# Textual analysis of movie reviews Team5

BY

**HARSH PATHAK**

**JIDAPA THADAJARASSIRI**

**KRUSHIKA TAPEDIA**

**PITCHAYA  WIRATCHOTISATIAN**

**PRINCE SHIVA CHAUDHARY**

# Emotions

http://webneel.com/daily/20-inside-out-characters
http://www.thecoli.com/threads/ios-emojis-degrade-and-simplify-human-expression-of-emotion.51174/
https://www.willbrattcounselling.com/blog-creating-difference/2015/1/12/your-emotions-arent-a-problem
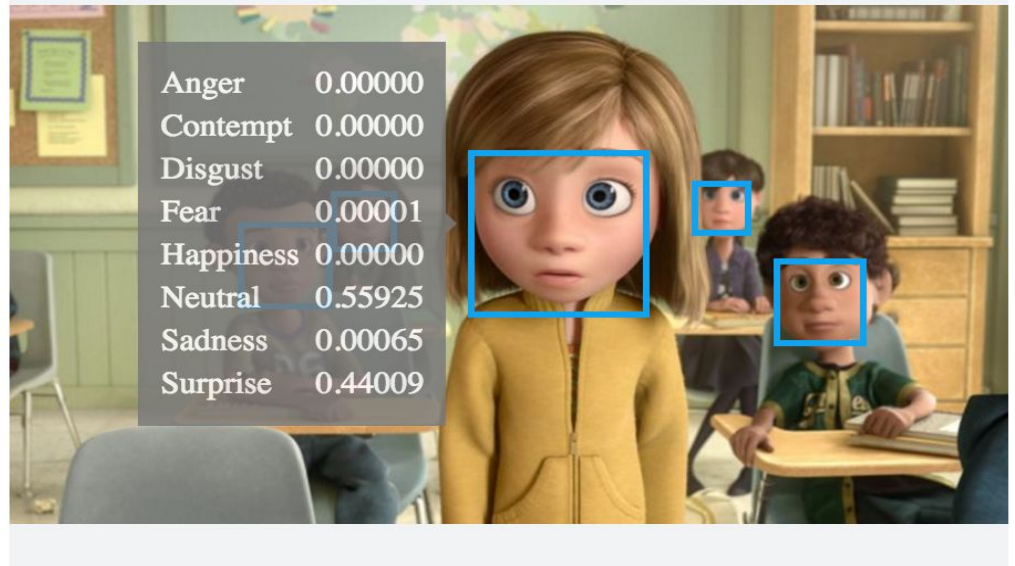
# Emotions are difficult to express in words

# Machine Learning to understand Emotions

**Two Approaches**
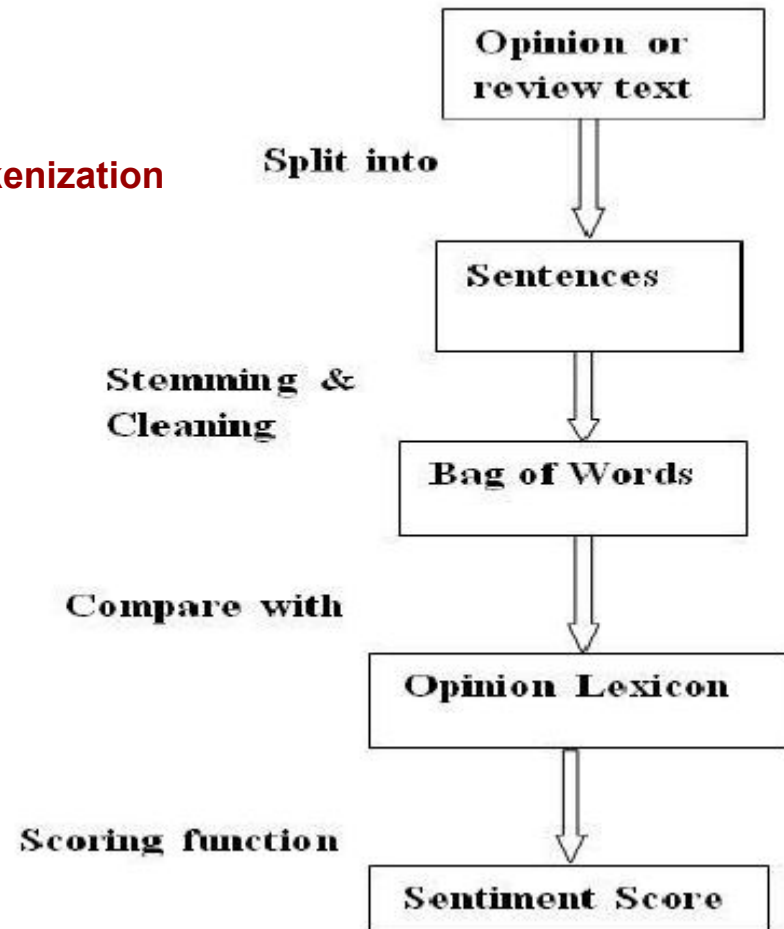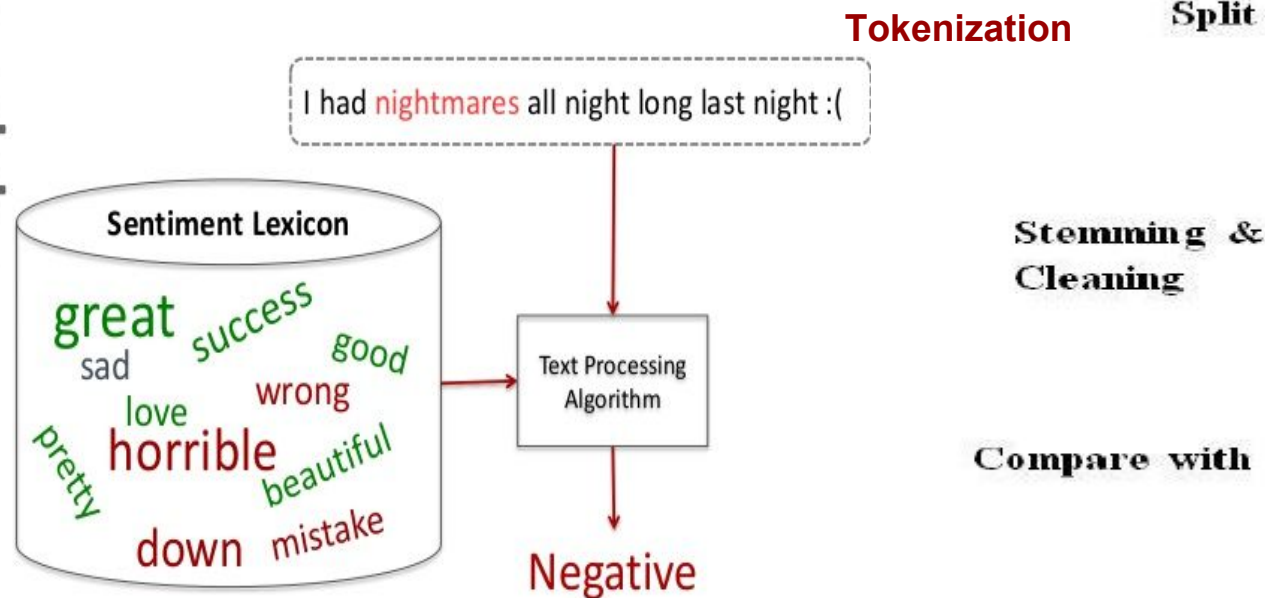
- Lexicon based approach

