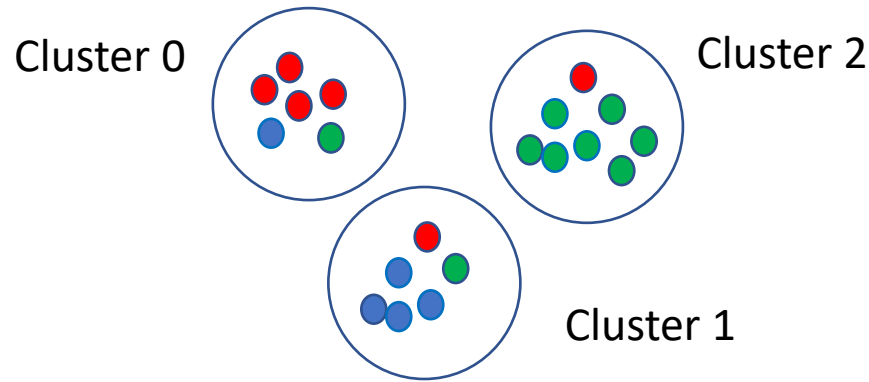


# Unsupervised Machine Learning with Python

# Section 10.1: Clustering Quality Measures

# Clustering Quality Measures

- Consider clustering example where each data point is assigned to a class



How can we quantify the clustering quality in this example?

- 3 actual classes (red, blue, green data points)
- 3 clusters are found, each with more than 1 class
- In perfect world clustering will identify clusters that have exactly 1 class
  - There should be red cluster, blue cluster, and green cluster

# Purity Measure

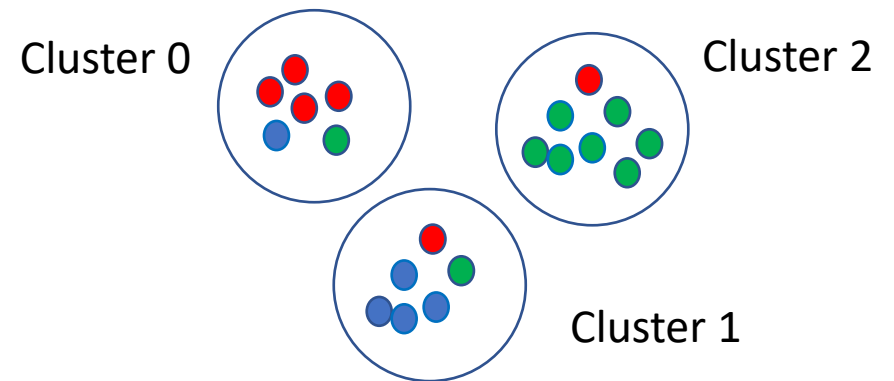
- Purity measures extent to which clusters contain a single class
- M is number of data points, C is set of clusters, D is set of classes
- For each cluster: determine maximum number of data points from any class
- Purity is sum of these maximums divided by total number of data points

$$P = \frac{1}{M} \sum_{c \in C} \max_{d \in D} |d \cap c|$$

- Purity satisfies  $0 < P \leq 1$  (P=1 for “perfect” clustering)

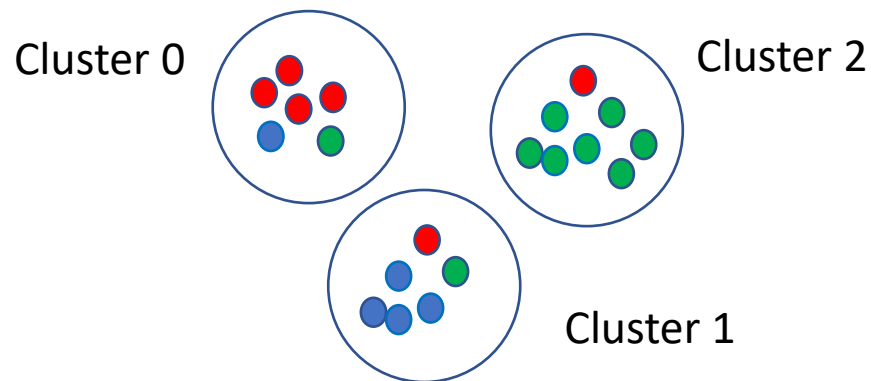
## Example

- 20 data points and 3 clusters
- 3 actual classes: red, blue, green
- Max number from any class:
  - Cluster 0: 4 red, Cluster 1: 4 blue, Cluster 2: 7 green
- $Purity = \frac{1}{20} (4 + 4 + 7) = 0.75$

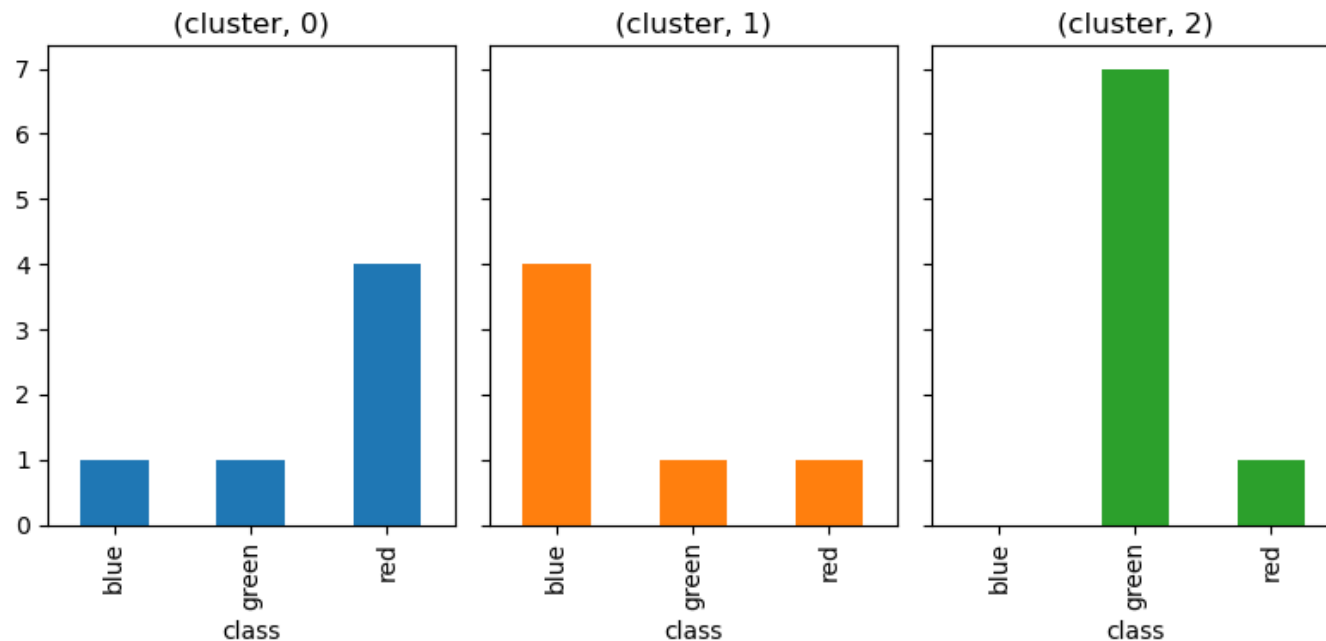


# Bar Chart

- Given data set:



- Can also represent clustering results using bar chart
  - In perfect world, there should non-zero bar for only 1 class for each chart



# Cluster Quality Code Design

Function	Input	Description
purity	cluster_assignment (1d numpy array) class_assignment (1d numpy array)	Computes purity value given cluster and class assignments Return: purity See <a href="#">UnsupervisedML/Examples/Section10/ClusteringQuality.ipynb</a>
plot_cluster_distribution	cluster_assignment (1d numpy array) class_assignment (1d numpy array) figsize (tuple) figrow (integer)	Creates bar charts given cluster and class assignments. figsize and figrow are used to configure the bar charts.  Return: nothing See: <a href="#">UnsupervisedML/Examples/Section10/ClusteringQuality.ipynb</a>

# 10.1 Clustering Quality DEMO

Jupyter Notebook located at:

- UnsupervisedML/Examples/Section10/ClusteringQuality.ipynb

Clustering Quality functions located at:

- UnsupervisedML/Code/Programs

Files to Review	Description
metrics.py	File containing purity and bar plot creation functions

Course Resources at:

- <https://github.com/satishchandrareddy/UnsupervisedML/>
- Stop video if you would like to implement code yourself first

