# Working With Text
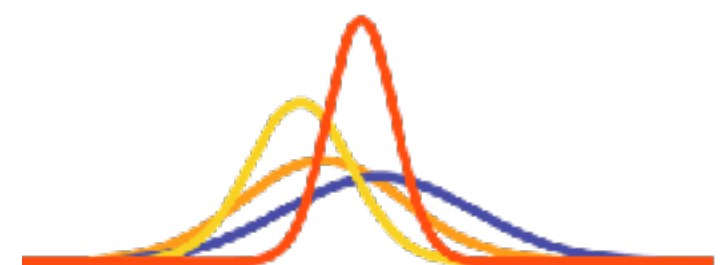
## Pittsburgh useR Group - 04/24/2018

Michael Patnik
AvenueFour Analytics
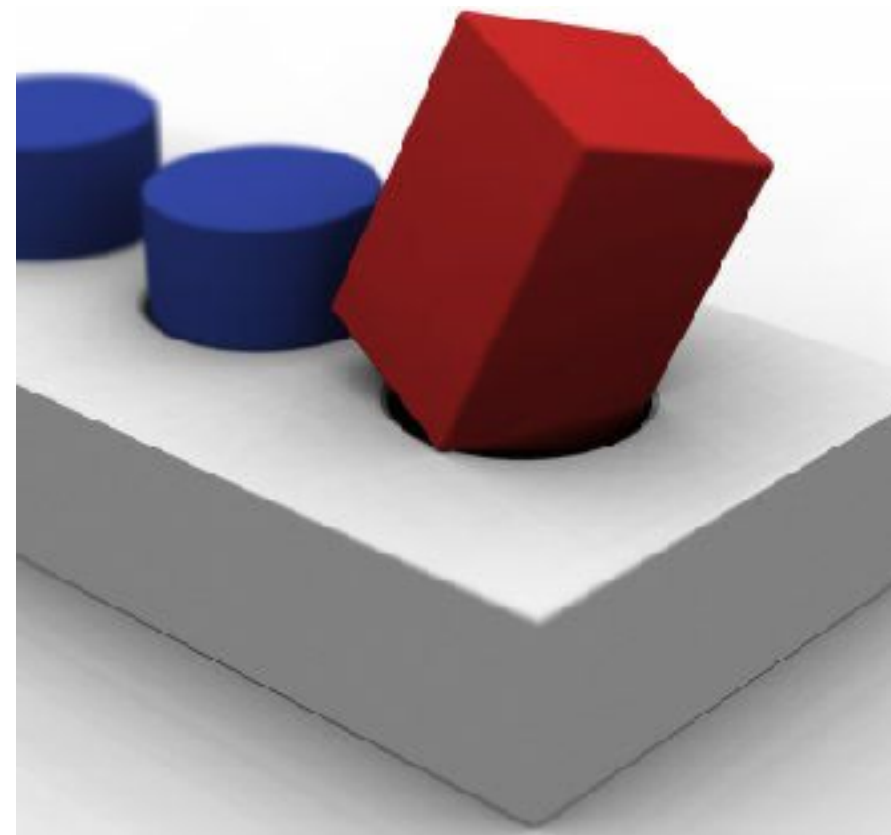
**Avenue 4 Analytics**

# Challenges



Text Is dirty!

Incompatible with most statistical models.

# Correctly Identifying Entities

# Similarity Algorithms

**Levenshtein Distance**

- minimum number of edits necessary to change one string into another

$$\text{lev}_{a,b}(i,j) = \begin{cases} \max(i,j) & \text{if} \min(i,j) = 0, \\ \min \begin{cases} \text{lev}_{a,b}(i-1,j)+1 \\ \text{lev}_{a,b}(i,j-1)+1 \\ \text{lev}_{a,b}(i-1,j-1)+1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

**Jaccard Index (using ngrams)**

- count of the intersection of ngrams divided by count of the union of ngrams

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

unigram

| C | O | L | D |

bigram

trigram

n-gram (n = 4)

