

20211006 Singular Value Decomposition

Wednesday, October 6, 2021 4:12 PM

- HW1
 - o A lot of back and forth about what packages you are allowed to use
 - Prof thought it was pretty clear
 - Some miscommunication has happened
 - Now some people used pandas for question 2 and now that is probably allowed
- Moving forward:
 - o HWs are pretty much open to whatever packages you want
 - o Otherwise, may put some piazza post specifying it
- You can implement all this stuff on your own
- Team lead role: check out piazza
- Posted monthly check-in poll
- HW grades are posted
- Midterm will have even more data, and will take even more time to code
 - o Need to get used to running code overnight
 - o Computation time needs to be factored into your planning
 - o Kaggle competition, will take about 10 days
 - Start as soon as its posted

Singular Value Decomposition

Boston University CS 506 - Lance Galletti

- Another set of unsupervised techniques

Recall

$$\text{n data points} \left\{ \begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1m} \\ \vdots & \ddots & \vdots & & \vdots \\ x_{i1} & \dots & x_{ij} & \dots & x_{im} \\ \vdots & & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{nm} \end{pmatrix} \right. \underbrace{\phantom{\begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1m} \\ \vdots & \ddots & \vdots & & \vdots \\ x_{i1} & \dots & x_{ij} & \dots & x_{im} \\ \vdots & & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{nm} \end{pmatrix}}}_{\text{m features}}$$

Goal

Examine this matrix and uncover its linear algebraic properties to:

1. Approximate A with a smaller matrix B that is easier to store but contains similar information as A
 2. Dimensionality Reduction / Feature Extraction
 3. Anomaly Detection & Denoising
-

- Supposing everything is numeric in our dataset
- Goals are 3:
 1. Smaller matrix, so it takes less storage and all algorithms and improve the performance
 - i. Maybe even denoise the dataset a little here
 2. Dimensionality reduction AND feature extraction
 3. anomaly detection and denoising

Linear Algebra Review

Definition: The vectors in a set $V = \{v_1, \dots, v_n\}$ are **linearly independent** if

$$a_1v_1 + \dots + a_nv_n = \mathbf{0}$$

can only be satisfied by $a_1 = 0$

Note: this means no vector in that set can be expressed as a **linear combination** of other vectors in the set.

- Turns out probability wasn't actually a prereq, but is linear algebra?
- Linearly independent:
 - o Each element can't be represented by a linear combination of the other vectors in the set
 - o The only special case is if all a is 0, but only allowed case

Linear Algebra Review

Definition:

The **determinant** of a square matrix A is a scalar value that encodes properties about the **linear mapping** described by A .

2x2:

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad \det(A) = ad - bc$$

- Determinant
 - o Expresses so much about a matrix

Linear Algebra Review

Definition:

The **determinant** of a square matrix A is a scalar value that encodes

Definition:

The **determinant** of a square matrix A is a scalar value that encodes properties about the **linear mapping** described by A.

3x3:

$$A = \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} \quad \det(A) = a \cdot \det \begin{pmatrix} e & f \\ h & i \end{pmatrix} - b \cdot \det \begin{pmatrix} d & f \\ g & i \end{pmatrix} + c \cdot \det \begin{pmatrix} d & e \\ g & h \end{pmatrix}$$

Row 1 col 1
Row 1 col 2
Row 1 col 3

- topic of this lab
- Do you see how to extend this to the nxn case?

Linear Algebra Review

Definition:

The **determinant** of a square matrix A is a scalar value that encodes properties about the **linear mapping** described by A.

n X n:

Can recursively compute it. How?

Linear Algebra Review

Property:

n vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ in an **n**-dimensional space are **linearly independent** iff the matrix **A**:

$$\mathbf{A} = [\mathbf{v}_1, \dots, \mathbf{v}_n] \quad (n \times n)$$

has non-zero determinant.

Q: Can **m > n** vectors in an **n**-dimensional space be linearly independent?

- No

Linear Algebra Review

Definition:

A **basis** B of a vector space (over a field F) is a **linearly independent** subset of V that **spans** V . B **spans** V if for every vector v in V it is possible to choose v_1, \dots, v_n in F and b_1, \dots, b_n in B such that:

$$v = v_1 b_1 + \dots + v_n b_n$$

Ex: North & East in 2d-plane

Linear Algebra Review

Definition:

The **rank** of a matrix A is the dimension of the vector space spanned by its column space. This is equivalent to the maximal number of linearly independent columns / rows of A .

Definition:

A matrix A is **full-rank** iff $\text{rank}(A) = \min(m, n)$

Note: Get the rank of a matrix through the **Gram-Schmidt process**

- We normally talk about columns not rows
- Data is normally far from full rank?
 - o No, generally pretty close to full rank
 - Depends on the size of the matrix
- All columns could be linearly independent and it would still be full rank
- Gram-Smidt process
 - o This is how you get the rank of the matrix
 - o Take one out at a time, eventually you run out of vectors

Approximation

In practice, matrices describing our dataset contain a lot of redundant information.

It would be great to capture all the information of our dataset in the least amount of space possible.

Approximation

To store an $n \times m$ matrix A requires storing $m \cdot n$ values.

However, if the rank of the matrix of A is k , A can be factored as

$$A = UV$$

where

U is $n \times k$
 V is $k \times m$

which requires storing $k(m + n)$ values.