# Unsupervised Machine Learning with Python

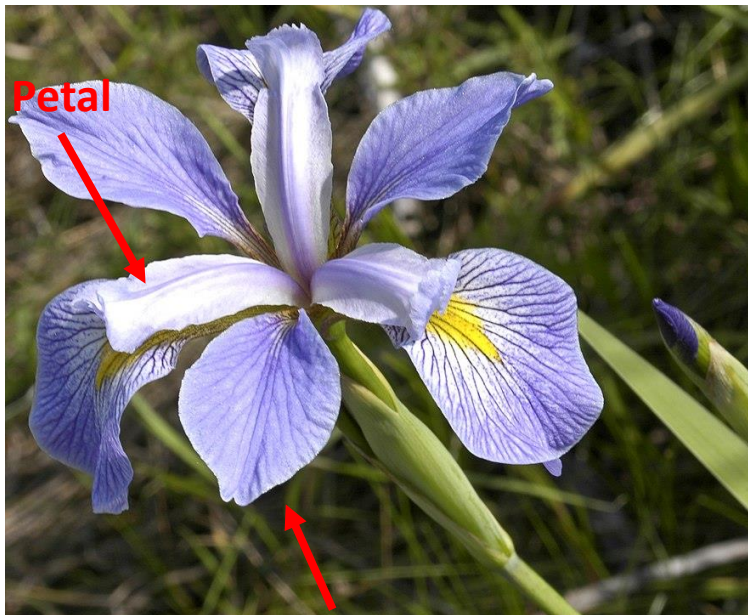


# Section 10.2: Clustering for Iris Flower Dataset

# Iris Flower Dataset

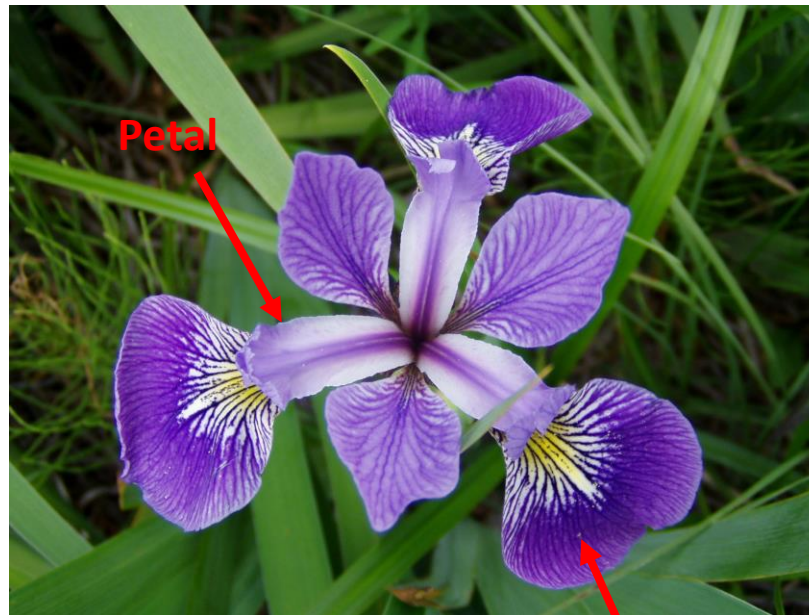
Three types of iris flowers in dataset

Iris Virginica



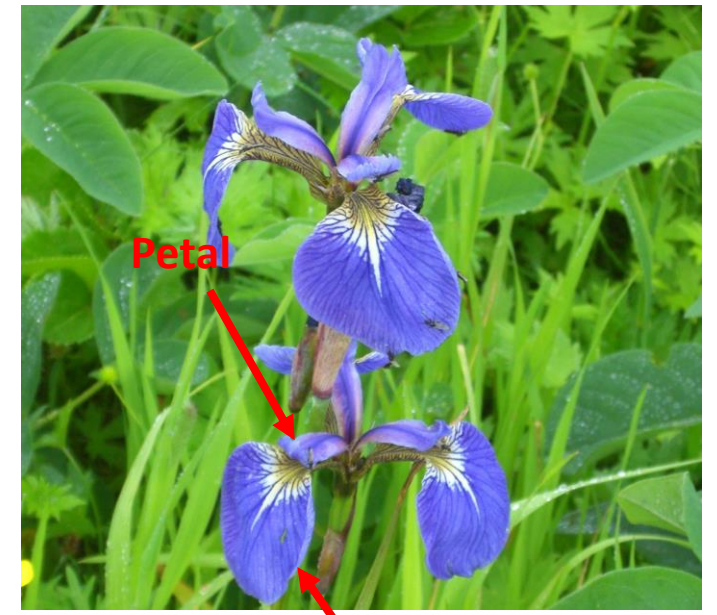
Sepal

Iris Versicolor



Sepal

Iris Setosa



Sepal

See UnsupervisedML\_Resources.pdf file for links  
Images reproduced here under Wikipedia Commons Copyright

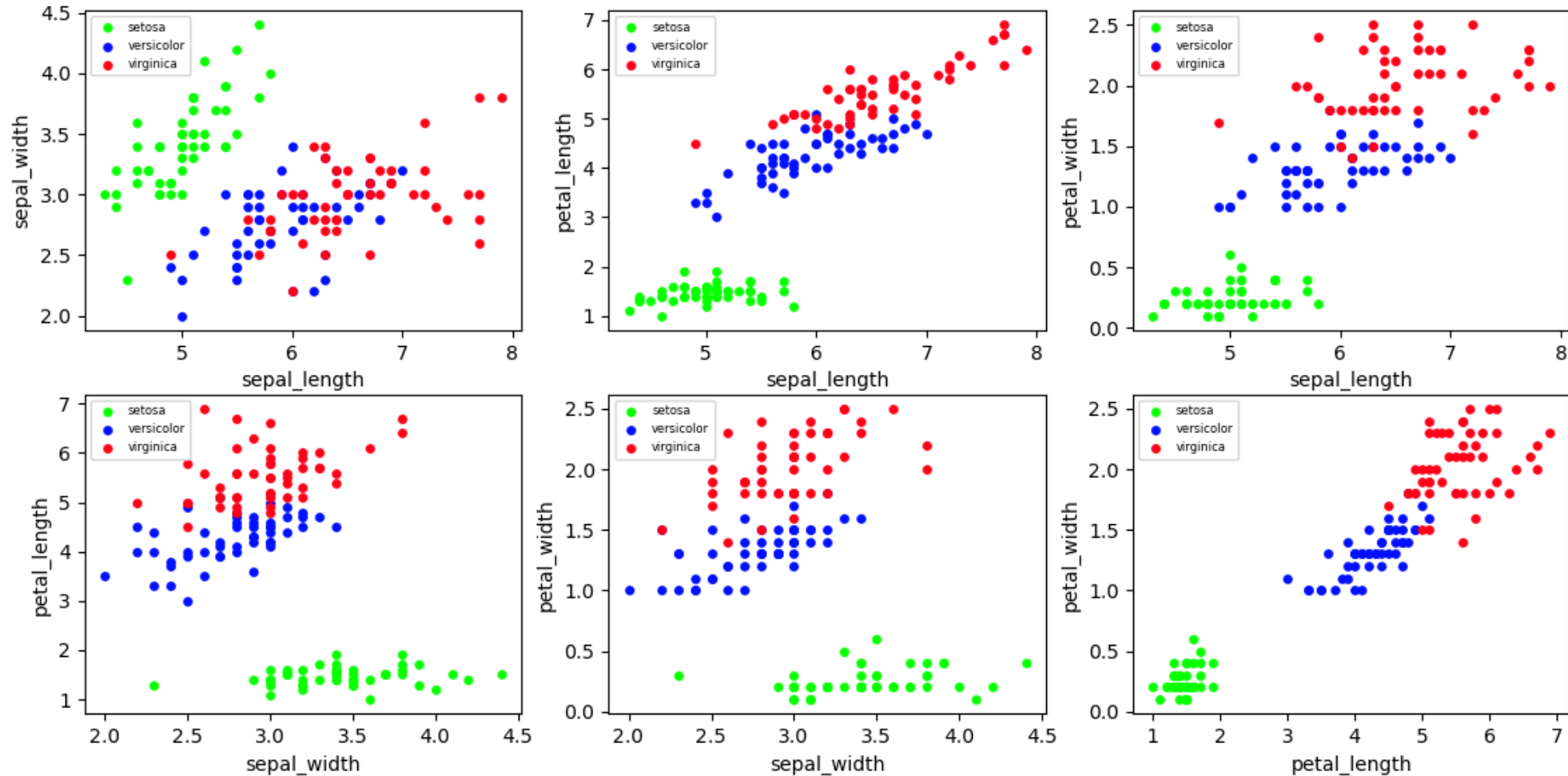
# Iris Dataset

- 50 samples each of 3 types of iris flower species: setosa, virginica, versicolor
- 4 features: sepal\_length, sepal\_width, petal\_length, petal\_width
- Dataset available at UCI Machine Learning Repository  
<https://archive.ics.uci.edu/ml/datasets/iris>
- File: Unsupervised/Clustering/Code/Data\_Iris/Iris.csv

