

New Media Data Analytics and Application Lecture 10: Text Mining and Data Visualization

Ting Wang

Outlines

• Text Mining

– Data Visualization using Python

• Data Mining Essentials







online text data mining based on natural language processing Text Mining

Now, we have data, how to mining it?





Case Description Motivations:

- To measure a news objectively
- To obtain new information efficiently

Methodologies:

- Describe a news report by quantitative method
- Technical integration by computer science, statistics and journalism







- 1. Download a news report
- 2. Word segmentation
- 3. Word tag extraction and statistical computing
- 4. Data visualization and news summarization



Step 1: Download a News Report

• Example: http://news.sina.com.cn/w/2018-05-21/doc-ihawmatz9906261.shtml

news.sina.com.cn/w/2018-05-21/doc-ihawmatz9906261.shtml

五大国为这事要被逼得联手 这次又和特朗普有关

A A

原标题:德媒:五大国这次要被逼得联手了

据德国《星期日世界报》20日报道,来自德国、法国、英国、俄罗斯和中国的外交官员正在协商一项 新协议,希望借此挽救2015年签署的伊核协议,并说服特朗普解除对伊朗的制裁。这些外交官还将于25 日在维也纳就此举行会议。不过路透社20日援引3名欧盟消息人士的话否认与会各方将讨论新协议。分析 认为,鉴于欧盟自知力量有限,因此有意与中俄共同商讨新协议,但短期内,这个目标并不现实。

Home / Europe, China, Russia discussing new deal for Iran

Europe, China, Russia discussing new deal for Iran

《星期日世界报》从欧盟高层人士获得的消息称,德国、法国、英国、俄罗斯和中国间的会谈定在下周末,但美国不会出席,伊朗官员是否参加还不得而知。会谈的目的是商讨美国退出伊朗核协议后的下一步进程。

报道称,新协议和2015年的伊核协议相似,但新增限制伊朗弹道导弹和地区角色的条款,未来还有可能增加对伊朗的财政援助内容。如果新协议能够达成,有助于说服特朗普解除对伊朗的制裁。

但3名曾参与阻止美国总统特朗普退出伊核协议谈判的欧盟消息人士20日晚些时候告诉路透社,上述 消息并不正确,"本周五的维也纳会议将讨论伊核协议的实施问题和细节。"德国外交部目前尚未就有关 消息予以回应。

News also can be obtained by web crawler or databases





Step 2: Word segmentation (1) Database Preparation

- Word Dictionary (required)
- Stop Word Dictionary (required)
- Dictionaries of Terms (optional)
- Word Chains (required if using N-gram)
- Part of Speech (optional)
- Word Sentiment (optional for Sentiment Analysis)





Step 2: Word segmentation (2)

Chinese Word Segmentation

MaxLen=10 #最大词长

def word seg fmm(content): #正向匹配

- FMM
- BMM
- N-gram



MM Len=MaxLen #动态切割词长 Seg_Content="" #返回的切割结果

> while len(content)>0: if content[0:Len] in WordMap: #词典中有匹配 Seg_Content=Seg_Content+content[0:Len]+" | " content=content[Len:] Len=MaxLen #print("Seg_Content1:"+Seg_Content) continue else: #词典中无匹配 Len=Len-1 if Len==1:#仅剩一个词还没匹配到 Seg_Content = Seg_Content + content[0:Len] + " | " content = content[Len:] Len = MaxLen #print("Seg_Content2:" + Seg_Content) return Seg_Content[:-1]

def word_seg_bmm(content): #逆向匹配 MaxLen=10 #最大词长 Len=MaxLen #动态切割词长 Seg_Content="" #返回的切割结果

while len(content)>0: if content[-Len:] in WordMap: #词典中有匹配 Seg_Content=content[-Len:]+" "+Seg_Content content=content[:-Len] Len=MaxLen #print("Seg_Content1:"+Seg_Content) continue else: #词典中无匹配 Len=Len-1 if Len==1:#仅剩一个词还没匹配到 Seg_Content = content[-Len:] + " | " + Seg_Content content = content[:-Len] Len = MaxLen #print("Seg_Content2:" + Seg_Content) return Seg_Content[:-1]



Step 2: Word segmentation (3)

- Tips for Chinese Word Segmentation
 - Initialization is very important
 - Segment in the memory (not hard disk or data bases) to <u>accelerate</u> the segmentation speed
 - Using "set" to store the dictionary, and "dict" for segmented words in Python
 - For Tag Analysis, a precise word segmentation is <u>unnecessary</u>