- You can factor a matrix of rank k into two matrices
 - o If k is small, how many values do you need to store?
 - o $n*k + k*m$ might be << $n*m$

Approximation

Goal:

Approximate \mathbf{A} with $\mathbf{A}^{(k)}$ (low-rank matrix) such that

1. $d(\mathbf{A}, \mathbf{A}^{(k)})$ is small
 2. k is small compared to m & n
-

Frobenius Distance

$$d_F(A, B) = \|A - B\|_F = \sqrt{\sum_{i,j} (a_{ij} - b_{ij})^2}$$

i.e. the pairwise differences in values of A and B

Approximation

Definition:

When $k < \text{rank}(A)$, the **rank-k approximation** of A (in the least squares sense) is

$$A^{(k)} = \arg \min_{\{B | \text{rank}(B)=k\}} d_F(A, B)$$

- Has the same dimensions, but with lower rank so we can factor things out and store as smaller matrices
- Infinite number of matrices with rank k , and all sorts of possible numbers in that matrix

Approximation

Definition:

The **Singular Value Decomposition** of a rank- r matrix A has the form

$$A = U \Sigma V^T$$

where

U is $n \times r$

The columns of **U** are orthogonal & unit length ($U^T U = I$)

V is $m \times r$

The columns of **V** are orthogonal & unit length ($V^T V = I$)

- **V** transpose is $r \times m$

Approximation

Definition:

The **Singular Value Decomposition** of a rank- r matrix A has the form

$$A = U \Sigma V^T$$

where

Approximation

Definition:

The **Singular Value Decomposition** of a rank-r matrix A has the form

$$A = U\Sigma V^T$$

where

$$\Sigma = \begin{pmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_r \end{pmatrix}$$

with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$

σ_i is the square root of the eigenvalues of $A^T A$ and are called **singular values**

$$A = U \Sigma V^T$$
$$d_F(A, A^{(k)})^2 = \sum_{i=k+1}^r \sigma_i^2$$

↑ k close to A, but little advantage?

- Diagonal matrix
- Singular values also have a particular size order

Approximation

Property:

$$d_F(A, A^{(k)})^2 = \sum_{i=k+1}^r \sigma_i^2$$

Note: the larger k is, the smaller the distance.

Approximation

Find $\mathbf{A}^{(k)}$ by decomposing \mathbf{A} :

$$\mathbf{A} = \begin{pmatrix} U_1 & U_2 \end{pmatrix} \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1 & V_2 \end{pmatrix}$$

$$\mathbf{A}^{(k)} = \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^\top$$

Where

\mathbf{U}_1 is $n \times k$

$\boldsymbol{\Sigma}_1$ is $k \times k$

\mathbf{V}_1 is $m \times k$

Approximation

$$\begin{array}{|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 0 \\ \hline 2 & 2 & 2 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 \\ \hline 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ \hline 0 & 0 & 0 & 3 & 3 \\ \hline 0 & 0 & 0 & 1 & 1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 0.18 & 0 \\ \hline 0.36 & 0 \\ \hline 0.18 & 0 \\ \hline 0.90 & 0 \\ \hline 0 & 0.53 \\ \hline 0 & 0.80 \\ \hline 0 & 0.27 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 9.64 & 0 \\ \hline 0 & 5.29 \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline 0.58 & 0.58 & 0.58 & 0 & 0 \\ \hline 0 & 0 & 0 & 0.71 & 0.71 \\ \hline \end{array}$$

\mathbf{U} $\boldsymbol{\Sigma}$ \mathbf{V}^\top

Rank 2

- Rank: of the dataset 2

Approximation

$$\begin{array}{|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 2 & 2 & 2 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 5 & 5 & 5 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 & 0 \\ \hline 0 & 0 & 0 & 3 & 3 & 0 \\ \hline 0 & 0 & 0 & 1 & 1 & 0 \\ \hline \end{array} \sim \begin{array}{|c|c|} \hline 0.18 & 0 \\ \hline 0.36 & 0 \\ \hline 0.18 & 0 \\ \hline 0.90 & 0 \\ \hline 0 & 0.53 \\ \hline 0 & 0.80 \\ \hline 0 & 0.27 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 9.64 & 0 \\ \hline 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline 0.58 & 0.58 & 0.58 & 0 & 0 \\ \hline 0 & 0 & 0 & 0.71 & 0.71 \\ \hline \end{array}$$

Approximation

$$\begin{array}{|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 2 & 2 & 2 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 5 & 5 & 5 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 & 0 \\ \hline 0 & 0 & 0 & 3 & 3 & 0 \\ \hline 0 & 0 & 0 & 1 & 1 & 0 \\ \hline \end{array} \sim \begin{array}{|c|c|} \hline 0.18 & 0 \\ \hline 0.36 & 0 \\ \hline 0.18 & 0 \\ \hline 0.90 & 0 \\ \hline 0 & 0.53 \\ \hline 0 & 0.80 \\ \hline 0 & 0.27 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 9.64 & 0 \\ \hline 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline 0.58 & 0.58 & 0.58 & 0 & 0 \\ \hline 0 & 0 & 0 & 0.71 & 0.71 \\ \hline \end{array}$$

- Rank two means we can represent the dataset with two features
- What these features mean and represent is not clear

Approximation

1	1	1	0	0
2	2	2	0	0
1	1	1	0	0
5	5	5	0	0
0	0	0	2	2
0	0	0	3	3
0	0	0	1	1

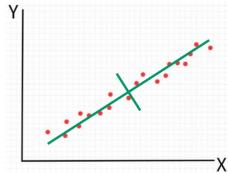
~

1	1	1	0	0
2	2	2	0	0
1	1	1	0	0
5	5	5	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

Demo

Approximation

The i^{th} singular vector represents the direction of the i^{th} most variance.



Singular Values express the importance / significance of a singular vector

- the i^{th} singular vector (the i^{th} column)
 - o Is in the direction of the i^{th} most variance
- $x = y$ is the first singular vector, with the most variance
- second most variance, is in the perpendicular
- So we basically want to change the coordinate system used
- Matrix U is orthonormal, all columns are perpendicular to each other

Approximation

To find the right k you can:

1. Look at the singular value plot to find the elbow point
 2. Look at the residual error of choosing different k
-

Demo

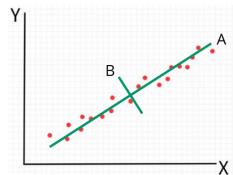
Principal Component Analysis

Idea: project the data onto a subspace generated from a subset of singular vectors / principal components.