M14				✕		✓		fx			
	A	B	C	D	E	F	G	H			
1		species_id	species	sepal_length	sepal_width	petal_length	petal_width				
2	0	1	setosa	5.1	3.5	1.4	0.2				
3	1	1	setosa	4.9	3	1.4	0.2				
4	2	1	setosa	4.7	3.2	1.3	0.2				
5	3	1	setosa	4.6	3.1	1.5	0.2				
6	4	1	setosa	5	3.6	1.4	0.2				
7	5	1	setosa	5.4	3.9	1.7	0.4				
8	6	1	setosa	4.6	3.4	1.4	0.3				
9	7	1	setosa	5	3.4	1.5	0.2				
10	8	1	setosa	4.4	2.9	1.4	0.2				
11	9	1	setosa	4.9	3.1	1.5	0.1				
12	10	1	setosa	5.4	3.7	1.5	0.2				
13	11	1	setosa	4.8	3.4	1.6	0.2				
14	12	1	setosa	4.8	3	1.4	0.1				
15	13	1	setosa	4.3	3	1.1	0.1				
16	14	1	setosa	5.8	4	1.2	0.2				
17	15	1	setosa	5.7	4.4	1.5	0.4				
18	16	1	setosa	5.4	3.9	1.3	0.4				
19	17	1	setosa	5.1	3.5	1.4	0.2				

# Iris Dataset

Iris Data



# Examples in this Section

## Example 1

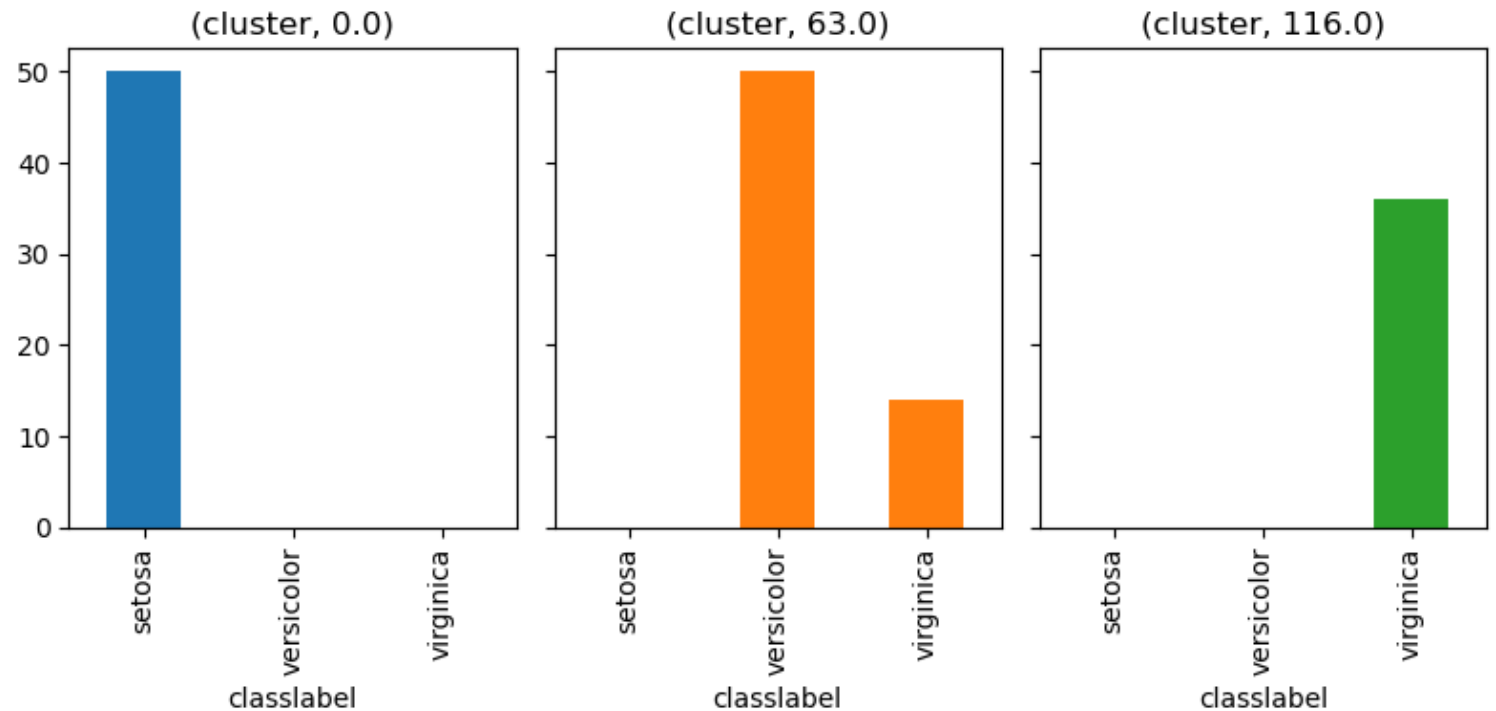
- Hierarchical Clustering for Iris dataset

## Example 2:

- Hierarchical Clustering for Iris dataset after using PCA to reduce dataset to 2 dimensions

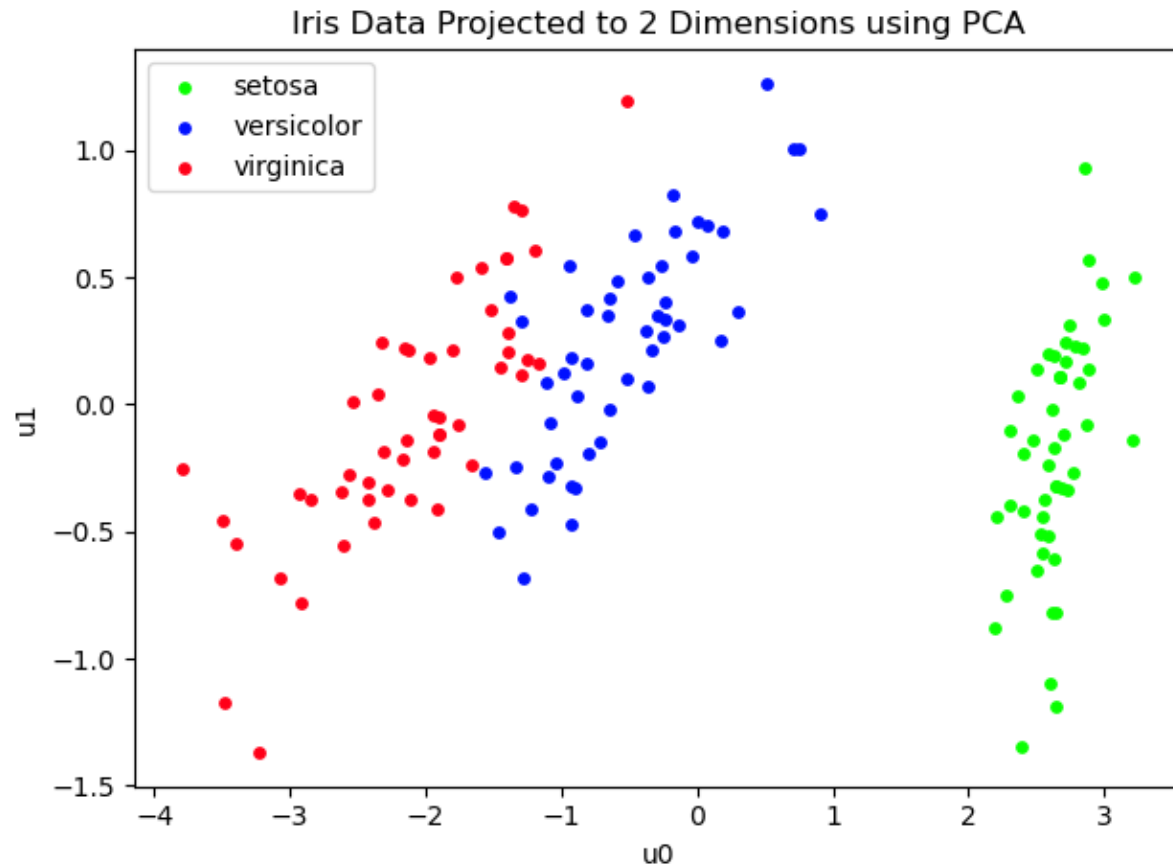
# Example 1: Clustering of Iris Dataset

- Dataset: Feature matrix  $X$  (4 dimensions x 150 data points)
- Algorithm: Hierarchical Clustering (stop at 3 clusters)
- Metrics:
  - Purity: 0.907
  - Davies-Bouldin: 0.66



