14:33 [error] 2234#0: *6 limiting connections by zone "addr", cli
nt: 192.168.0.25, server: tecmintlovesnginx.com, request: "GET /index.html HTTP
.0", host: "192.168.0.25"
2015/12/17 18:14:33 [error] 2234#0: *8 limiting connections by zone "addr", cli
nt: 192.168.0.25, server: tecmintlovesnginx.com, request: "GET /index.html HTTP
.0", host: "192.168.0.25"
root@jessie:~#
```



# Some Outlier Detection Algorithms

- ▶ Density-based OD
- ▶ One-Class SVM
- ▶ ISOLATION Forest
- ▶ Angle Based OD

# Thanks!

**Any questions?**

You can find me at:

- ▶ [Masoudkaviani.ir](http://Masoudkaviani.ir)
- ▶ [Chistio.ir](http://Chistio.ir)



## **Image Refs:**

<https://storybydata.com/datacated-challenge/text-mining-for-better-insights/>

<https://screencrush.com/the-matrix-reloaded-15-anniversary-defense/>

<https://www.imdb.com/title/tt0435998/>

<https://www.zeliosanalytics.com/blog/how-to-predict-customer-churn>

<https://www.tecmint.com/nginx-web-server-security-hardening-and-performance-tips/>

<https://scotthelme.co.uk/monitoring-http-2-usage-in-the-wild/>