We want to project onto the components that capture most of the variance / information in the data.

- The reasoning: we would like to project our data in a way to capture the most variance

Principal Component Analysis



Which principal component should we project on?

- Project data onto A

Demo

Are there datasets this will perform particularly poorly on?

(Demos all at end of notes)

Latent Semantic Analysis

Inputs are documents. Each word is a feature. We can represent each document by:

- The presence of the word (0 / 1)
- Count of the word (0, 1, ...)
- Frequency of the word ($n_i / \sum n_i$)
- TfIDf

$$\text{tf} \cdot \text{idf}$$

$\log \left(\frac{\text{number of documents}}{\text{number of documents that contain the term}} \right)$

- Presence of word 0 or 1
 - o Does not capture how much the word is used
 - o "This document is not science" will use the word science once, even though the document is not about science
- Even better than count an frequency: log inverse document frequency * term frequency
 - o Tells you frequently used in the document, factoring in how rarely used elsewhere
 - "the" is not representative of a unique topic if everyone is using it

Latent Semantic Analysis

	data	information	retrival	brain	lung
CS-paper-1	1	1	1	0	0
CS-paper-2	2	2	2	0	0
CS-paper-3	1	1	1	0	0
CS-paper-4	5	5	5	0	0
Med-paper-1	0	0	0	2	2
Med-paper-2	0	0	0	3	3
Med-paper-3	0	0	0	1	1

- CS papers vs medical papers
- Two linearly independent columns that represent our dataset
- Imagine you don't have access to which are CS and which are med papers:
 - o Want to identify how many topics are in this dataset

Latent Semantic Analysis

$$\begin{array}{|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 0 \\ \hline 2 & 2 & 2 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 \\ \hline 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ \hline 0 & 0 & 0 & 3 & 3 \\ \hline 0 & 0 & 0 & 1 & 1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 0.18 & 0 \\ \hline 0.36 & 0 \\ \hline 0.18 & 0 \\ \hline 0.90 & 0 \\ \hline 0 & 0.53 \\ \hline 0 & 0.80 \\ \hline 0 & 0.27 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 9.64 & 0 \\ \hline 0 & 5.29 \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline 0.58 & 0.58 & 0.58 & 0 & 0 \\ \hline 0 & 0 & 0 & 0.71 & 0.71 \\ \hline \end{array}$$

- Now we can associate meanings to each of these matrices

Latent Semantic Analysis

CS concept MD concept


$$\begin{array}{|c|c|} \hline 0.18 & 0 \\ \hline 0.36 & 0 \\ \hline 0.18 & 0 \\ \hline 0.90 & 0 \\ \hline 0 & 0.53 \\ \hline 0 & 0.80 \\ \hline 0 & 0.27 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 9.64 & 0 \\ \hline 0 & 5.29 \\ \hline \end{array} \times \begin{array}{|c|c|c|c|c|} \hline 0.58 & 0.58 & 0.58 & 0 & 0 \\ \hline 0 & 0 & 0 & 0.71 & 0.71 \\ \hline \end{array}$$

- 1st column: learned cs concept
- 2nd: MD concept

Latent Semantic Analysis

CS concept	MD concept	doc-to-concept similarity
0.18	0	
0.36	0	
0.18	0	
0.90	0	
0	0.53	
0	0.80	
0	0.27	

X **X**

9.64	0	0.58	0.58	0.58	0	0
0	5.29	0	0	0	0.71	0.71

- How similar is a given document to a given concept?
 - o One interpretation

Latent Semantic Analysis

doc-to-concept
similarity matrix

0.18	0	
0.36	0	
0.18	0	
0.90	0	
0	0.53	
0	0.80	
0	0.27	

X **X**

9.64	0	0.58	0.58	0.58	0	0
0	5.29	0	0	0	0.71	0.71

Latent Semantic Analysis

doc-to-concept
similarity matrix

0.18	0
0.36	0
0.18	0
0.90	0
0	0.53
0	0.80
0	0.27

$$\begin{matrix} & \text{X} \\ \text{X} & \end{matrix} \times \begin{matrix} 9.64 & 0 \\ 0 & 5.29 \end{matrix} \times \begin{matrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{matrix}$$

"strength" of the CS concept

- Strength of the concept

Latent Semantic Analysis

doc-to-concept
similarity matrix

0.18	0
0.36	0
0.18	0
0.90	0
0	0.53
0	0.80
0	0.27

$$\begin{matrix} & \text{X} \\ \text{X} & \end{matrix} \times \begin{matrix} 9.64 & 0 \\ 0 & 5.29 \end{matrix} \times \begin{matrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{matrix}$$

"strength" of the
each concept

Latent Semantic Analysis

doc-to-concept
similarity matrix

0.18	0
0.36	0
0.18	0
0.90	0
0	0.53
0	0.80
0	0.27

"strength" of the
each concept

X

X

9.64	0
0	5.29

X

0.58	0.58	0.58	0	0
0	0	0	0.71	0.71

term-to-concept similarity

- How similar is a particular word to the concept
- How useful is a particular word to the concept

Latent Semantic Analysis

doc-to-concept
similarity matrix

0.18	0
0.36	0
0.18	0
0.90	0
0	0.53
0	0.80
0	0.27

"strength" of the
each concept

X

X

term-to-concept similarity
matrix

9.64	0
0	5.29

Demo

Anomaly Detection

Define $\mathbf{O} = \mathbf{A} - \mathbf{A}^{(k)}$

The largest rows of \mathbf{O} could be considered anomalies

- To do anomaly detection, you need to ?