Step 3: Word Tag Extraction and Statistical Computing

- str.split() for all tags
- Discarding One-Char tags
- Discarding Stop-Word tags
- Select tags whose term frequencies are larger than a threshold (for example >2)
- Other statistical computing



Step 4: Data Visualization and News Summarization



Data Visualization using Python

- Necessity:
 - NumPy (Computing Package)
 - Scipy (Scientific Computing Package)
 - Pillow(Image)
 - Matplotlib (Diagram Package)
 - wordcloud (Word Cloud Package)
 - Some packages also need some other required packages





Result

正向匹配 (FMM) 结果:

原 标题 : 德 媒 : 五大 国 这次 要 被 逼 得 联手 了 | | | | | ┃据|德国|《|星期日|世界|报|》|2|0|日报|道|, |来自|德国|、|法国|、|英国|、|俄罗斯|和|中国|的|外交 官|员|正在|协商|一|项|新|协议|, |希望|借|此|挽救|2|0|1|5|年|签署|的|伊|核|协议|, |并|说服|特|朗|普|解除|对|伊朗|的|制裁|。|这些|外交官|还|将|于|2|5|日|在|维也纳 |就此|举行|会议|。|不过|路透社|2|0|日|援引|3|名|欧盟|消息|人士|的话|否认|与会|各方|将|讨论|新|协议|。|分析|认为|,|鉴于|欧盟|自知|力量|有限|,|因此|有意|与|中|俄 | | 《 | 星期日 | 世界 | 报 | 》 | 从 | 欧盟 | 高层 | 人士 | 获得 | 的 | 消息 | 称 | , | 德国 | 、 | 法国 | 、 | 英国 | 、 | 俄罗 共同商讨新协议,但短期内,这个目标并不现实。 斯|和|中国|间|的|会谈|定|在下|周末|,|但|美国|不会|出席|,|伊朗|官员|是否|参加|还|不得而知|。|会谈|的|目的|是|商讨|美国|退出|伊朗|核|协议|后|的|下|一|步|进程|。|| 人士|2|0|日|晚|些|时候|告诉|路透社|,|上述|消息|并不|正确|,| "|本周|五|的|维也纳|会议|将|讨论|伊|核|协议|的|实施|问题|和|细节|。|"|德国|外交部|目前|尚未|就|有关 【虽然】维也纳】会议】的】具体】议题】尚】不明】确】,】但】 "】为】挽救】伊 | 核 | 协议】】,】五】国】正】组成】联合】阵线】 "】。】 "】德国】之】声】"】20】 消息 予以 回应 。 日|称|,|计划|中|的|会议|显示|欧盟|致力|于|确保|伊|核|协议|得以|继续|执行|,|即便|这|意味着|他们|要|在|脱离|美国|的|情况|下|,|与|莫斯科|、|北京|和|德黑兰|展开|合 |卡塔尔|半岛|电视台|2|0|日|称|,|自|5|月|8|日|特|朗|普|宣布|退出|伊|核|协议|以来|,|欧洲|和|德黑兰|相互|谨慎|接近|,|双方|声明|遵守|协议|的|要 作 【同时】监测】彼此】的【行为】,/【以】确保】履行【承诺】。【欧洲】国家【表示】将【尽力】保持【伊朗【石油】和】投资】的【流动】,/【但】同时【也】承认【这】并不【容易】。【伊朗【原子能】机构【负责 人 | 萨 | 利 | 希 | 表示 | , | 如果 | 欧洲 | 国家 | 未能 | 保留 | 协议 | , | 伊朗 | 有多 | 种 | 选择 | , | 包括 | 恢复 | 提炼 | 浓缩铀 | 至 | 纯度 | 2 | 0 | % | , | 并称 | 欧盟 | 只有 | 几 | 个 | 星期 | 的 | 时间 | 来 | 履行 | 其 | 承 | 而 | 《 | 星期日 | 世界 | 报 | 》 | 认为 | , | 之所以 | 要 | 寻找 | 新 | 途径 | , | 是因为 | 欧洲 | 官员 | 知道 | , | 欧洲 | 企业 | 在 | 美国 | 的 | 新 | 制裁 | 背景 | 下 | 难以 | 在 | 伊朗 | 进行 | 商 业|活动]。|欧盟|希望|伊朗|知道|,|只要|后者|遵守|伊|核|协议|,|欧盟|就|愿意|为|德黑兰|注资|。|欧盟|高级|官员|认为|,|布鲁塞尔|就|美国|的|制裁|措施|所|采取|的|对策 |, |对| "|伊朗|经济|的|积极|影响|非常|有限|" |, |因此|有|必要|与|中|俄|缔结|新|的|协议|。||||||| | |不过|, |中国|社会科学|院|西亚|非洲|所|副|研究员|王|凤|2|0| 日|对|《|环球|时报|》|记者|表示|,|各方|在|短期|内|就|伊|核|问题|达成|新|的|协议|并不|现实|。|因为|研发|弹道导弹|一直|是|伊|核计划|的|内容|,|很|难|要求|伊朗|停止 | 面对 | 美国 | 的 | 强势 | , | 欧盟 | 应该 | 怎么办 | ? | 美国 | 《 | 商业 | 内幕 |》|2|0|日|称|,|欧盟|可以|签署|一个|变动|极|小|的|协议|,|以|绥靖|特|朗|普|,|然后|坐等|他|任期|结束|。