- Machine learning based approach



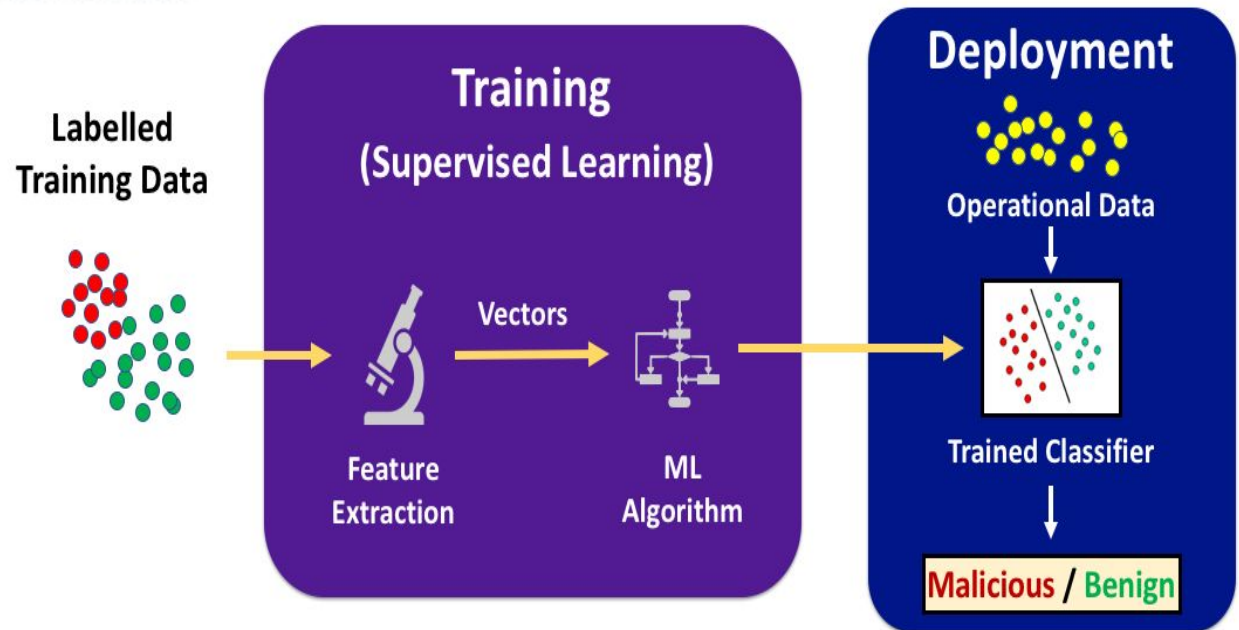| Anger | 0.00000 |
| Contempt | 0.00000 |
| Disgust | 0.00000 |
| Fear | 0.00001 |
| Happiness | 0.00000 |
| Neutral | 0.55925 |
| Sadness | 0.00065 |
| Surprise | 0.44009 |

http://thenextweb.com

# Lexicon-based

**Approach**

**Tokenization**

I had nightmares all night long last night :(

Sentiment Lexicon

great success
sad        good
love   wrong
pretty  horrible  beautiful
down  mistake