Demo

**lسا.py is one of
most important
files for the
midterm!!!!**

```
1 import numpy as np
2 from sklearn import metrics
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 from sklearn.datasets import fetch_20newsgroups
6 from sklearn.feature_extraction.text import
    TfidfVectorizer
7 from nltk.stem.snowball import SnowballStemmer
8 from nltk.tokenize import word_tokenize, sent_tokenize
9
10 # THIS IS PROBABLY ONE OF THE MOST IMPORTANT FILES FOR
#     THE MIDTERM
11 # First preprocess your data. You should always do
#     this
12 # Reducing prefixes and suffixes to try to condense
#     it farther
13 # Limitation: how you are handling stripping on prefix
#     /suffix and ignoring context
14 # Transform, center data, plotted singular values
15 # Already know there are 3 concepts, so want to do
#     kmeans with 3 clusters to see how well you can cluster
#     based on
16 # this
17
18 categories = ['comp.os.ms-windows.misc', 'sci.space',
    'rec.sport.baseball']
19 news_data = fetch_20newsgroups(subset='train',
    categories=categories)
20 vectorizer = TfidfVectorizer(stop_words='english',
    min_df=4,max_df=0.8)
21
22 stemmed_data = [ " ".join(SnowballStemmer("english",
    ignore_stopwords=True).stem(word)
    for sent in sent_tokenize(message))
    for word in word_tokenize(sent))
    for message in news_data.data]
23
24
25
26
27 dtm = vectorizer.fit_transform(stemmed_data)
28 terms = vectorizer.get_feature_names()
29 centered_dtm = dtm - np.mean(dtm, axis=0)
30
```

```

File - C:\Users\Wolfs\PycharmProjects\cs506lec20211006\lsa.py
31 u, s, vt = np.linalg.svd(centered_dtm)
32 plt.xlim([0,50])
33 plt.plot(range(1,len(s)+1),s)
34 plt.show()
35
36 ag = []
37 max = len(u)
38 for k in range(1,25):
39     vectorsk = u[:, :k] @ np.diag(s[:k])
40     kmeans = KMeans(n_clusters=3, init='k-means++',
41                      max_iter=100, n_init=10, random_state=0)
41     kmeans.fit_predict(vectorsk)
42     labelsk = kmeans.labels_
43     ag.append(metrics.v_measure_score(labelsk,
44                                         news_data.target)) # closer to 1 means better
44                                         clustering
44
45 plt.plot(range(1,25),ag)
46 plt.ylabel('Agreement',size=20)
47 plt.xlabel('No of Prin Comps',size=20)
48 plt.show()
49
50

```

Page 2 of 2

```

Traceback (most recent call last):
  File "C:/Users/Wolfs/PycharmProjects/cs506lec20211006/lsa.py", line 20, in <module>
    stemmed_data = [" ".join(SnowballStemmer("english", ignore_stopwords=True).stem(word)
  File "C:/Users/Wolfs/PycharmProjects/cs506lec20211006/lsa.py", line 21, in <listcomp>
    for sent in sent_tokenize(message)
  File "C:/Users/Wolfs/anaconda3/envs/cs506lec20211006/lib/site-packages/nltk/tokenize._init_.py", line 106, in sent_tokenize
    tokenizer = load(f"tokenizers/punkt/{language}.pickle")
  File "C:/Users/Wolfs/anaconda3/envs/cs506lec20211006/lib/site-packages/nltk\data.py", line 750, in load
    opened_resource = _open(resource_url)
  File "C:/Users/Wolfs/anaconda3/envs/cs506lec20211006/lib/site-packages/nltk\data.py", line 870, in _open
    return find(path_, path + ["]").open()
  File "C:/Users/Wolfs/anaconda3/envs/cs506lec20211006/lib/site-packages/nltk\data.py", line 583, in find
    raise LookupError(resource_not_found)
LookupError:
*****
Resource punkt not found.
Please use the NLTK Downloader to obtain the resource:

>>> import nltk
>>> nltk.download('punkt')

For more information see: https://www.nltk.org/data.html

Attempted to load tokenizers/punkt/english.pickle

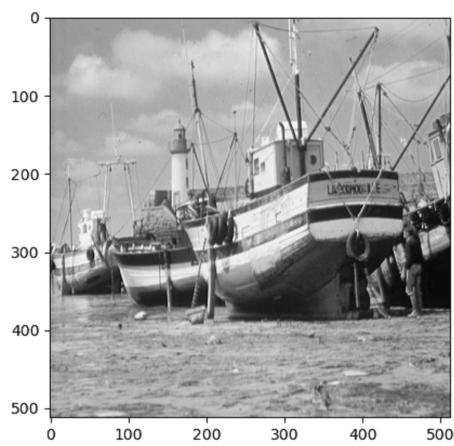
Searched in:
- C:\\Users\\Wolfs\\nltk_data'
- 'C:\\\\Users\\\\Wolfs\\\\anaconda3\\\\envs\\\\cs506lec20211006\\\\nltk_data'
- 'C:\\\\Users\\\\Wolfs\\\\anaconda3\\\\envs\\\\cs506lec20211006\\\\share\\\\nltk_data'
- 'C:\\\\Users\\\\Wolfs\\\\anaconda3\\\\envs\\\\cs506lec20211006\\\\lib\\\\nltk_data'
- 'C:\\\\Users\\\\Wolfs\\\\AppData\\\\Roaming\\\\nltk_data'
- 'C:\\\\nltk_data'
- 'D:\\\\nltk_data'
- 'E:\\\\nltk_data'
- ...
*****

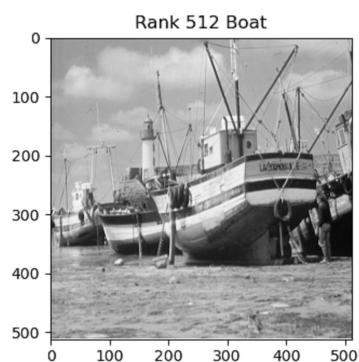
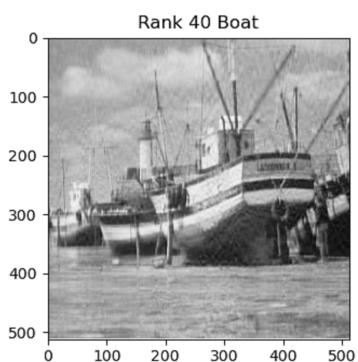
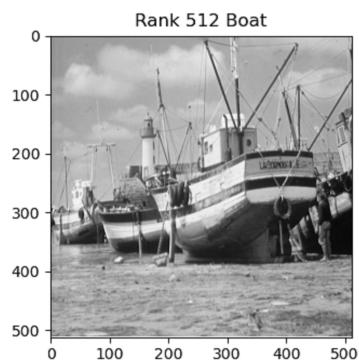
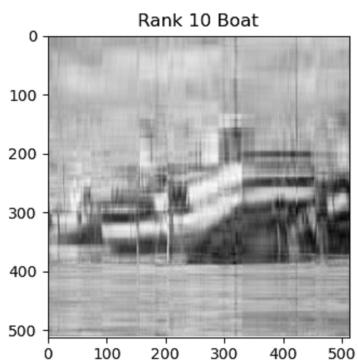
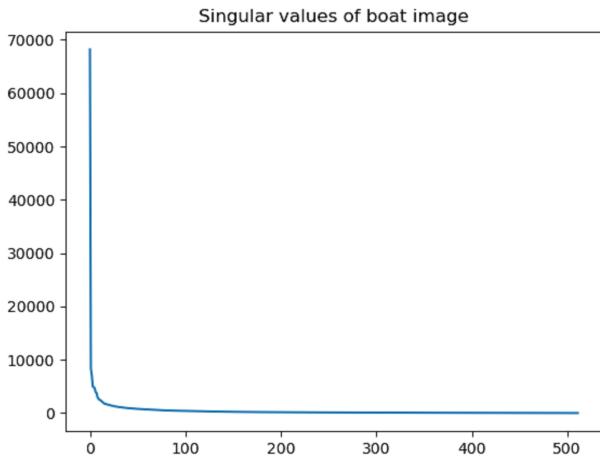
```

```
1 import numpy as np
2 import matplotlib.cm as cm
3 import matplotlib.pyplot as plt
4
5 # Compression
6 # Image on screen is rarely what you see on screen
7
8 # How big the singular values are represents how
# important that feature is to the dataset
9 # you can see we have very high singular values before
# it shoots down, so you can cut off at particular
# point
10 # So he decided to cut off at rank 10 (but rank 40
# seems more appropriate)
11 # With only 10 vectors, you can see boat still
12
13 boat = np.loadtxt('data/boat.dat')
14 plt.figure()
15 _ = plt.imshow(boat,cmap = cm.Greys_r)
16 plt.show()
17
18 u,s,vt=np.linalg.svd(boat,full_matrices=False)
19 _ = plt.plot(s)
20 plt.title('Singular values of boat image')
21 plt.show()
22
23 RANK = 10
24
25 for RANK in [10, 40]:
# construct a rank-40 version of the boat
26     scopy = s.copy()
27     scopy[RANK:]=0
28     boatApprox = u.dot(np.diag(scopy)).dot(vt)
29     #
30     plt.figure(figsize=(9,6))
31     plt.subplot(1,2,1)
32     plt.imshow(boatApprox,cmap = cm.Greys_r)
33     plt.title('Rank '+str(RANK)+' Boat')
34     plt.subplot(1,2,2)
35     plt.imshow(boat,cmap = cm.Greys_r)
36     plt.title('Rank 512 Boat')
```

```
38     _ = plt.subplots_adjust(wspace=0.5)
39     plt.show()
40
```