# PCA for Iris Dataset

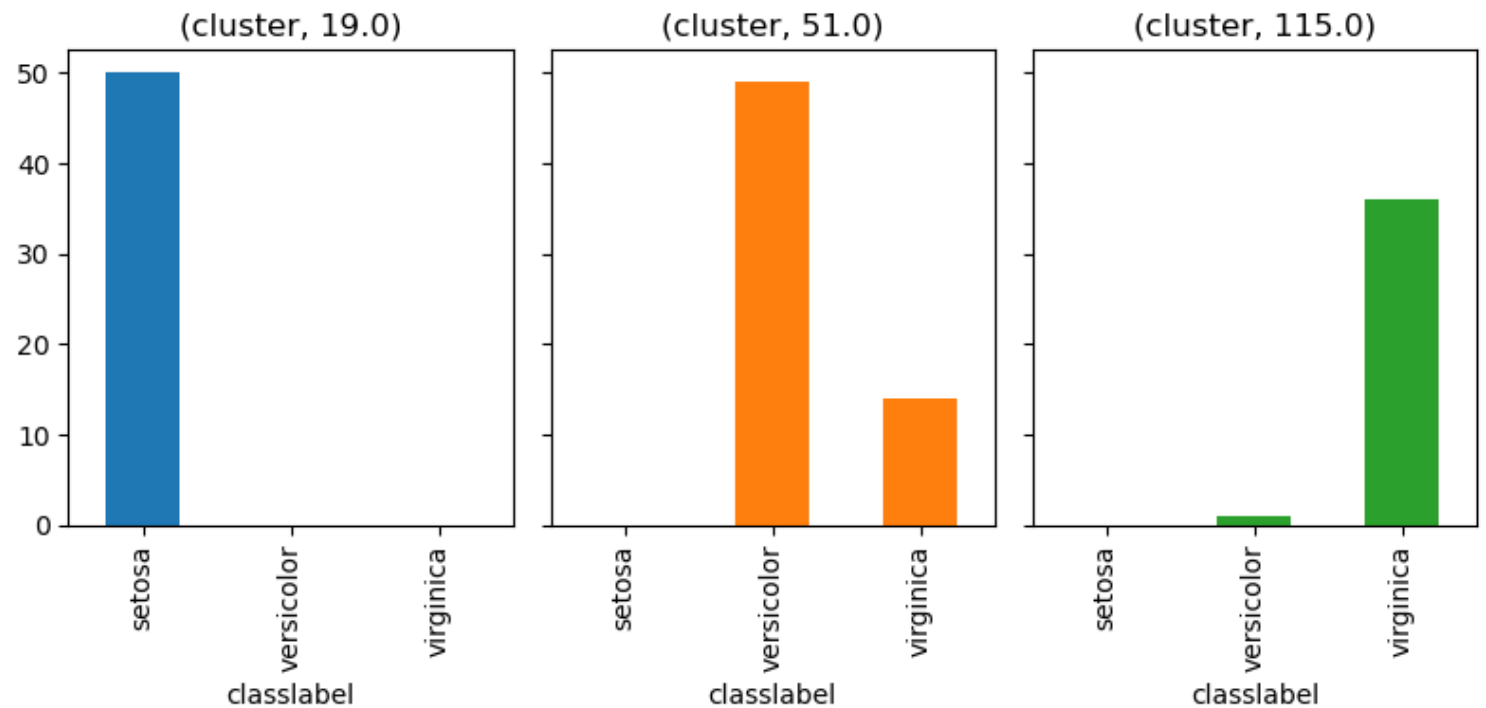
- Project data from 4 dimensions to 2 dimension using PCA
- Variance capture is 97.8%



- New basis vectors/features  $u_0$  and  $u_1$  do not correspond to actual measurable quantities, such as sepal width or length or petal width or length

# Example 2: Clustering for Iris Dataset using PCA

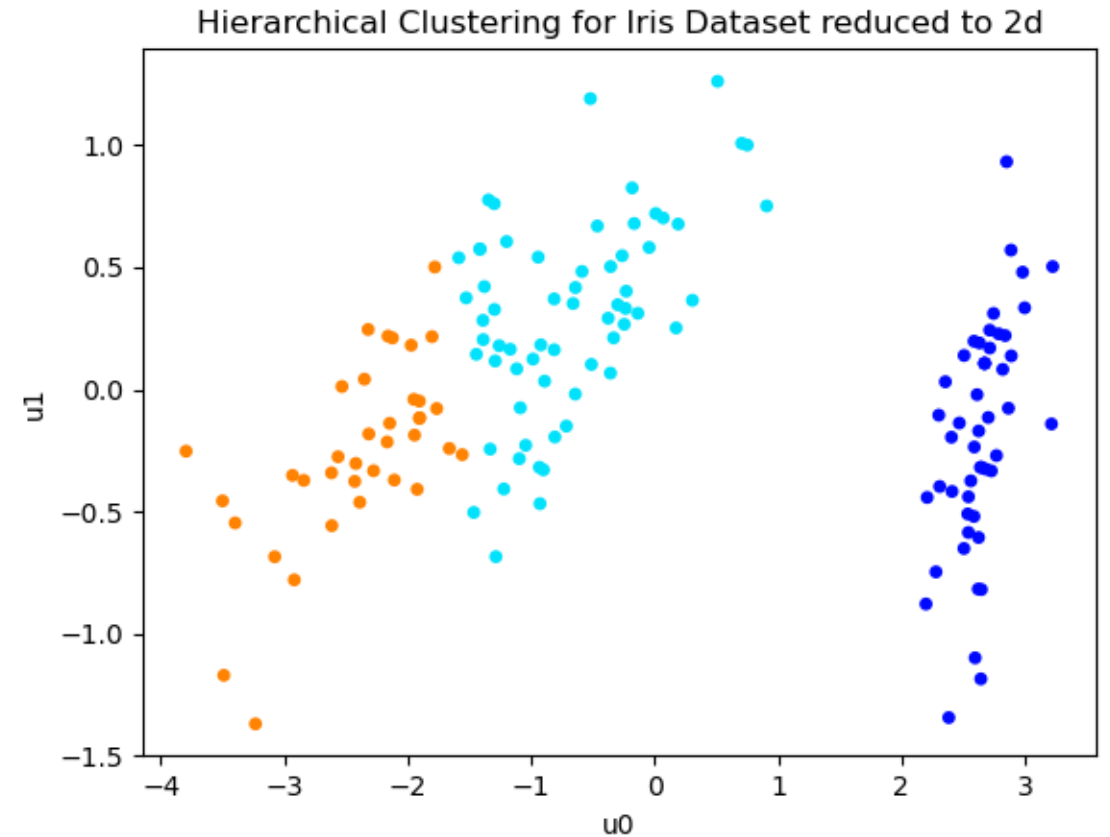
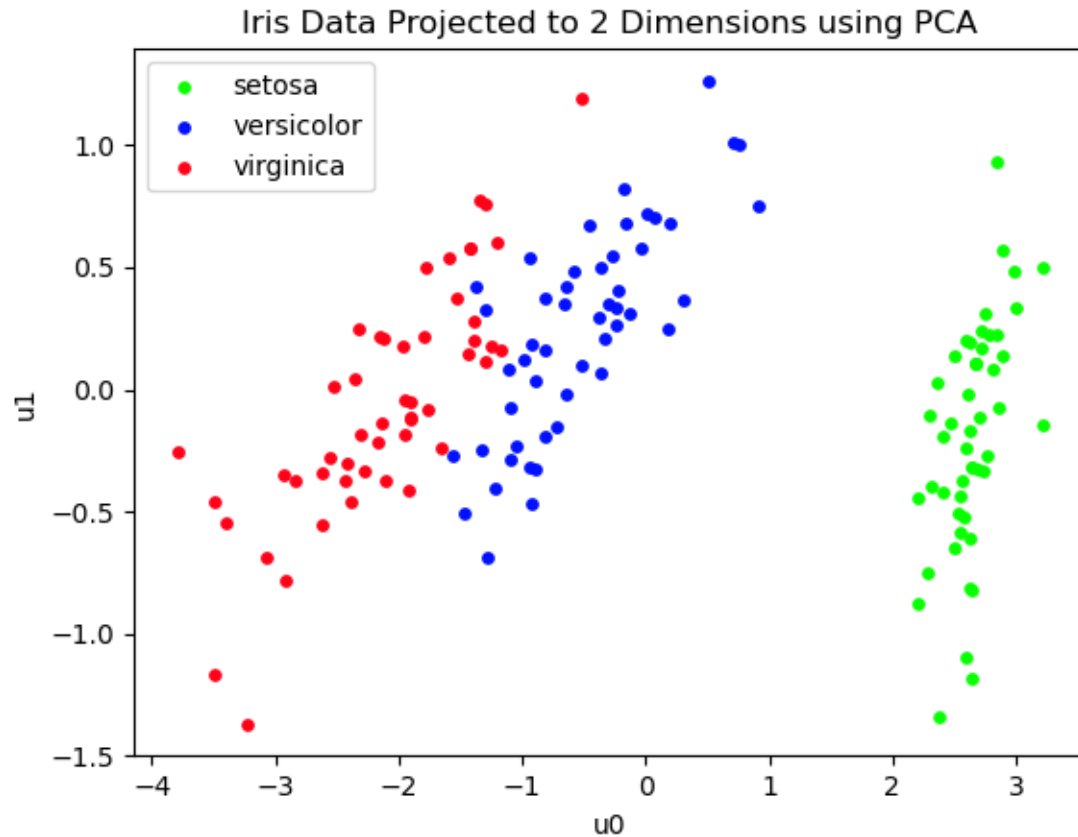
- Dataset: reduced dimension feature matrix R (2 x 150)
- Algorithm: Hierarchical Clustering (stop at 3 clusters)
- Metrics:
  - Purity: 0.90
  - Davies-Bouldin: 0.66