Text Processing Algorithm

Negative

Opinion or review text

Split into

Sentences

Stemming & Cleaning

Bag of Words

Compare with

Opinion Lexicon

Scoring function

Sentiment Score

# Machine Learning Approach
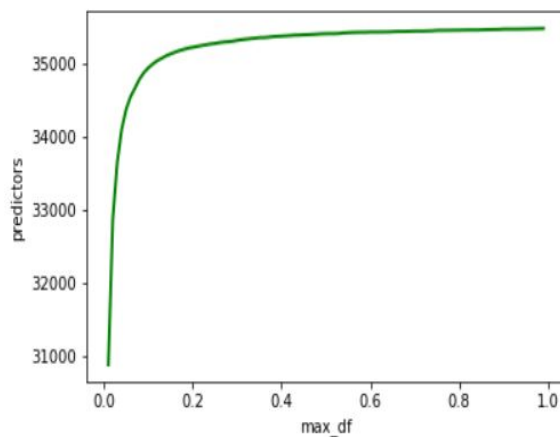
# Introduction

## I. Sentiment Analysis
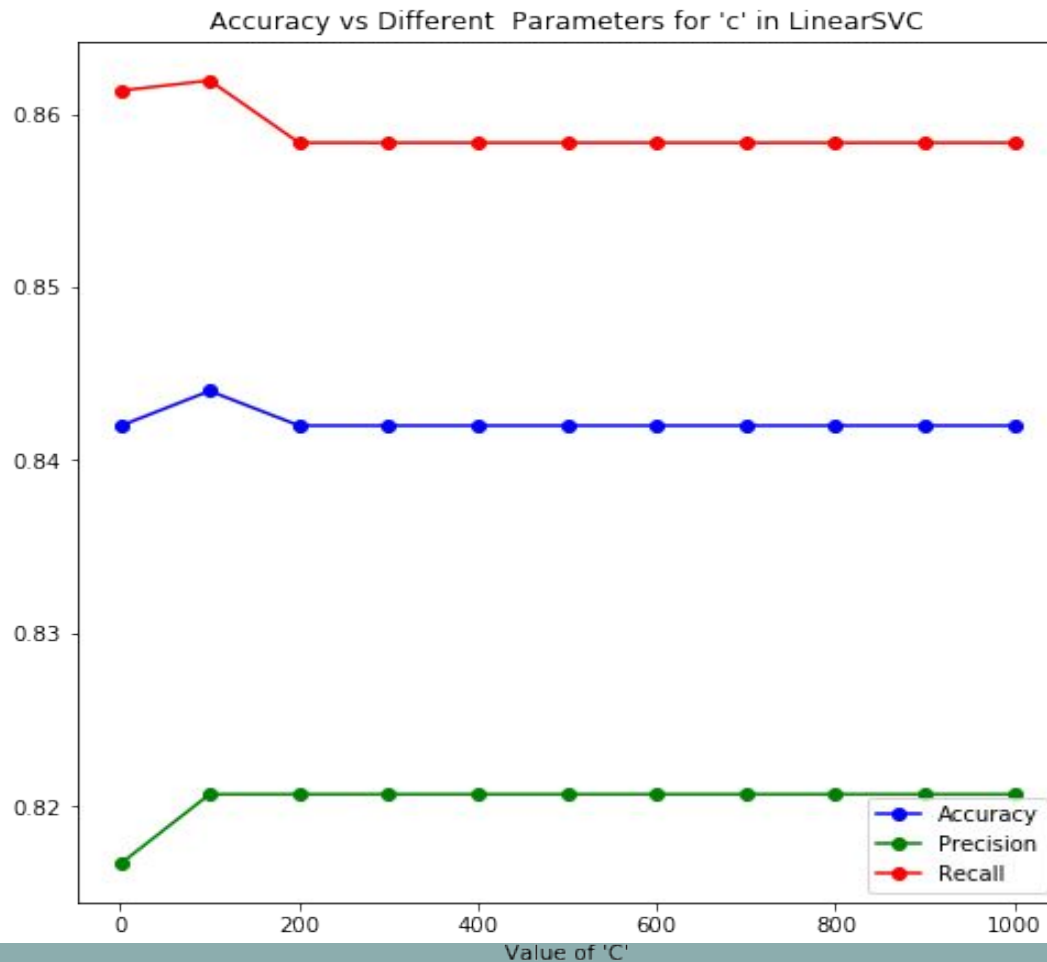
## II. TfidfVectorizer

### min_df

### max_df

### ngram_range

# Machine Learning Algorithms (LinearSVC)

## Performance of LinearSVC : Parameter C



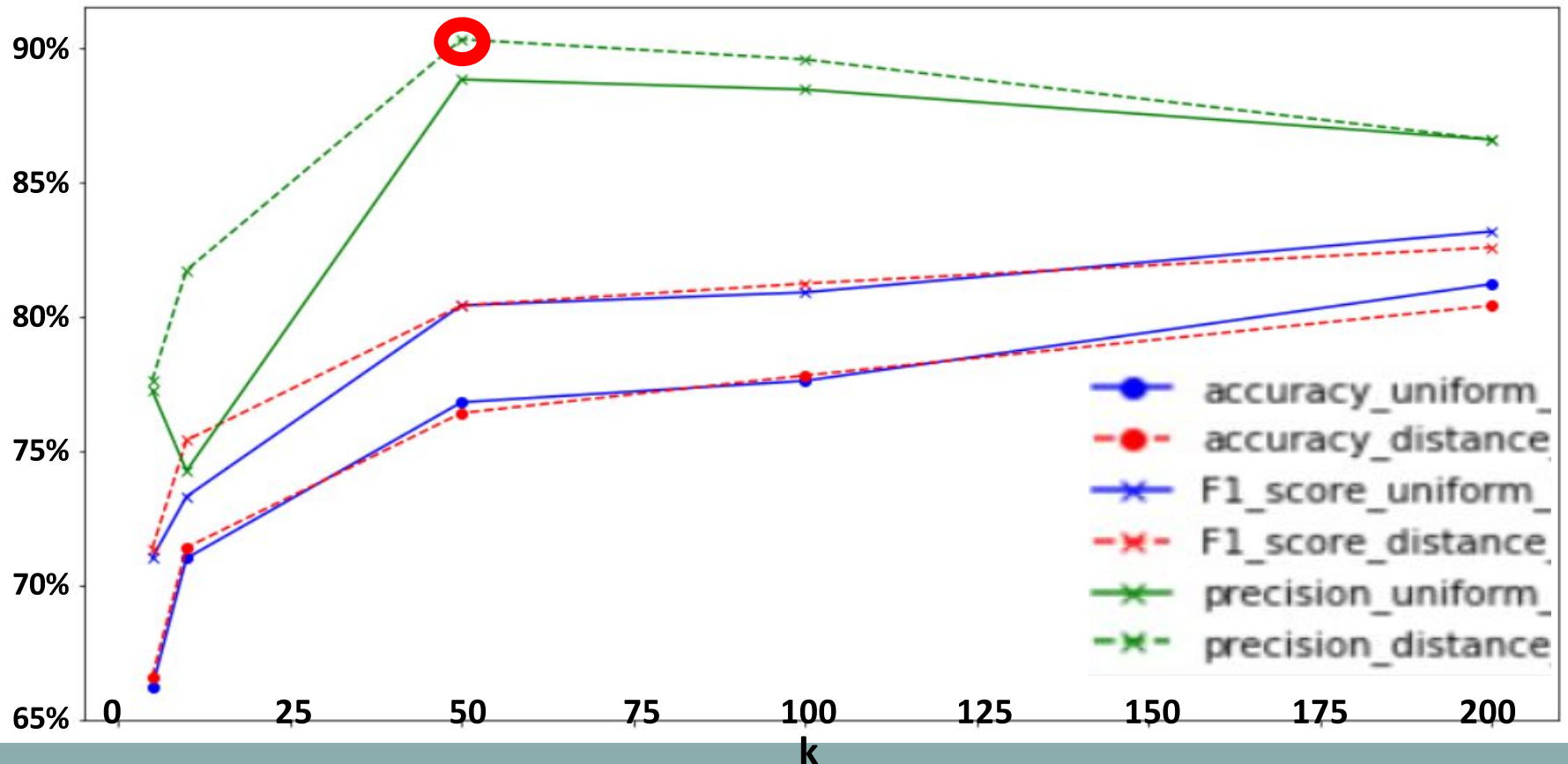Accuracy vs Different Parameters for 'c' in LinearSVC

# Machine Learning Algorithms (kNN)

## Performance of KneighborClassifier

Parameter :  k = 5, 10, 50, 100, 200

Parameter : weight function = 'uniform' or 'distance'

# Finding the right plot (1)

## Step 1: Data Preprocessing
- Remove stopwords : NLTK
- Extract stem-words : snowball
- Create TF-IDF vector matrix : TfidfVectorizer