Page 2 of 2





```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 # May want to center the matrix before doing stuff
6 # This is a trick you may want to do when pre-
7 # processing your dataset
8
9 with open('data/AbileneFlows/odnames','r') as f:
10     odnames = [line.strip() for line in f]
11 dates = pd.date_range('9/1/2003', freq='10min', periods=
12                         1008)
13 Atraf = pd.read_table('data/AbileneFlows/X', sep=' ', 
14                       header=None, names=odnames, engine='python')
15 Atraf.index = dates
16 print(Atraf.shape)
17
18 print("rank of our matrix is ", np.linalg.matrix_rank(
19     Atraf))
20
21 u,s,vt = np.linalg.svd(Atraf)
22
23 fig = plt.figure(figsize=(6,4))
24 plt.plot(range(1,1+len(s)),s)
25 plt.xlabel(r'$k$',size=20)
26 plt.ylabel(r'$\sigma_k$',size=20)
27 _ = plt.title(r'Singular Values of $A$',size=20)
28 plt.show()
29
30 fig = plt.figure(figsize=(6,4))
31 Anorm = np.linalg.norm(Atraf)
32 plt.plot(range(1,21),s[0:20]/Anorm)
33 plt.xlabel(r'$k$',size=20)
34 plt.ylabel(r'$\sigma_k$',size=20)
35 plt.show()
36
37 fig = plt.figure(figsize=(6,4))
38 Anorm = np.linalg.norm(Atraf)
39 err = np.cumsum(s[::-1]**2)
40 err = np.sqrt(err[::-1])
```

```
38 plt.plot(range(1,21),err[:20]/Anorm)
39 plt.xlim([0,20])
40 plt.xlabel(r'$k$',size=16)
41 plt.ylabel(r'relative F-norm error',size=16)
42 _ = plt.title(r'Relative Error of rank-$k$ approximation to $A$',size=16)
43 plt.show()
44
45
46
```