C O L D

# Levenshtein Distance

'K C Leung' → ' C Leung' → 'C Leung'     **2** ✅

# Levenshtein Distance

'Derek K C Leung' → 'Derek K C Leun' → 'Derek K C Leu'

'Derek K C ' ← 'Derek K C L' ← 'Derek K C Le'

'Derek K C' → 'LDerek K C' → 'LeDerek K C'

'LeungDerek K C' ← 'LeunDerek K C' ← 'LeuDerek K C'

'Leung,Derek K C' → 'Leung, Derek K C'  **13** ✗

# Jaccard Index (using ngrams)

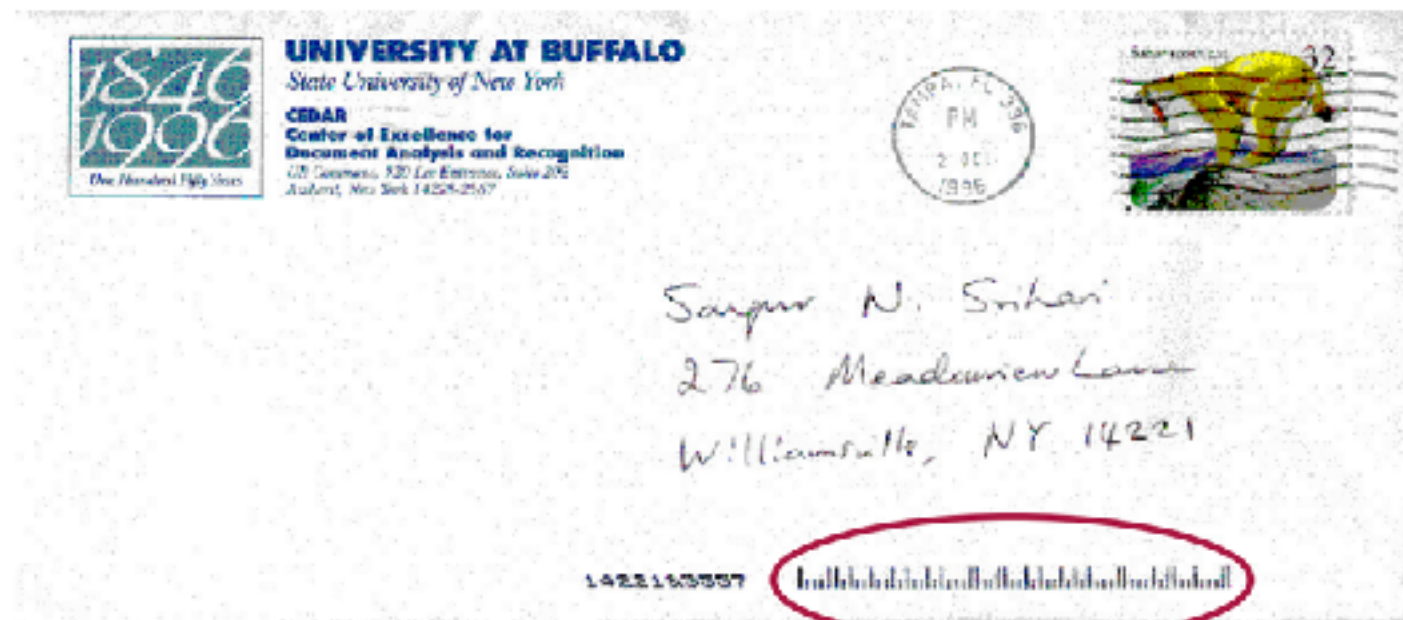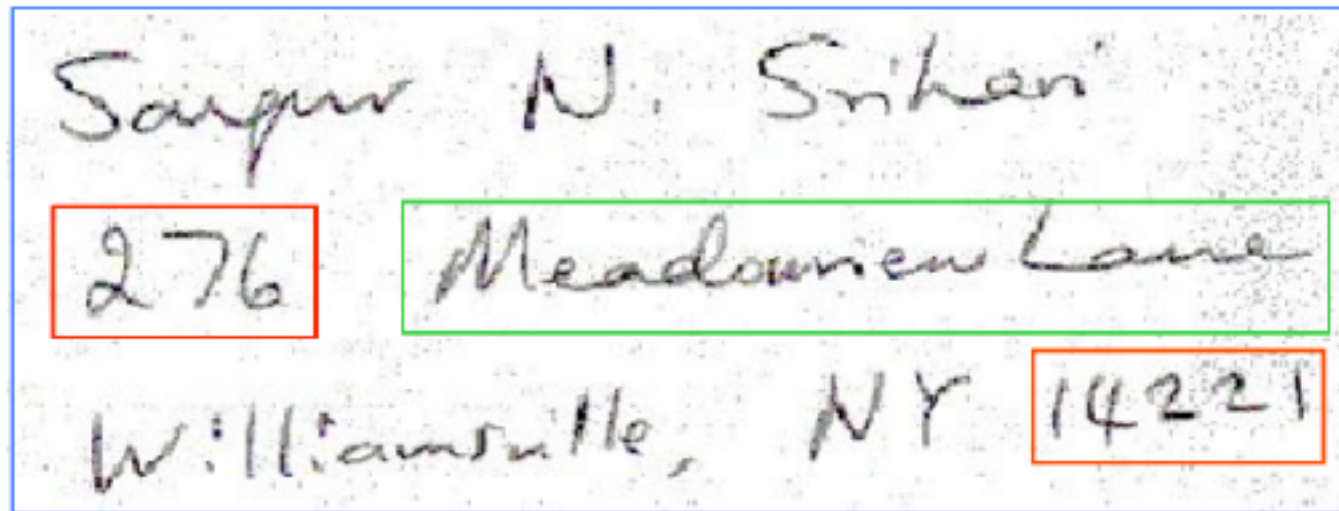| | | |
|---|---|---|
| 'K C Leung' | 'K ', ' C', 'C ', ' L', 'Le', 'eu', 'un','ng' | Intersection:<br>'C ', ' L', 'Le', 'eu', 'un','ng'<br><br>Union:<br>'K ', ' C', 'C ', ' L', 'Le', 'eu', 'un','ng'<br><br>Score:<br>6 / 8 = 0.75 |
| 'C Leung' | 'C ', ' L', 'Le', 'eu', 'un','ng' | |
| 'Derek K C Leung' | 'De', 'er', 're', 'ek', 'k ', ' K', 'K ', ' C', 'C ', ' L', 'Le', 'eu', 'un','ng' | Intersection:<br>'De', 'er', 're', 'ek', 'k ', ' K', 'K ', ' C', 'Le', 'eu', 'un','ng'<br>Union:<br>'De', 'er', 're', 'ek', 'k ', ' K', 'K ', ' C', 'C ', ' L', 'Le', 'eu', 'un','ng', 'g,', ', ', ' D'<br>Score:<br>12 / 17 = 0.71 |
| 'Leung, Derek K C' | 'Le', 'eu', 'un','ng', 'g,', ', ', ' D', 'De', 'er', 're', 'ek', 'k ', ' K', 'K ', ' C' | |