## Step 2: Feature Selection
- Logistic regression with lasso
- Linear model with lasso
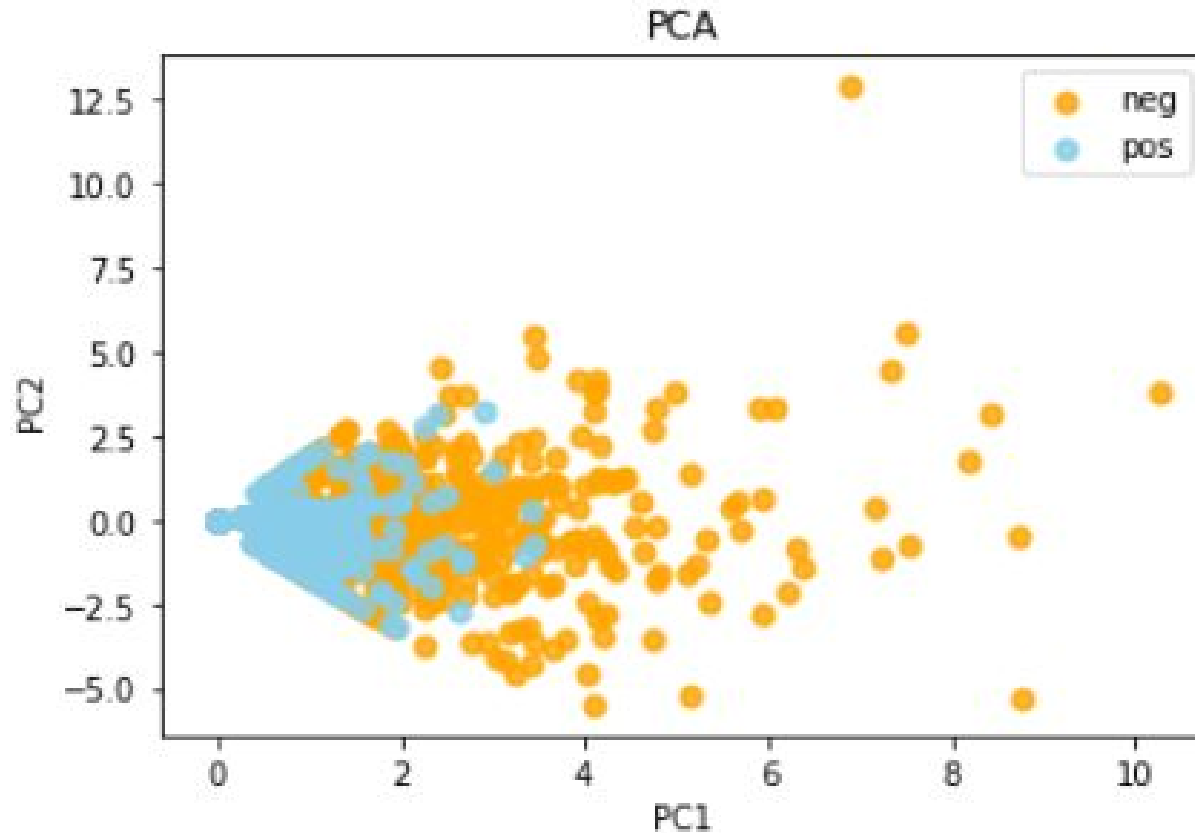- TruncatedSVD: first 2 components, 256 components

## Step 3: Methodologies
1. Feature selection: TruncatedSVD, Lasso
2. Clustering: K-Means clustering, Hierarchical clustering
3. Ensemble learning: RandomTreesEmbedding
4. Manifold learning: MDS, Isomap, Spectral decomposition, Locally Linear Embedding, t-SNE

# Finding the right plot (2)

## Method 1: TruncatedSVD



The first two PCs of features in LinearLassoIndex
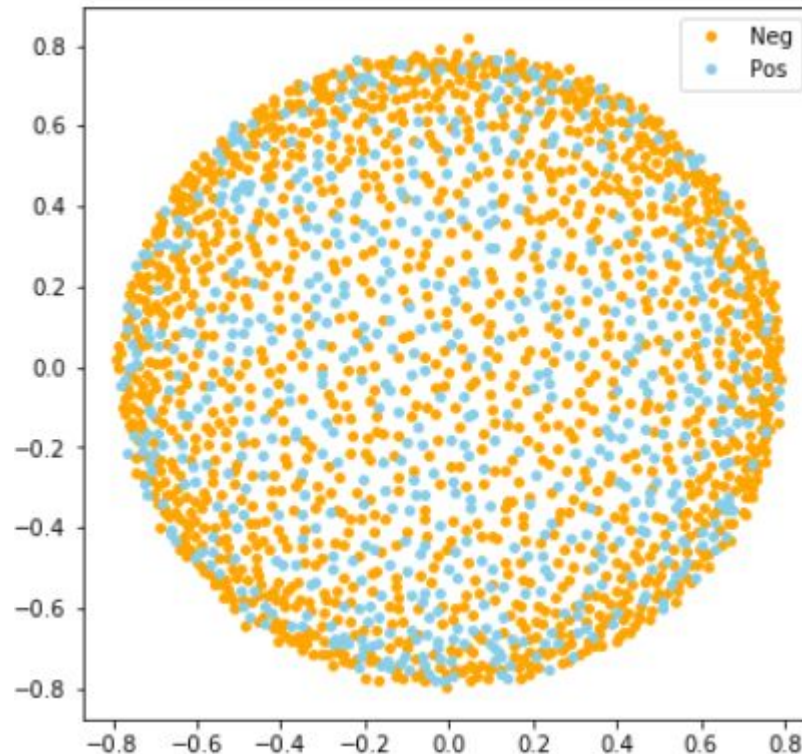
# Finding the right plot (3)

## Method 2: Lasso



Parallel coordinates plot of features from LinearLassoIndex

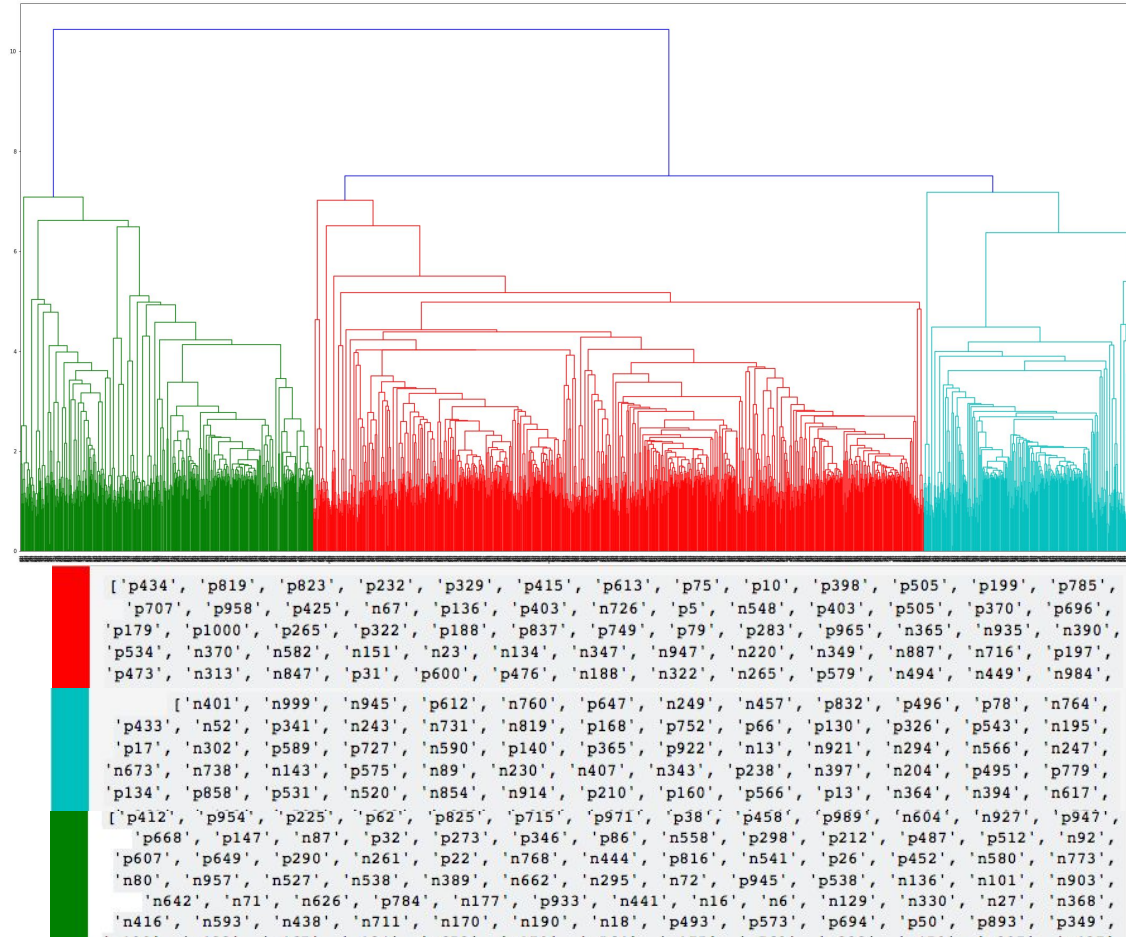# Finding the right plot (4)