# Clustering for Iris Dataset using PCA

- Comparison of Class and Clustering Results



# iris class Code Design

method	Input	Description
<code>__init__</code>		Constructor for iris class – saves data directory Return: nothing
<code>load</code>		Loads all 150 samples and corresponding class labels from iris dataset  Return: X (2d numpy array), class_label (1d numpy array) See <a href="#">UnsupervisedML/Examples/Section02/Pandas.ipynb</a>
<code>plot</code>		Creates scatter plots showing classes as a function of all possible 2 variable combinations of sepal width, sepal length, petal width & petal length  Return: nothing See <a href="#">UnsupervisedML/Examples/Section02/MatplotlibAdvanced.ipynb</a>

# Iris Clustering Code Walkthrough

Code and data located at:

- UnsupervisedML/Code/Programs
- UnsupervisedML/Code/Data\_Iris

Files to Review	Description
iris.csv	Iris dataset
data_iris.py	Class for loading and processing iris data
plot_data.py	Functions for creating basic scatter plots
casestudy_iris.py	Driver for iris clustering
casestudy_iris_pca.py	Driver for iris clustering using pca to reduce dimension

Course Resources at:

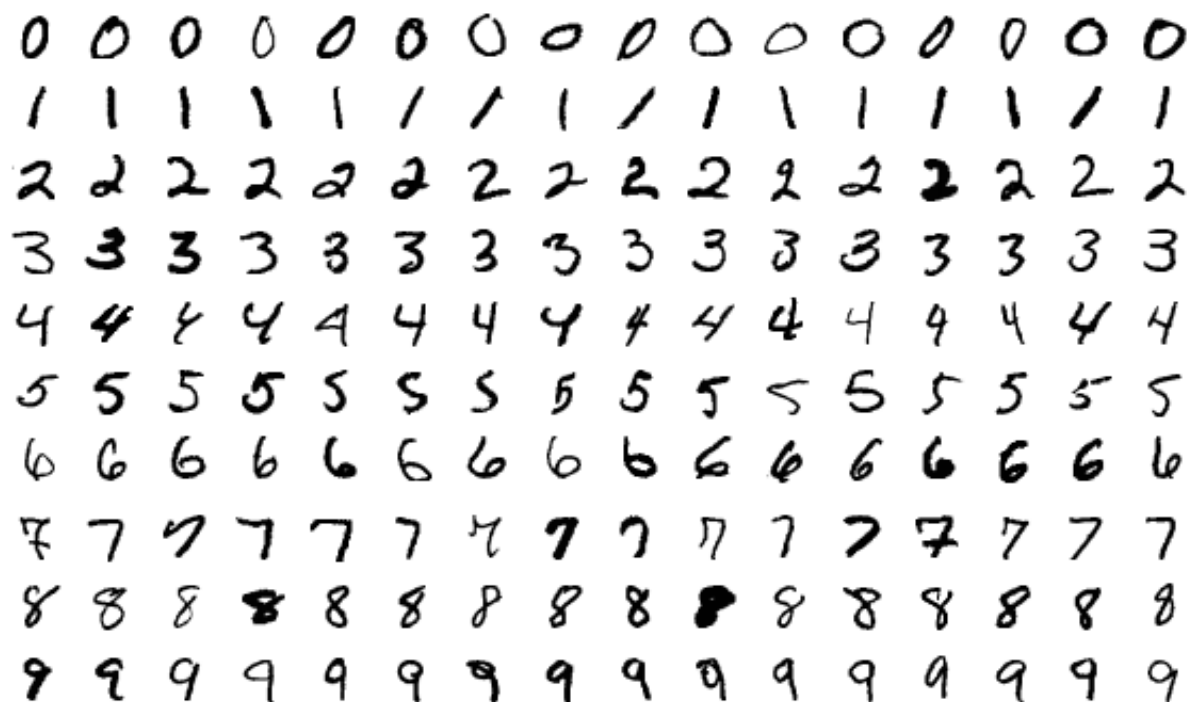
- <https://github.com/satishchandrareddy/UnsupervisedML/>
- Stop video if you would like to implement code yourself first

# Unsupervised Machine Learning with Python

# Section 10.3: Clustering for MNIST Digits Dataset

# MNIST Digits Dataset

- Thousands of handwritten digit images with 28x28 resolution
- Data Source: <http://yann.lecun.com/exdb/mnist/>
- Used extensively for testing machine learning algorithms



Collage of 160 individual digit images

By Josef Steppan - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=64810040>

# Examples in this Section

Clustering problem:

- Employ clustering algorithm/PCA to group images in MNIST Dataset
- See how well algorithm creates clusters with the same digits

Example 1

- K Means Clustering for MNIST Dataset

Example 2:

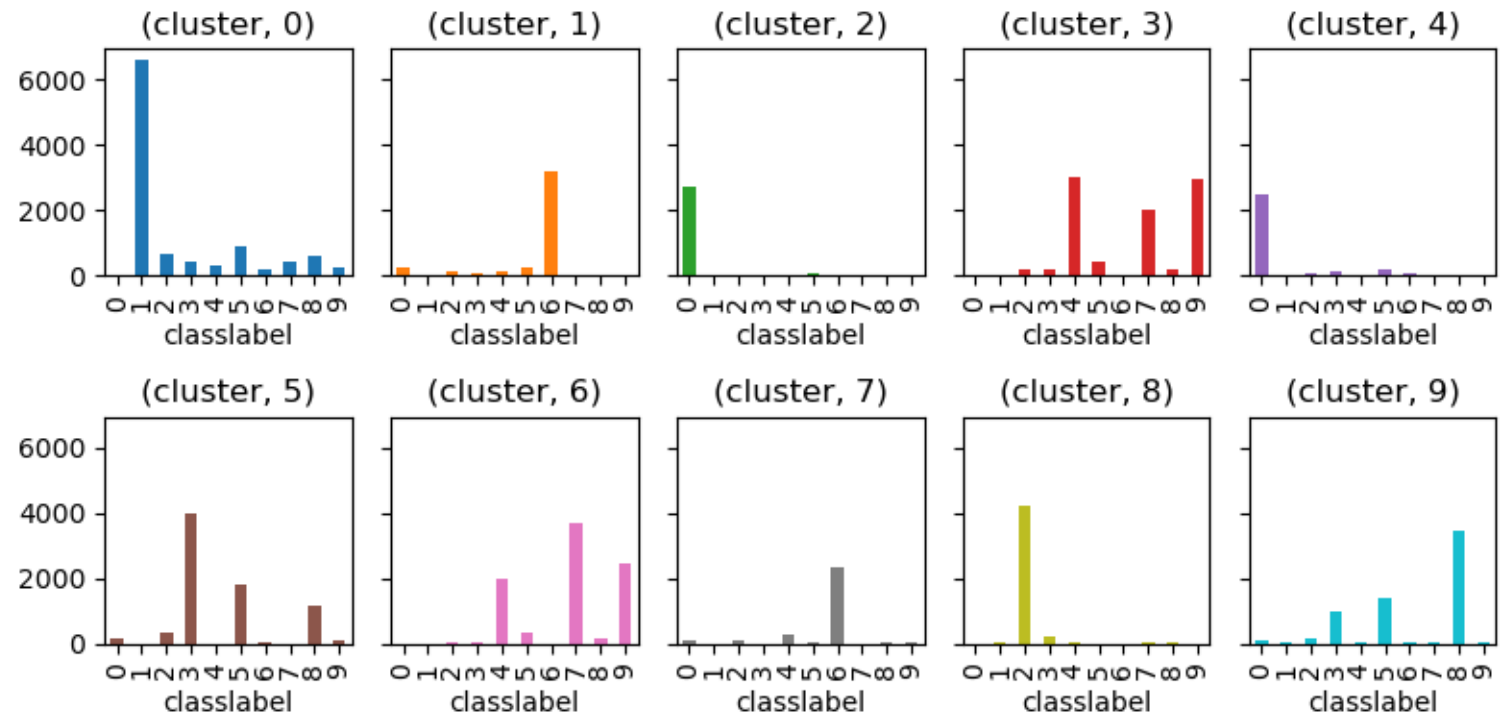
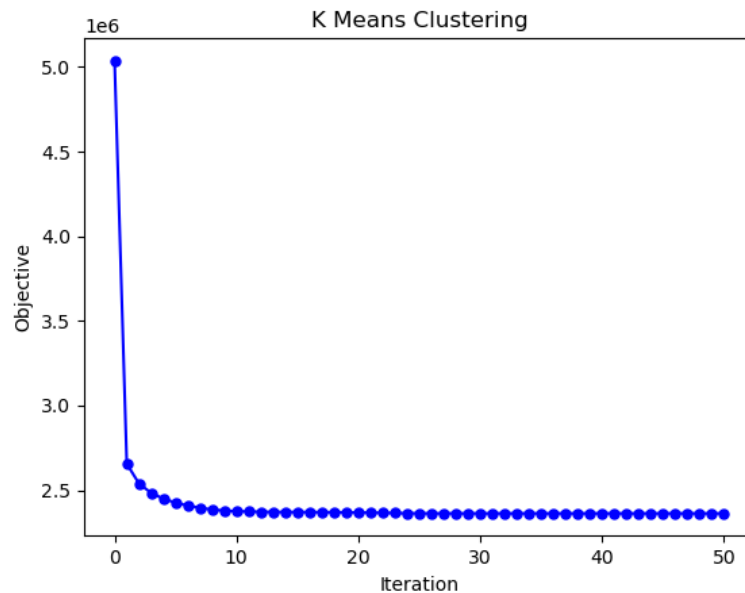
- K Means Clustering for MNIST Dataset after using PCA to capture 90% of variance (reduces number of dimensions to 87)