# Leverage Domain Knowledge and Other Data Elements

# Case Study: USPS



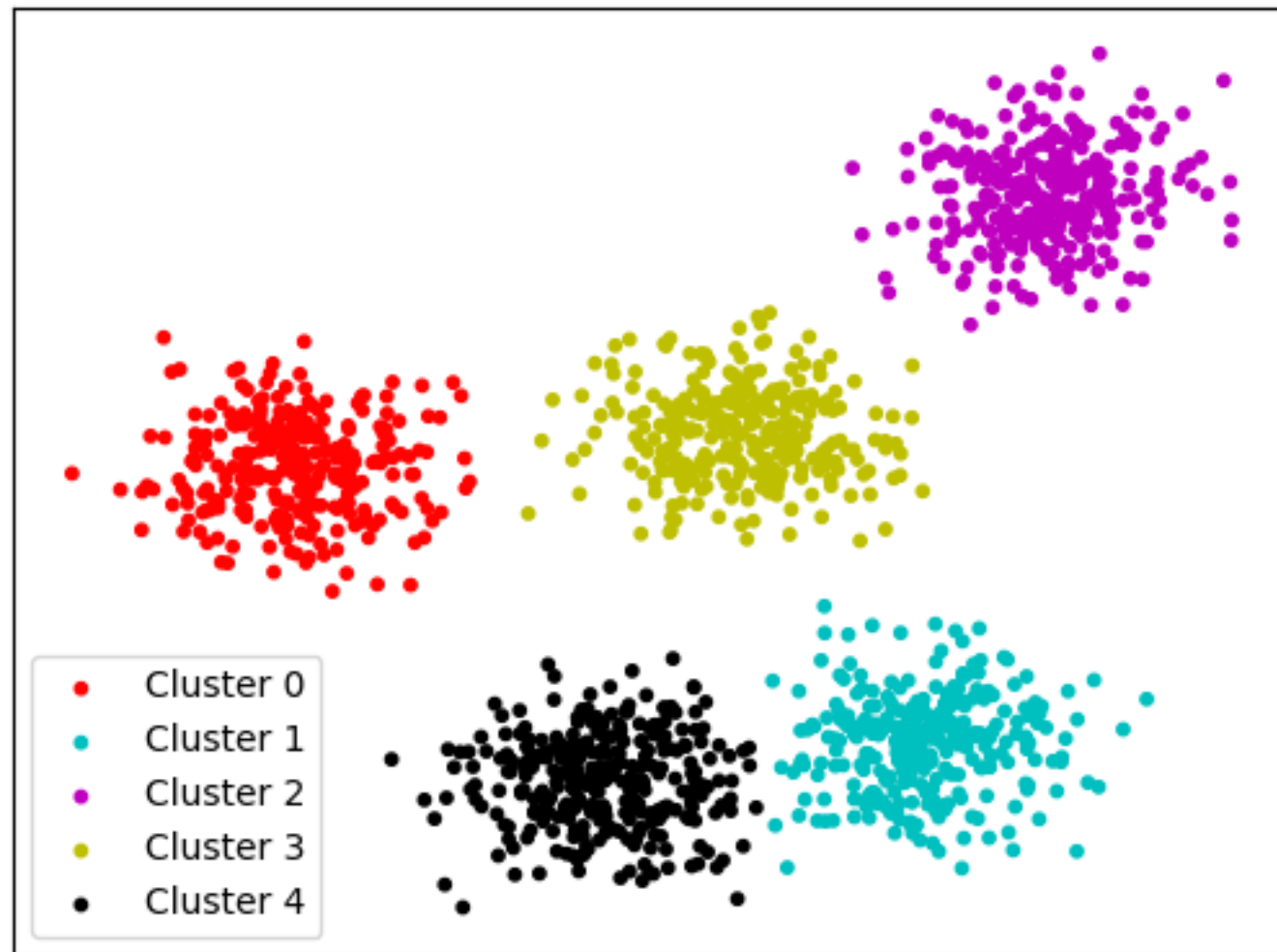Postnet Bar Code representing Delivery Point

# Feature Engineering - Clustering

# Text to Vector

$$w_{i,j} = tf_{i,j} \times \log \left( \frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of $i$ in $j$

$df_i$ = number of documents containing $i$

$N$ = total number of documents

# Example: Greyhound Comments
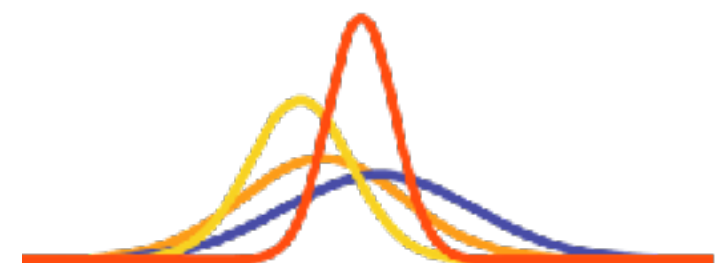
# Thank You Questions?

**Michael Patnik**
**mike@avenuefour.com**

**Avenue 4 Analytics**