## Method 3: K-Means Clustering



K-mean of the MDS(n_components =2 ) of the distance matrix where n_clusters = 2
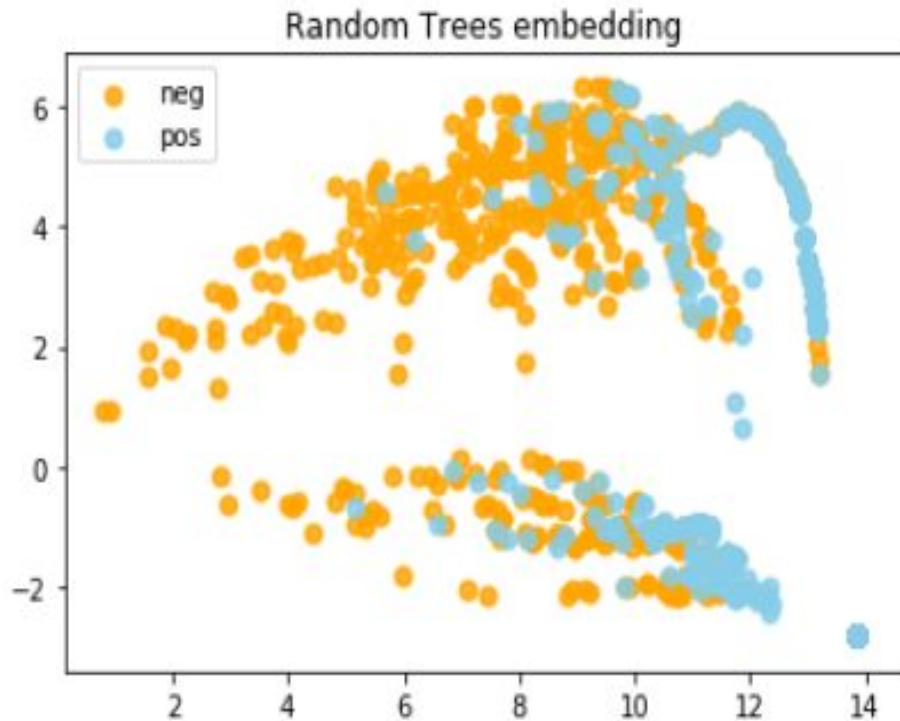
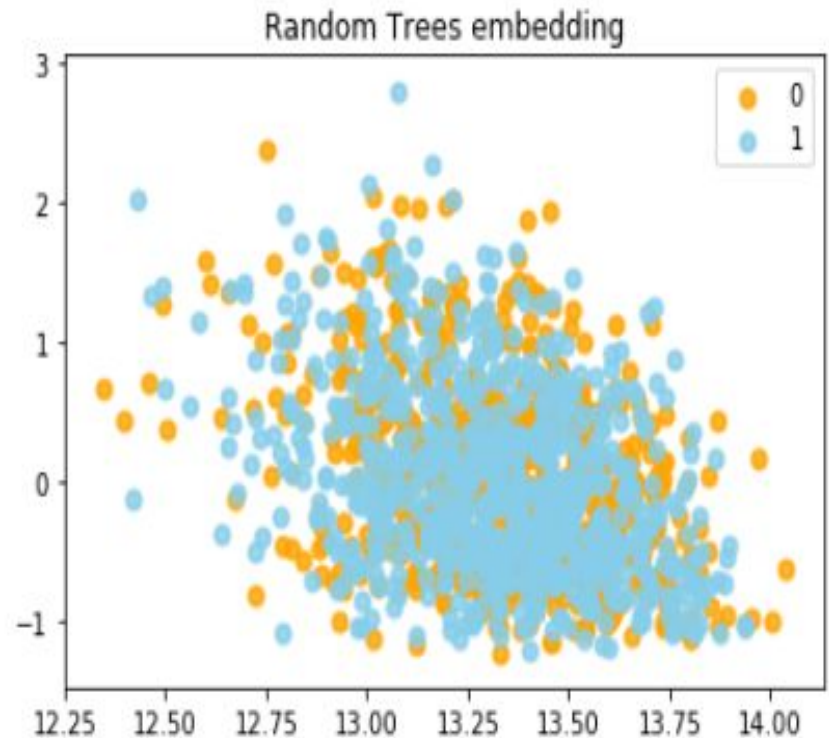# Finding the right plot (5)