Example 3:

- Gaussian MM Clustering for MNIST Dataset after using PCA to capture 90% of variance (reduces number of dimensions to 87)

# Example 1: K Means Clustering

- Dataset: Feature matrix X (784 dimensions x 60000 images)
- Algorithm: K Means with 10 clusters, kmeans++ for initialization, 100 iterations maximum, tolerance of  $10^{-4}$
- Metrics:
  - Purity: 0.596
  - Davies-Bouldin: 2.82
  - Clustering Time: 196 seconds

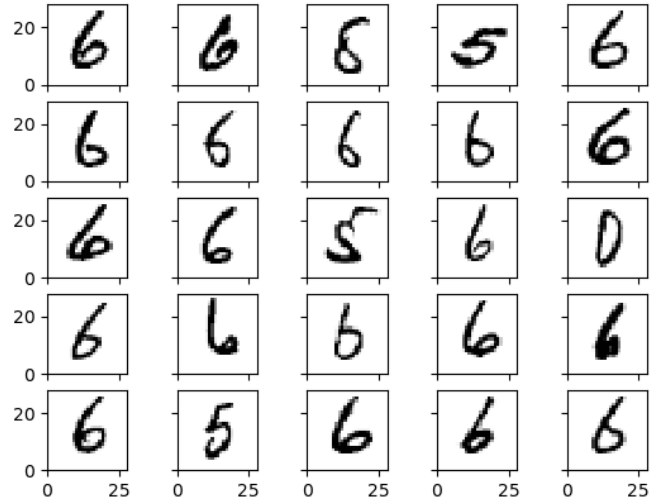




# Example 1: K Means Clustering

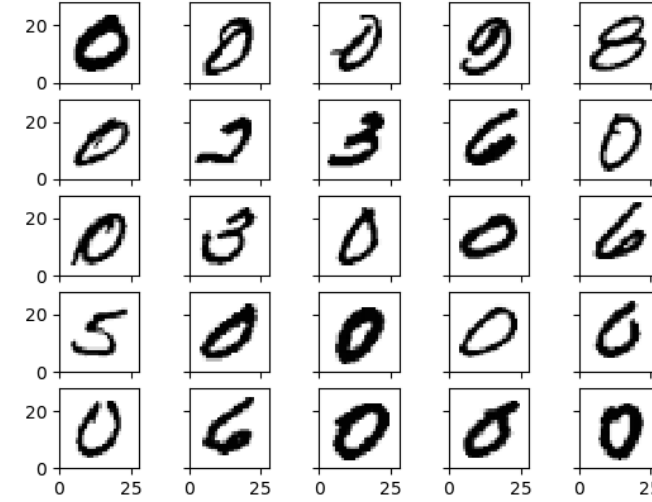
Images of Sample MNIST Digits

Cluster 1



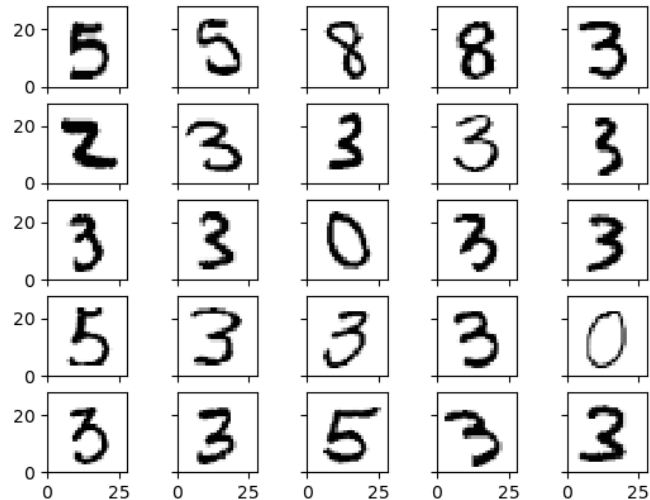
Images of Sample MNIST Digits

Cluster 4



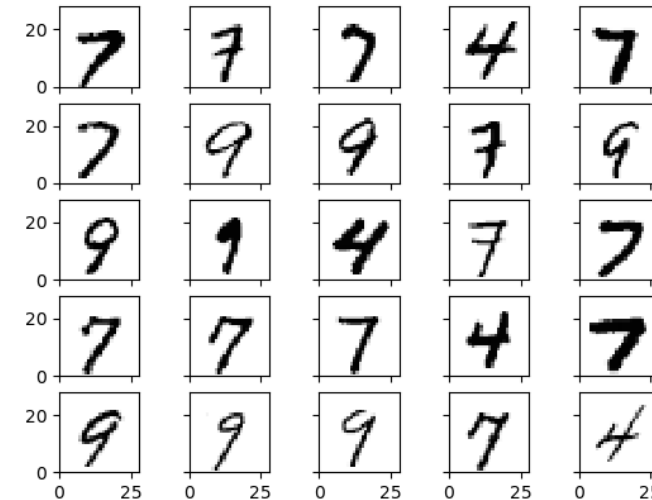
Images of Sample MNIST Digits

Cluster 5



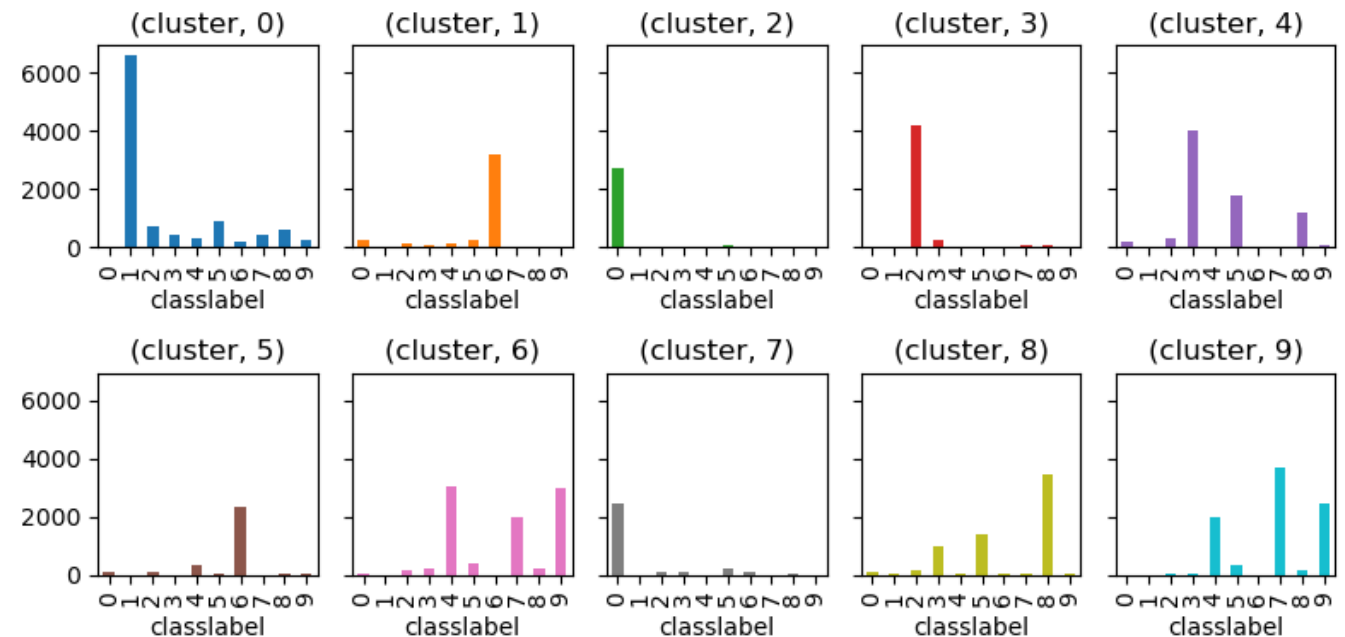
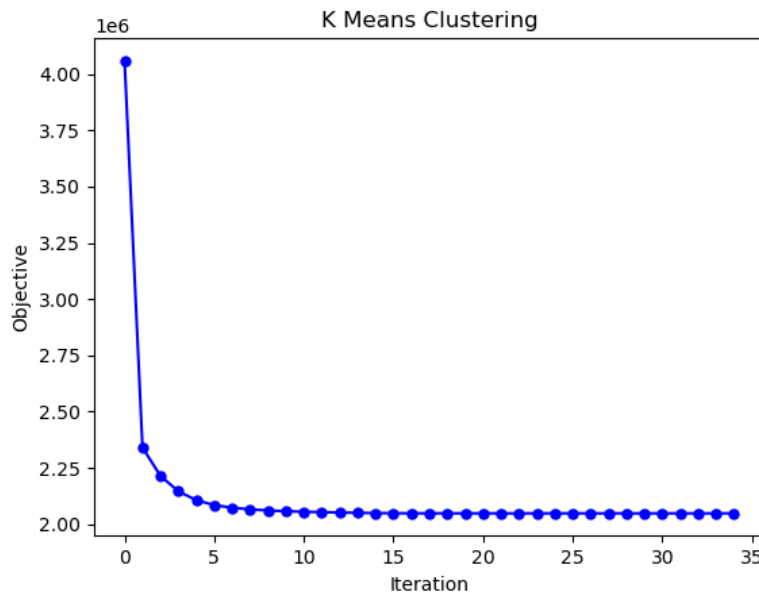
Images of Sample MNIST Digits

Cluster 6



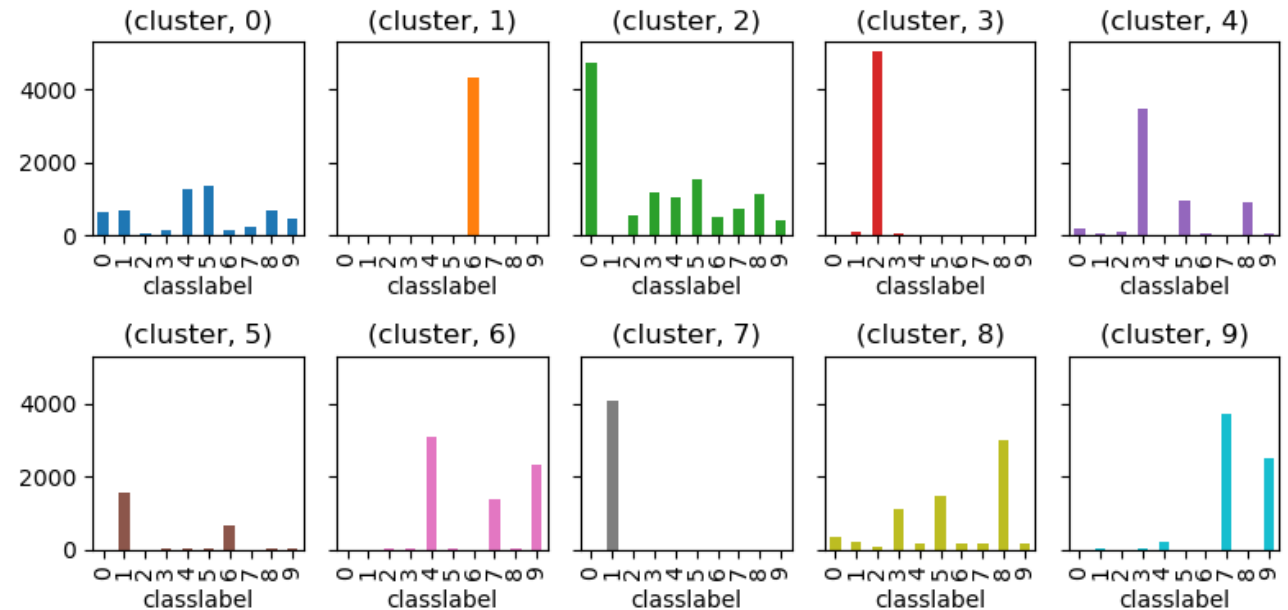
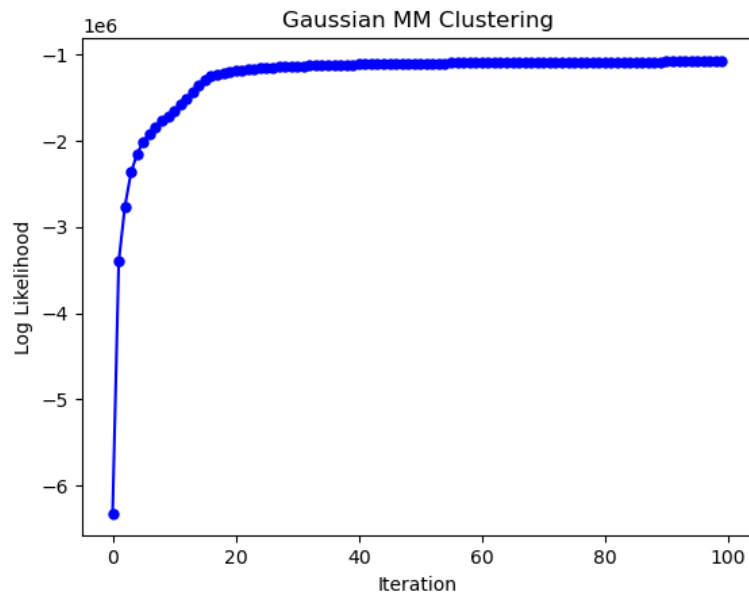
# Example 2: K Means Clustering with PCA

- Dataset: apply PCA with 90% variance capture to 60000 images resulting in feature matrix R (87 dimensions x 60000 images)
- Algorithm: K Means with 10 clusters, kmeans++ for initialization, 100 iterations maximum,  $10^{-4}$  tolerance
- Metrics:
  - Purity: 0.596
  - Davies-Bouldin: 2.82
  - PCA Time: 7.6 seconds + Clustering Time: 16.2 seconds



## Example 3: GaussianMM Clustering with PCA

- Dataset: apply PCA with 90% variance capture to 60000 images resulting in feature matrix R (87 dimensions x 60000 images)