Conclusions

本文与伊朗问 题有关,可能跟武 器和制裁有关, 起 决定力量的应该是 德国、 美国、 中国 和伊朗。 欧洲与此 消息关系较大。







machine learning approaches for data mining Data Mining Essentials



- Data Mining is the power for producing highquality journalism.
- Data Mining is an interdisciplinary subfield of computer science, and statistics.





Social Demands

•Data production rate has increased dramatically (**Big Data**) and we are able store much more data

E.g., purchase data, social media data, cell phone data

•Businesses and customers need <u>useful</u> or <u>actionable</u> knowledge to gain insight from raw data for various purposes

- It's not just searching data or databases



The process of extracting useful patterns from raw data is known as Knowledge Discovery in Databases (KDD)



KDD from Data Bases





Data 数据

Continuous Data 连续型数据

 Regression

Discrete Data 离散型数据
 – Classification





Data Feature (1) 数字特征

Feature also called as Measurement, Attribute

- Nominal 名词性
 - Operations:
 - Mode (most common feature value), Equality Comparison
 - E.g., {male, female}
- Ordinal 序数性
 - Feature values have an intrinsic order to them, but the difference is not defined
 - Operations:
 - same as nominal, feature value rank
 - E.g., {Low, medium, high}



Data Feature (2) 数字特征

- Interval 间隔性
 - Operations:
 - Addition and subtractions are allowed whereas divisions and multiplications are not
 - E.g., 3:08 PM, calendar dates
- Ratio 比例性
 - Operations:
 - divisions and multiplications are allowed
 - E.g., Height, weight, money quantities



Data Quality 数据质量

- Noise 噪声数据
 - Noise is the distortion of the data
- Outliers 异常值
 - Outliers are data points that are considerably different from other data points in the dataset
- Missing Values 缺失值
 - Missing feature values in data instances
 - Solution:
 - Remove instances that have missing values
 - Estimate missing values, and
 - Ignore missing values when running data mining algorithm
- Duplicate data 重复数据



- Data Preprocessing (1)
 数据预处理
- Aggregation 聚合
 - It is performed when multiple features need to be combined into a single one or when the scale of the features change
 - Example: image width , image height -> image area (width x height)
- Discretization 离散化
 - From continues values to discrete values
 - Example: money spent -> {low, normal, high}



- Data Preprocessing (2) 数据预处理
- Feature Selection 特征选择
 - Choose relevant features
- Feature Extraction 特征提取
 - Creating new features from original features
 - Often, more complicated than aggregation
- Sampling 取样
 - Random Sampling
 - Sampling with or without replacement
 - Stratified Sampling: useful when having class imbalance
 - Social Network Sampling



Machine Learning 机器学习

- Supervised Learning
 - Classification
 - Regression
- Unsupervised Learning
 - Clustering
 - Dimensional Reduction





Supervised Machine Learning 有监督学习







Prediction Result with Labeled Discrete Value

- KNN(K-Nearest Neighbors) KI临近原则
- Linear Classifier 线性分类器
- Neural Networks 神经网络
- Support Vector Machine 支撑向量机
- Decision Tree 决策树