**Method 4: Hierarchical Clustering**

# Finding the right plot (6)

## Method 5: RandomTreesEmbedding



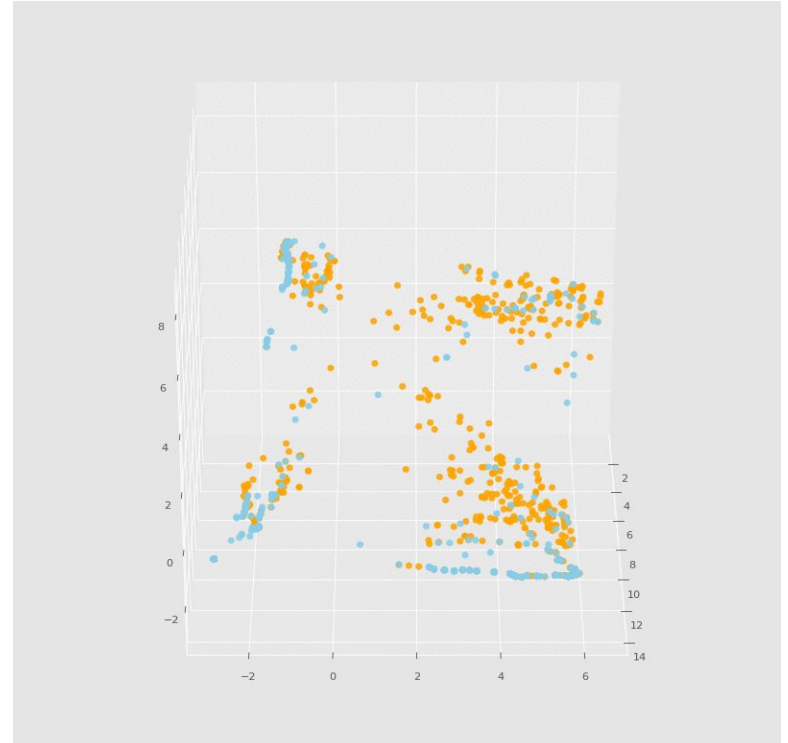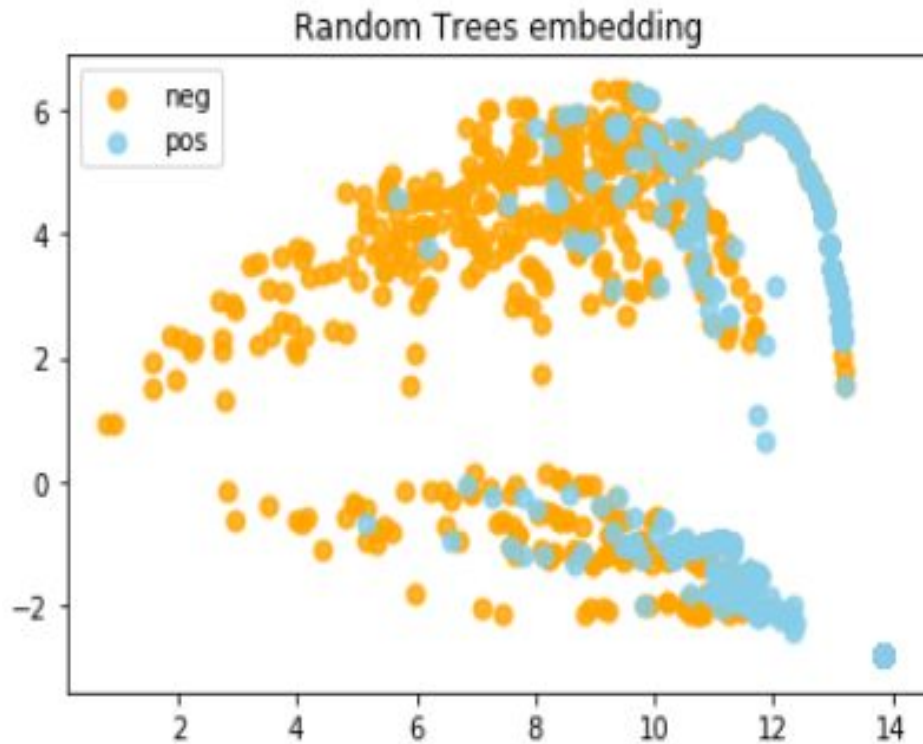n_estimators = 200, max_depth = 5, and features in LinearLassoIndex

TruncatedSVD(256)

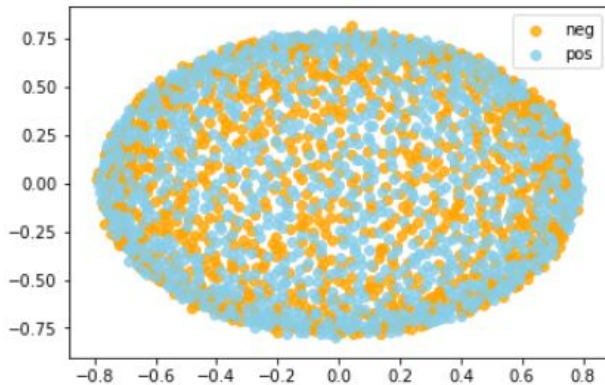# Finding the right plot (6)

## Method 5: RandomTreesEmbedding



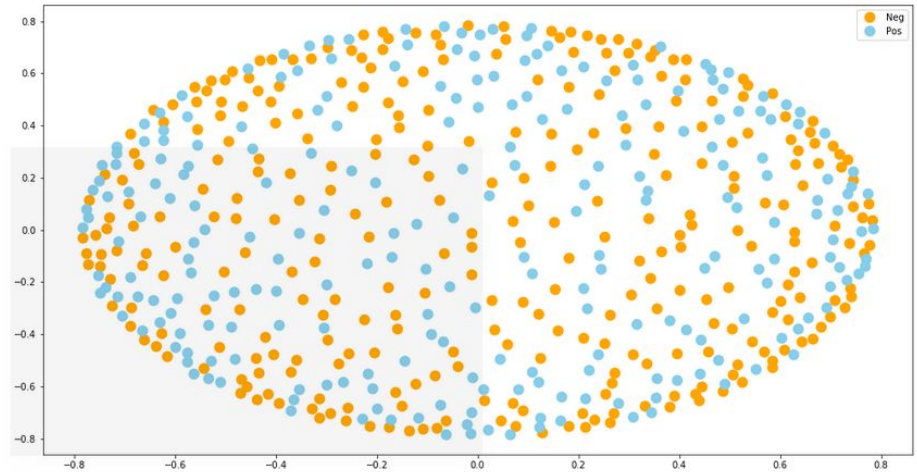n_estimators = 200, max_depth = 5, and features
in LinearLassoIndex

# **Finding the right plot (7)**

## **Method 6: MDS**
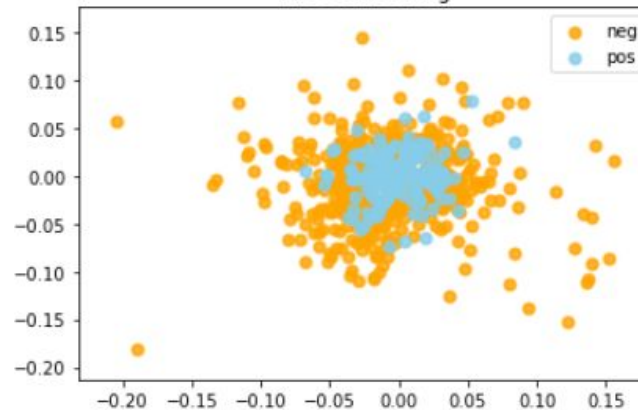
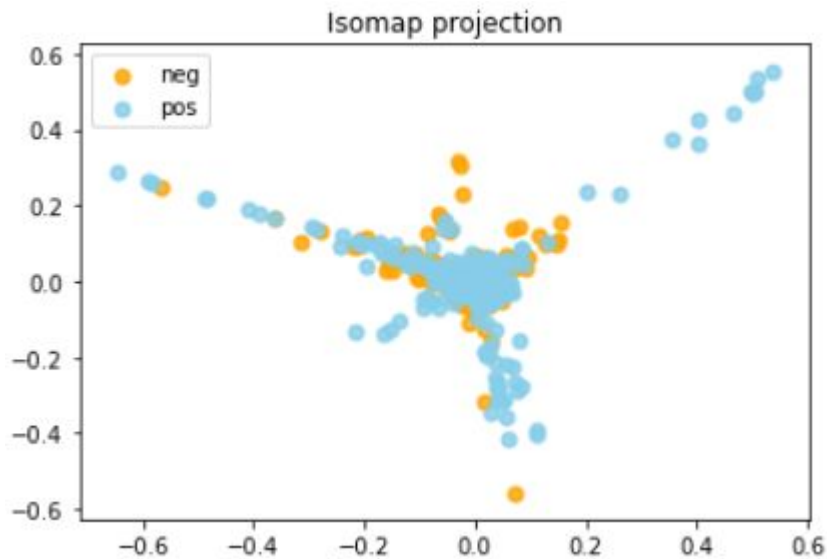Distance matrix computed
from cosine similarity