- Algorithm: GaussianMM with 10 clusters, kmeans++ for initialization, 100 iterations maximum,  $10^{-4}$  tolerance
- Metrics:
  - Purity: 0.574
  - Davies-Bouldin: 3.19
  - PCA Time: 6.6 seconds + Cluster Time: 172 seconds



# Comments

- K Means Clustering

- Achieved 60% Purity result for clustering for 60000 image MNIST dataset
- Clustering results similar with/without PCA as measured by Purity and Davies-Bouldin values
- Using PCA can significantly reduce clustering time
- I have found many more iterations are required for convergence using “random” compared to “kmeans++” initialization

- Gaussian MM Clustering

- Don't get convergence after 100 iterations even after applying PCA to reduce dimensions to 87
- Approach is slow for large numbers of dimensions
- I am finding that method is not stable for other values of variance capture – numerical issues because determinant of covariance matrix is close to 0
- In exercises you will investigate using spherical Gaussian MM approach

# MNIST Clustering Code Walkthrough

Code and data located at:

- UnsupervisedML/Code/Programs
- UnsupervisedML/Code/Data\_MNIST

Files to Review	Description
MNIST_train_set1_30K.csv MNIST_train_set2_30K.csv MNIST_valid_10K.csv	MNIST train and valid datasets
data_mnist.py	Code for loading and processing MNIST data
casestudy_mnist.py	Driver for MNIST clustering

Course Resources at:

- <https://github.com/satishchandrareddy/UnsupervisedML/>
- Stop video if you would like to implement code yourself first