は まままま した ようトロネナッジ SHANCHALINTERNATIONAL STUDIES UNIVERSITY

Regression (1) 回归

Prediction Result with Unlabeled Continuous Value





Regression (2)回归 Nonlinear Regression 非线性回归计算 • Linearization 线性化方法

1. Transformation 变形法

$$y=ae^{bx}U$$
 $ightharpoonup \ln\left(y
ight)=\ln\left(a
ight)+bx+u$

2. Segmentation 分割法

split up into classes or segments and *linear* regression can be performed per segment





Unsupervised Machine Learning



machine learning task of inferring a function to describe hidden structure from <u>unlabeled data</u>





Clustering <u></u> *聚类*

- **Clustering Goal:** Group together similar items
- Clustering algorithms group together similar items

 The algorithm does not have examples showing how the samples should be grouped together (unlabeled data)
 Y integration

Similarity Computing (1) 相似度计算

-The most popular (dis)similarity measure for continuous features are **Euclidean Distance** and **Pearson Linear Correlation**



 $d(X,Y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$



Similarity Computing (2) 相似度计算	$X = (x_1, x_2, \dots, x_n)$
X and Y are n Dimensional Vectors	$Y = (y_1, y_2, \dots, y_n)$

Measure Name	Formula	Description
Mahalanobis	$d(X,Y) = \sqrt{(X-Y)^T \Sigma^{-1} (X-Y)}$	X, Y are features vec-
		tors and Σ is the co-
		variance matrix of the
		dataset
Manhattan $(L_1 \text{ norm})$	$d(X,Y) = \sum_{i} x_i - y_i $	X, Y are features vec-
		tors
L_p -norm	$d(X,Y) = \left(\sum_{i} x_{i} - y_{i} ^{n}\right)^{\frac{1}{n}}$	X, Y are features vec-
		tors

Once a distance measure is selected, instances are grouped using it.



Pearson Linear Correlation 皮尔逊线性相关

Correlation Coefficient 相关系数 "x = El

$$ho_{X,Y} = rac{\mathrm{cov}(X,Y)}{\sigma_X\sigma_Y}$$

 $egin{aligned} &\sigma_Y^2 = \mathrm{E}[(Y-\mathrm{E}[Y])^2] = \ &\mathrm{E}[(X-\mu_X)(Y-\mu_Y)] \end{aligned}$

$$\mu_X = E[X]$$

$$\mu_Y = E[Y]$$

$$\sigma_X^2 = E[(X - E[X])^2] = E[X^2] - [E[X]]^2$$

$$\sigma_X^2 = E[(Y - E[Y])^2] = E[Y^2] - [E[Y]]^2$$

Deletione hetureen

$$\mathrm{E}[(X-\mu_X)(Y-\mu_Y)] = \mathrm{E}[(X-\mathrm{E}[X])(Y-\mathrm{E}[Y])] = \mathrm{E}[XY] - \mathrm{E}[X]\mathrm{E}[Y]$$

Where, cov is the covariance $\sigma_{\rm c}$ is the standard deviation

$$\mathrm{cov}(X,Y) = \mathrm{E}[(X-\mu_X)(Y-\mu_Y)]$$

$$\rho_{X,Y} = \frac{\mathbf{E}[XY] - \mathbf{E}[X] \mathbf{E}[Y]}{\sqrt{\mathbf{E}[X^2] - [E[X]]^2} \sqrt{\mathbf{E}[Y^2] - [E[Y]]^2}}$$



Film Ranking Correlation

Superman was rated 3 by Mick LaSalle and 5 by Gene Seymour, so it is placed at (3,5) on the chart.



上海小型法大 Stagen Number International Studies UNIV Conclusion: Films recommended to Lisa, also can be recommended to Jack.



- PCA is a statistical procedure *converts* a set of observations of possibly <u>correlated variables</u> *into* a set of values of linearly <u>uncorrelated variables</u> called principal components.
- 2. The number of principal components is <u>less than or equal to</u> the number of original variables.
- This transformation is defined in such a way that <u>the first principal component</u> has *the largest possible variance*, and each <u>succeeding component</u> in turn has <u>the highest variance possible under the constraint</u> that it is orthogonal to the preceding components.