TruncatedSVD(256)
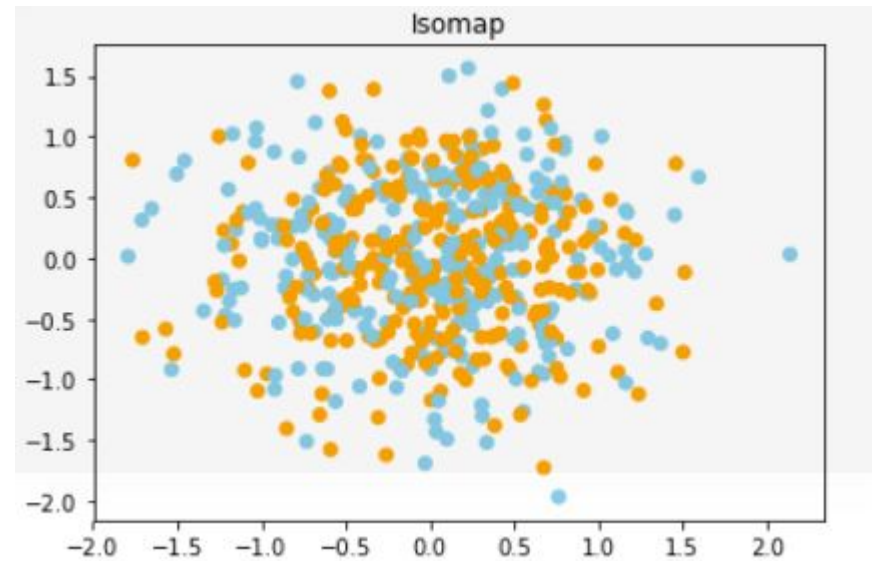


MDS applied to features in
LinearLassoIndex

# Finding the right plot (8)

## Method 7: Isomap Projection



n_neighbors = 100, n_components = 2, and
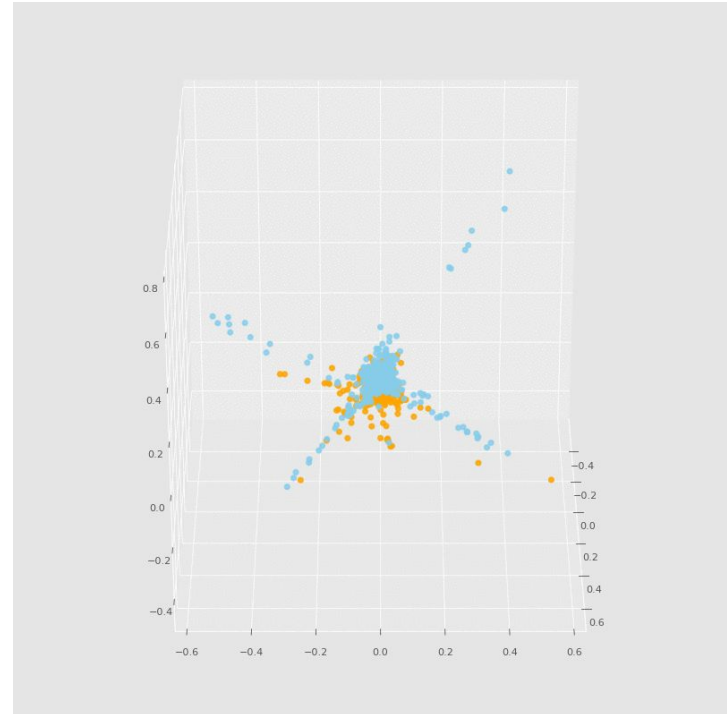features in LogRegLassoIndex



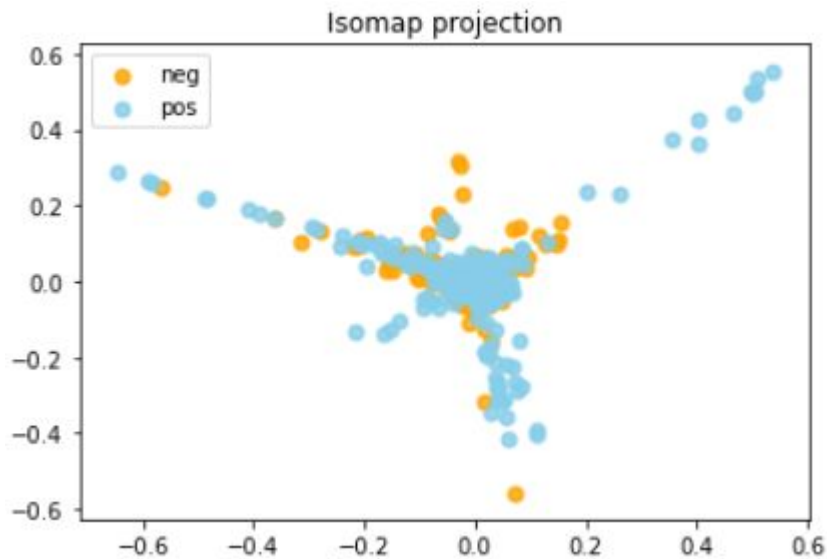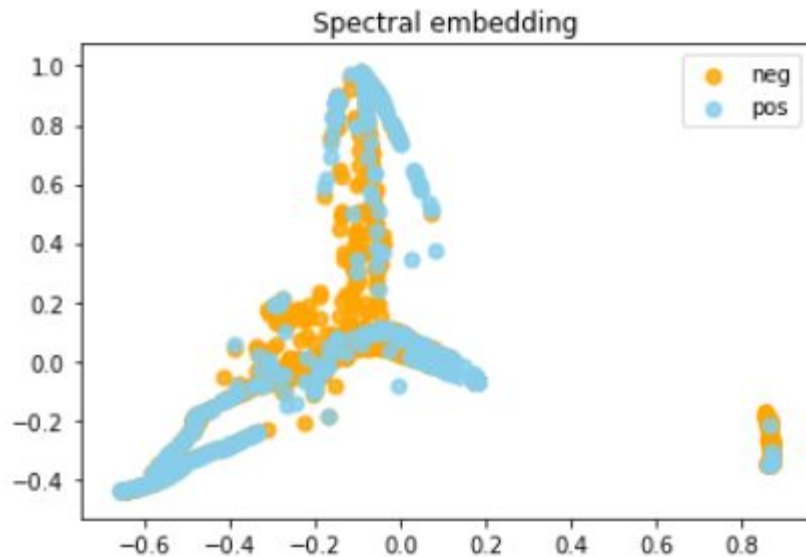TruncatedSVD(256)