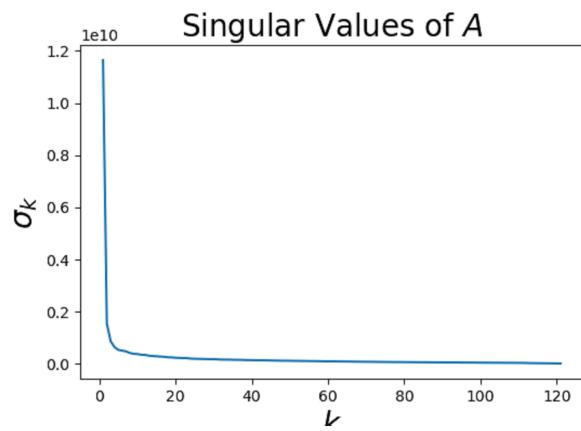
```

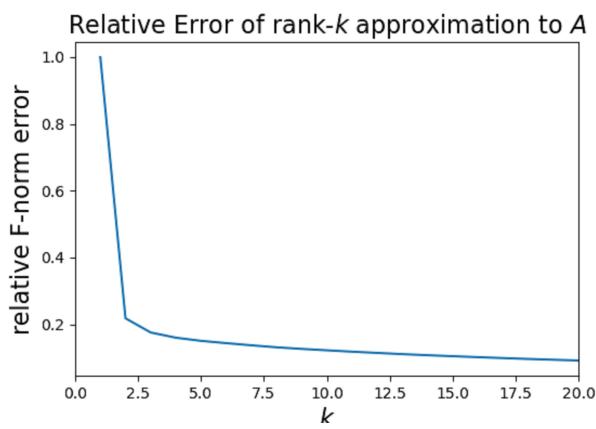
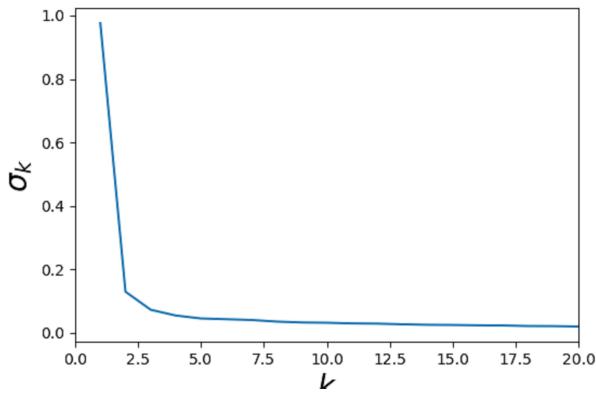
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 # May want to center the matrix before doing stuff
6 # This is a trick you may want to do when pre-
  processing your dataset
7
8 with open('data/AbileneFlows/odnames','r') as f:
9     odnames = [line.strip() for line in f]
10 dates = pd.date_range('9/1/2003', freq='10min', periods=
  1008)
11 Atraf = pd.read_table('data/AbileneFlows/X', sep=' ', 
  header=None, names=odnames, engine='python')
12 Atraf.index = dates
13 print(Atraf.shape)
14
15 print("rank of our matrix is ", np.linalg.matrix_rank(
  Atraf))
16
17 u,s,vt = np.linalg.svd(Atraf)
18
19 fig = plt.figure(figsize=(6,4))
20 plt.plot(range(1,1+len(s)),s)
21 plt.xlabel(r'$k$',size=20)
22 plt.ylabel(r'$\sigma_k$',size=20)
23 _ = plt.title(r'Singular Values of $A$',size=20)
24 plt.show()
25
26 fig = plt.figure(figsize=(6,4))
27 Anorm = np.linalg.norm(Atraf)
28 plt.plot(range(1,21),s[0:20]/Anorm)
29 plt.xlim([0,20])
30 plt.xlabel(r'$k$',size=20)
31 _ = plt.ylabel(r'$\sigma_k$',size=20)
32 plt.show()
33
34 fig = plt.figure(figsize=(6,4))
35 Anorm = np.linalg.norm(Atraf)
36 err = np.cumsum(s[::-1]**2)
37 err = np.sqrt(err[::-1])

```

```
File - C:\Users\Wolfs\PycharmProjects\cs506lec20211006\relative-error.py
38 plt.plot(range(1,21),err[:20]/Anorm)
39 plt.xlim([0,20])
40 plt.xlabel(r'$k$',size=16)
41 plt.ylabel(r'relative F-norm error',size=16)
42 _ = plt.title(r'Relative Error of rank-$k$ approximation to $A$',size=16)
43 plt.show()
44
45
46
```