Ref. http://www.cnblogs.com/SCUJIN/p/5965946.html







Books and Chapters (1) https://item.jd.com/11983227.html Chapter 1-2 Machine Learning Package Installation Machine Learning Theory Foundations





Books and Chapters (2) https://item.jd.com/11803260.html Chapter 5 Data Mining Essentials

Online Reference: <u>http://www.public.asu.edu/~huanliu/</u>





Books and Chapters (3) https://item.jd.com/11676691.html Python Data Visualization





Books and Chapters (4) https://item.jd.com/11667512.html Programming Collective Intelligence





Python Extension Packages

http://www.lfd.uci.edu/~gohlke/pythonlibs/

 \leftarrow \rightarrow X (i) www.lfd.uci.edu/~gohlke/pythonlibs/

Unofficial Windows Binaries for Python Extension Packages

by Christoph Gohlke, Laboratory for Fluorescence Dynamics, University of California, Irvine.

This page provides 32- and 64-bit Windows binaries of many scientific open-source extension packages for the official <u>CPvthon distribution</u> of the <u>Pvthon</u> programming language.

The files are unofficial (meaning: informal, unrecognized, personal, unsupported, no warranty, no liability, provided "as is") and made available for testing and evaluation purposes.

If downloads fail reload this page, enable JavaScript, disable download managers, disable proxies, clear cache, and use Firefox. Please only download files manually as needed.

Most binaries are built from source code found on PVPI or in the projects public revision control systems. Source code changes, if any, have been submitted to the project maintainers or are included in the packages.

Refer to the documentation of the individual packages for license restrictions and dependencies.

Use pip version 8 or newer to install the downloaded .whl files. This page is not a pip package index.

Many binaries depend on numpy-1.11+nkl and the Microsoft Visual C++ 2008 (x64, x86, and SP1 for CPython 2.6 and 2.7), Visual C++ 2010 (x64, x86, for CPython 3.3 and 3.4), or the Visual C++ 2015 (x64 and x86 for CPython 3.5 and 3.6) redistributable packages.

Install <u>numpv+mkl</u> before other packages that depend on it.

The binaries are compatible with the official CPython distribution on Windows >=6.0. Chances are they do not work with custom Python distributions included with Blender, Maya, ArcGIS, OSGeo4W, ABAQUS, Cygwin, Pythonxy, Canopy, EPD, Anaconda, WinPython etc. Many binaries are not compatible with Windows XP or Wine.

The packages are ZIP or 7z files, which allows for manual or scripted installation or repackaging of the content.

The files are provided "as is" without warranty or support of any kind. The entire risk as to the quality and performance is with you.

Index by date: greenlet pygresql netcdf4 lxml pyamg jupyter cython liblinear cobra pybox2d fastcluster vlfd sfepy pytables h5py grako fonttools pymol pygame pyflux matplotlib spacy cytoolz apsw chainer mathutils veusz mercurial pyeda numpy cvxopt pywavelets pymongo gr persistent aichttp pycobc twisted ets vtk pocketsphinx simpleaudio pyaudio soundevice fisx tensorflow multiprocess libsbml cvxcanon spectrum pyvrm197 ta-lib pythonmagick pyzmq triangle pgmagick ujson yappi pyfltk mod_wsgi pyfftw py_gd pyviennacl python-ldap openpiv pyx mpi4py pyephem pyemd planar mysqlclient xxhash zarr regex ode spyder lsqfit fann2 fisher ffnet entropy autopy slycot sparsesvd scs ecos sasl twainmodule dulwich datrie cx_oracle cyordereddict coverage cdecimal cartopy blz bigfloat aspell-python simpleparse milk menpo marisa-trie llist setproctitle hddm hmmlearn seqlearn jsonlib rtree rtmidi-python udunits heatmap scikit-umfpack scikits.vectorplot kwant tinyarray rpy2 fina cx_freeze operov netifaces multineat basemap py-earth pulp mlpy reportlab pyminuit pymetis python-snappy python-lzo python-levenshtein python-lz4 pystemmer



Data Visualization in Python

- <u>http://it.sohu.com/20151119/n427117609.shtml</u>
- <u>http://www.oschina.net/translate/python-data-visualization-libraries</u>





Using WordCloud

- http://blog.csdn.net/tanzuozhev/article/details/50789226
- <u>https://www.oschina.net/code/snippet_2294527_56155</u>

Chinese Display

- http://blog.csdn.net/u012705410/article/details/47379957



Provided Repositories for Social Mining

- http://socialcomputing.asu.edu
- <u>http://snap.Stanford.edu</u>
- <u>https://github.com/caesar0301/awesome-public-datasets</u>







The End of Lecture 10

Thank You



http://www.wangting.ac.cn