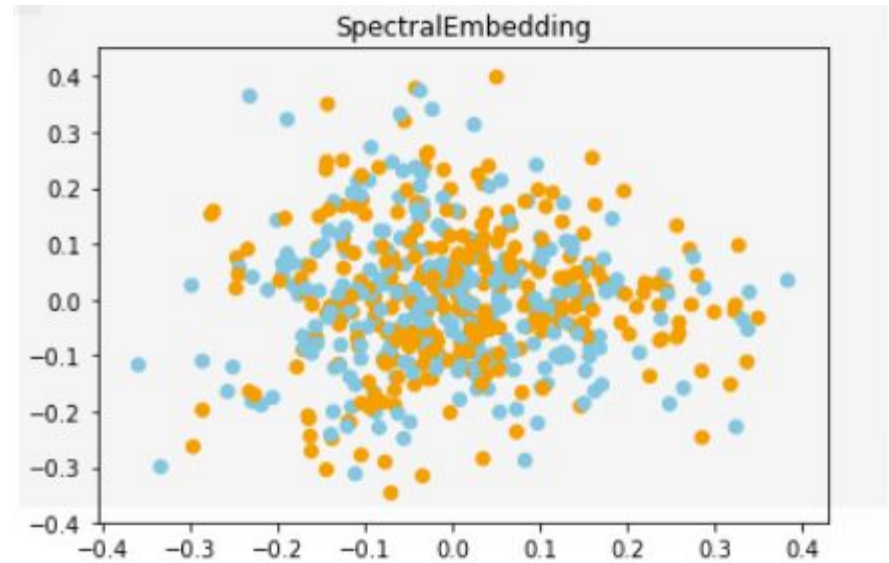# Finding the right plot (8)

## Method 7: Isomap



n_neighbors = 100, n_components = 2, and features in LogRegLassoIndex

## Method 8: Spectral Embedding



eigen_solver = 'arpack' and
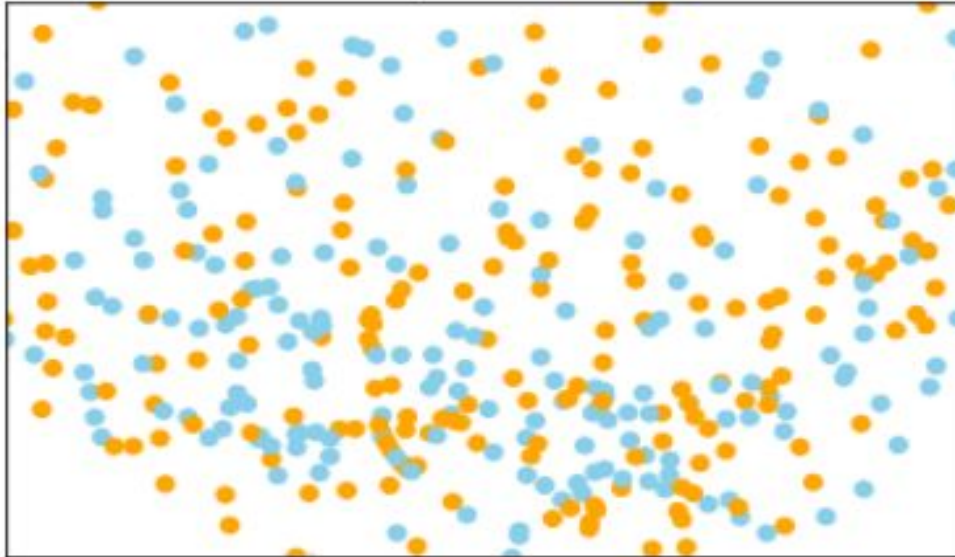features in LinearLassoIndex

TruncatedSVD(256)

Apply t-SNE with init = 'pca' to the data matrix

## Method 9: Locally Linear Embedding

manifold.locally_linear_embedding(data_svd, n_neighbors=2, n_components=2)
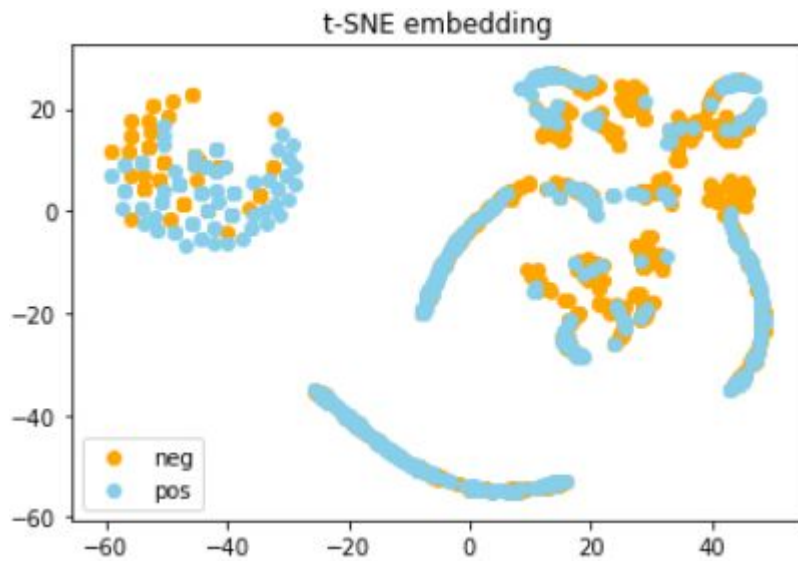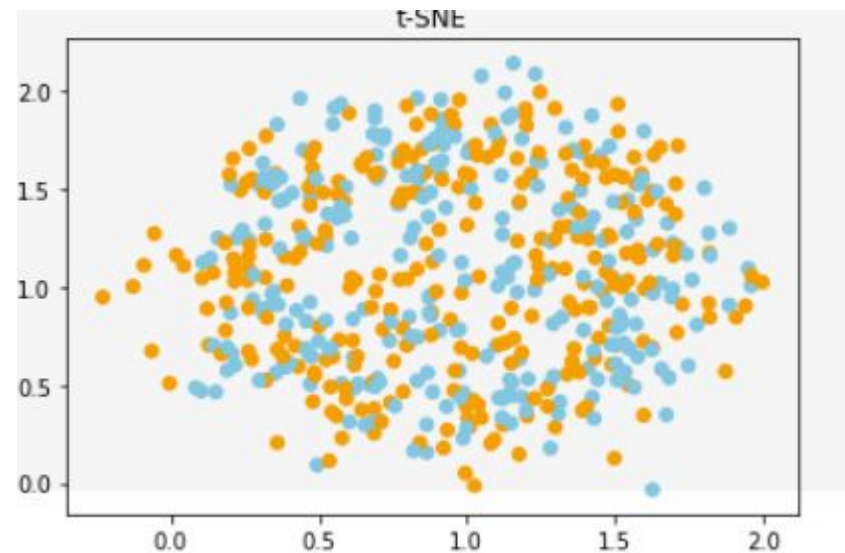


Projected data

# Finding the right plot (11)

## Method 10: t-SNE

init = 'pca' and features in LinearLassoIndex

TruncatedSVD(256)

# Thank you

# Q & A