Page 2 of 2

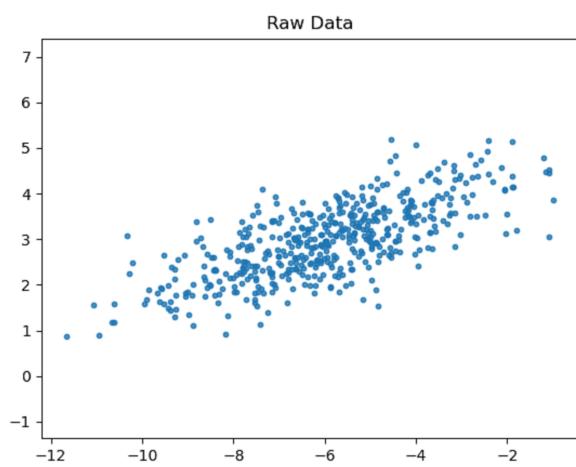




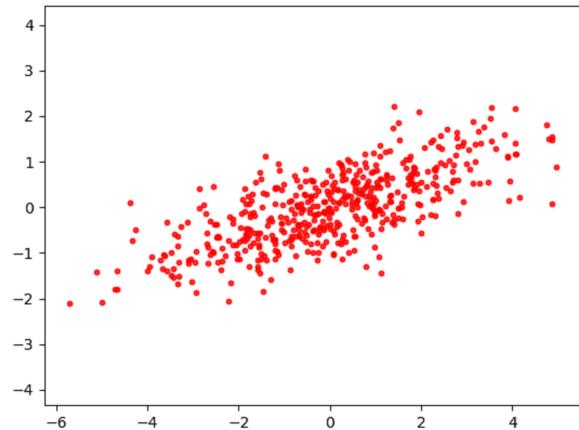
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # You want to normalize your data first, some features
5 # may have scales that are pretty incompatible
6 # Mean center your data
7 # We are using numpy as our linear algebra library
8 # Singular values: 2 values (we only had two
9 # dimensions)
10 # The points in blue in the last plot are the rank 1
11 # approximation of the points in red projected onto the
12 # green line
13 # The feature from this column (the first, the
14 # strongest) may more accurately represent the dataset
15 # So then we end up with the dataset with the blue
16 # points
17
18 n_samples = 500
19 C = np.array([[0.1, 0.6], [2., .6]])
20 X = np.random.randn(n_samples, 2) @ C + np.array([-6,
21 3])
22 plt.axis('equal')
23 plt.scatter(X[:, 0], X[:, 1], s=10, alpha=0.8)
24 plt.title("Raw Data")
25 plt.show()
26
27 Xc = X - np.mean(X, axis=0)
28 plt.axis('equal')
29 plt.scatter(Xc[:, 0], Xc[:, 1], s=10, alpha=0.8, color
30 ='r')
31 plt.title("Mean Centered Data")
32 plt.show()
33
34 scope = s.copy()
35 scope[1] = 0.
36 reducedX = u @ np.diag(scope) @ vt
```

```
35 plt.axis('equal')
36 plt.scatter(Xc[:,0],Xc[:,1], color='r')
37 plt.scatter(reducedX[:,0], reducedX[:,1])
38 endpoints = np.array([[-10],[10]]) @ vt[[0],:]
39 _ = plt.plot(endpoints[:,0], endpoints[:,1], 'g-')
40 plt.show()
41
42
```

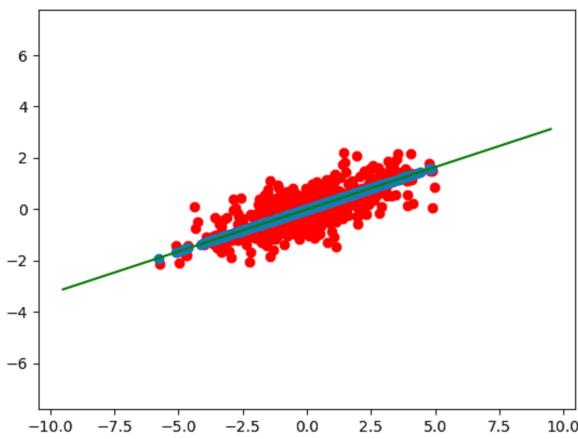
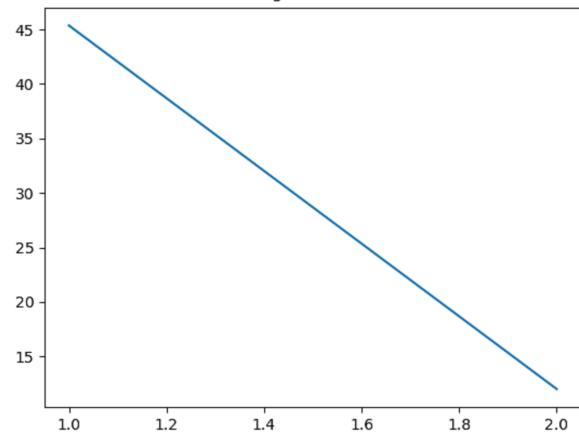
Page 2 of 2



Mean Centered Data



Singular Values



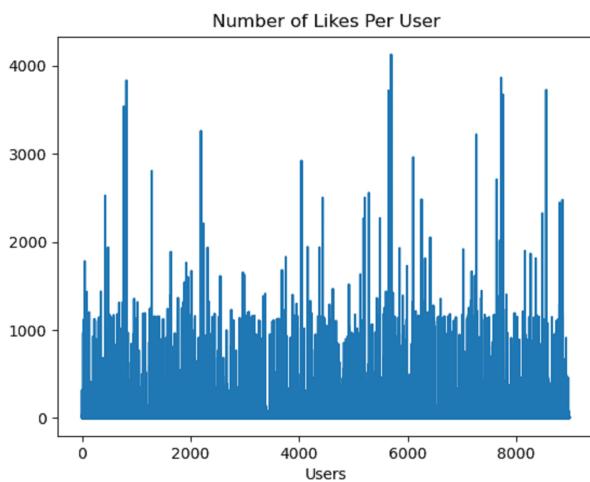
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # The most anomalous users are not what is expected
5 # Singular value decomposition can pinpoint anomalous
6 # users really quickly, but could be wrong
7
8 data = np.loadtxt('data/spatial_data.txt')
9 # This data consists of the number of 'Likes'
10 # during a six month period, for each of 9000
11 # users across the 210 content categories that
12 # Facebook assigns to pages.
13
14 FBSpatial = data[:,1:]
15 FBSnorm = np.linalg.norm(FBSpatial, axis=1, ord=1)
16 plt.plot(FBSnorm)
17 _ = plt.xlabel('Users')
18 plt.show()
19
20 u,s,vt = np.linalg.svd(FBSpatial, full_matrices=False)
21 plt.plot(s)
22 _ = plt.title('Singular Values of Spatial Like Matrix')
23 plt.show()
24
25 scopy = s.copy()
26 scopy[25:] = 0.
27 N = u @ np.diag(scopy) @ vt
28 O = FBSpatial - N
29 Onorm = np.linalg.norm(O, axis=1)
30 anomSet = np.argsort(Onorm)[-30:]
31 plt.plot(Onorm)
32 plt.plot(anomSet, Onorm[anomSet], 'ro')
33 _ = plt.title('Norm of Residual (rows of O)')
34 plt.show()
35
36 anomSet = np.argsort(Onorm)[-30:]
37 plt.plot(FBSnorm)
38 plt.plot(anomSet, FBSnorm[anomSet], 'ro')
39 _ = plt.title('Top 30 Anomalous Users - Total Number')
```

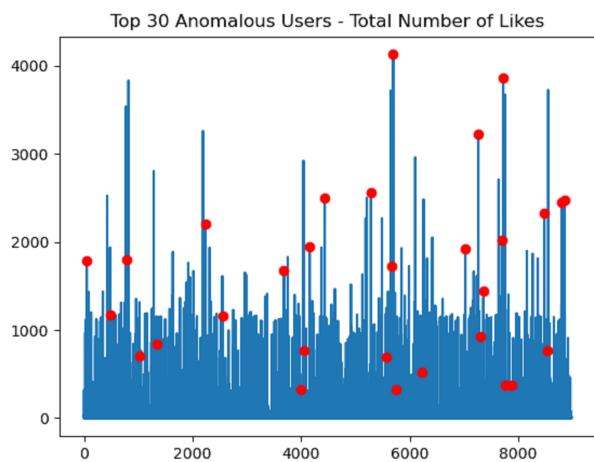
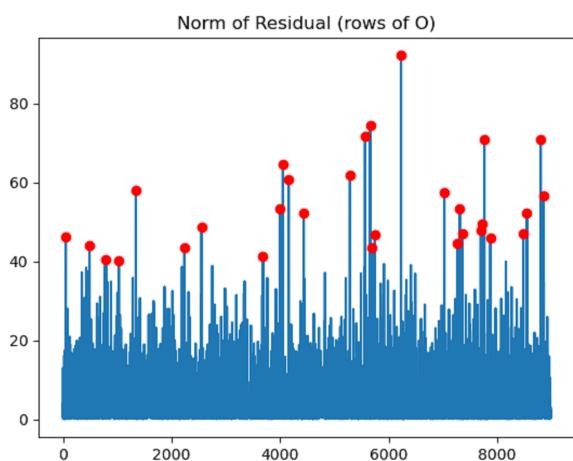
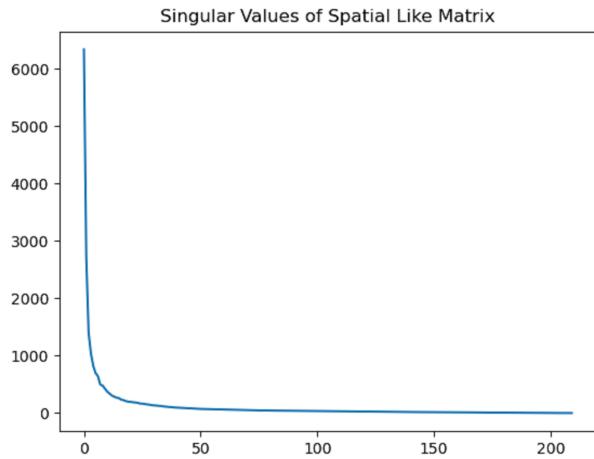
```

File - C:\Users\Wolfs\PycharmProjects\ics506lec20211006\svd-anomaly-det.py
39 of Likes')
40 plt.show()
41
42 # anomalous users
43 plt.figure(figsize=(9,6))
44 for i in range(1,10):
45     ax = plt.subplot(3,3,i)
46     plt.plot(FBSpatial[anomSet[i-1],:])
47     plt.xlabel('FB Content Categories')
48 plt.subplots_adjust(wspace=0.25,hspace=0.45)
49 _ = plt.suptitle('Nine Example Anomalous Users',size=20)
50 plt.show()
51
52 # normal users
53 set = np.argsort(Onorm)[0:7000]
54 # that have high overall volume
55 max = np.argsort(FBSnorm[set])[::-1]
56 plt.figure(figsize=(9,6))
57 for i in range(1,10):
58     ax = plt.subplot(3,3,i)
59     plt.plot(FBSpatial[set[max[i-1]],:])
60     plt.xlabel('FB Content Categories')
61 plt.subplots_adjust(wspace=0.25,hspace=0.45)
62 _ = plt.suptitle('Nine Example Normal Users',size=20)
63 plt.show()
64
65

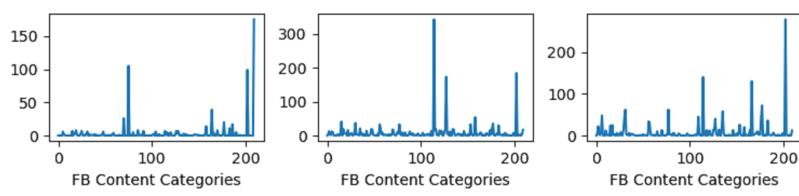
```

Page 2 of 2

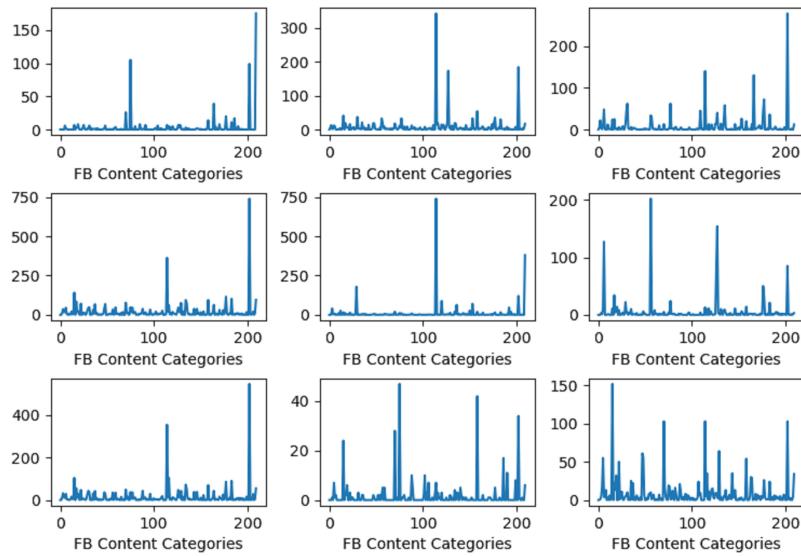




Nine Example Anomalous Users



Nine Example Anomalous Users



Nine Example Normal Users