# Unsupervised Machine Learning with Python

# Section 10.4: Clustering for Text Documents

# BBC Text Dataset

- 2225 BBC articles in 5 categories: sport, business, tech, entertainment, politics
- Dataset: <https://www.kaggle.com/yufengdev/bbc-fulltext-and-category>
- File: UnsupervisedML/Code/Data\_Text/bbc-text.csv
- Use Tfidf vectorizer in sklearn
- 12915 words in dictionary -> 12915 x 2225 feature matrix

	A	B	C	D	E	F	G	H	I	J
1	category	text								
2	tech	tv future in the hands of viewers with home theatre systems plasma high-definition tvs and di								
3	business	worldcom boss left books alone former worldcom boss bernie ebbers who is accused of over								
4	sport	tigers wary of farrell gamble leicester say they will not be rushed into making a bid for andy fa								
5	sport	yeading face newcastle in fa cup premiership side newcastle united face a trip to ryman premie								
6	entertainm	ocean s twelve raids box office ocean s twelve the crime caper sequel starring george clooney								
7	politics	howard hits back at mongrel jibe michael howard has said a claim by peter hain that the tory le								
8	politics	blair prepares to name poll date tony blair is likely to name 5 may as election day when parlian								
9	sport	henman hopes ended in dubai third seed tim henman slumped to a straight sets defeat in his ra								
10	sport	wilkinson fit to face edinburgh england captain jonny wilkinson will make his long-awaited retu								
11	entertainm	last star wars not for children the sixth and final star wars movie may not be suitable for youn								
12	entertainm	berlin cheers for anti-nazi film a german movie about an anti-nazi resistance heroine has drawi								
13	business	virgin blue shares plummet 20% shares in australian budget airline virgin blue plunged 20% after								
14	business	crude oil prices back above \$50 cold weather across parts of the united states and much of eur								



# Examples in this Section

Clustering problem:

- Employ clustering algorithm/PCA to group articles in BBCText dataset
- How well can algorithm create clusters of articles in the same category?

Example 1

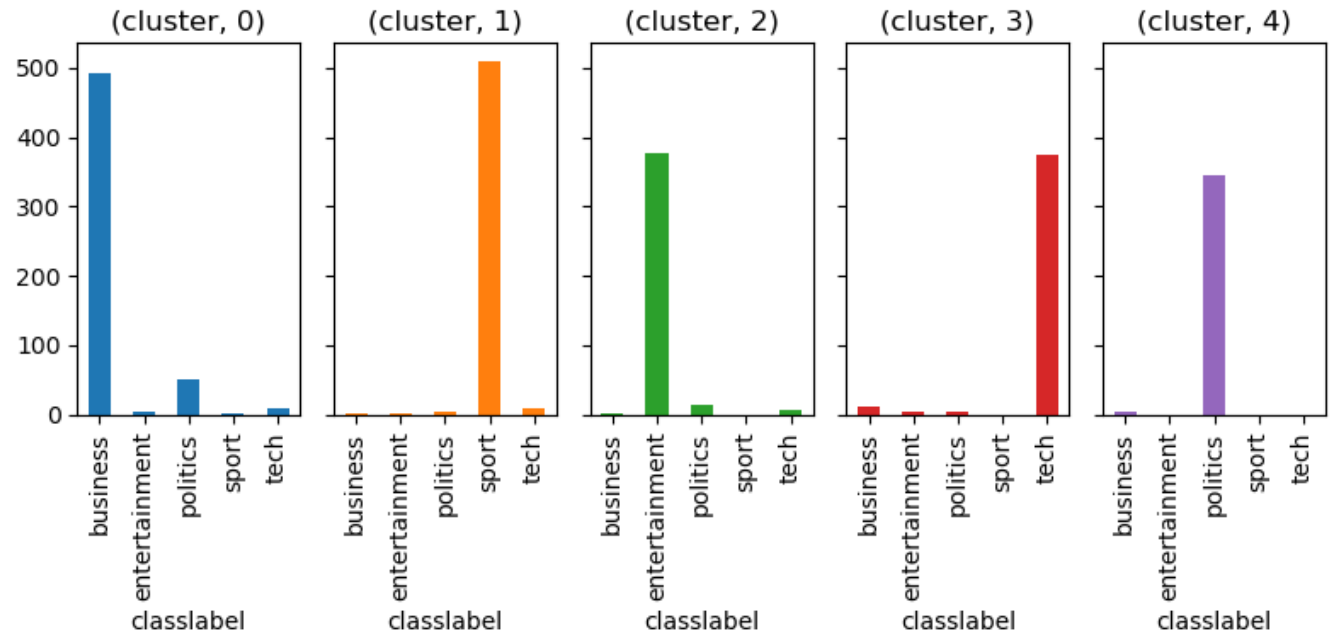
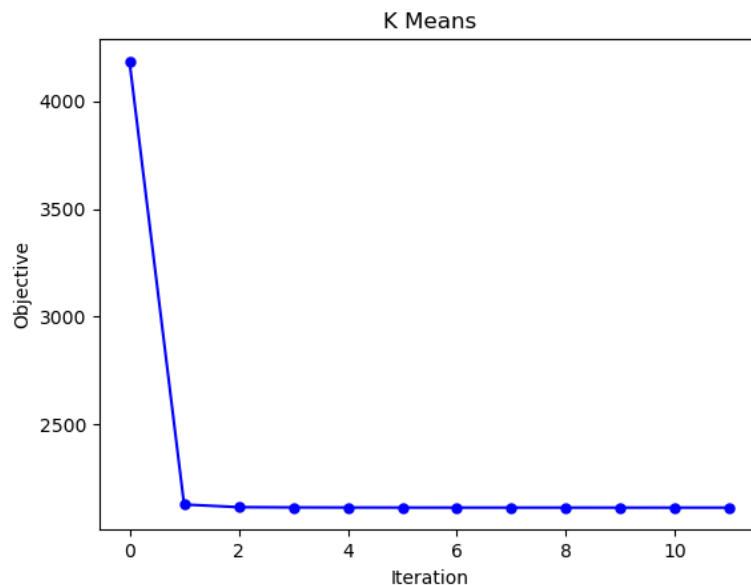
- K Means Clustering for BBCText dataset

Example 2:

- K Means Clustering for BBCText dataset after using PCA to reduce dimensions and still capture 100% of variance

# Example 1: K Means Clustering

- Dataset: Feature matrix X (12915 dimensions x 2225 data points)
- Algorithm: K Means with 5 clusters, random initialization, 50 iterations maximum, tolerance of  $10^{-4}$
- Metrics:
  - Purity: 0.943
  - Davies-Bouldin: 8.06
  - Clustering Time: 13.8 seconds



# Example 1: Wordclouds for Clusters

### Cluster 0



## Cluster 1



### Cluster 2



### Cluster 3

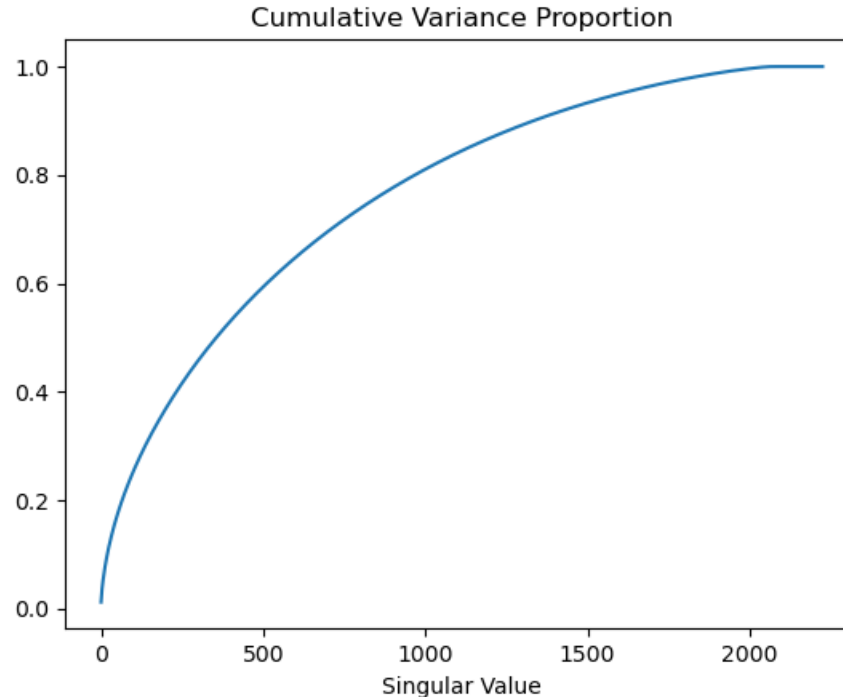


## Cluster 4



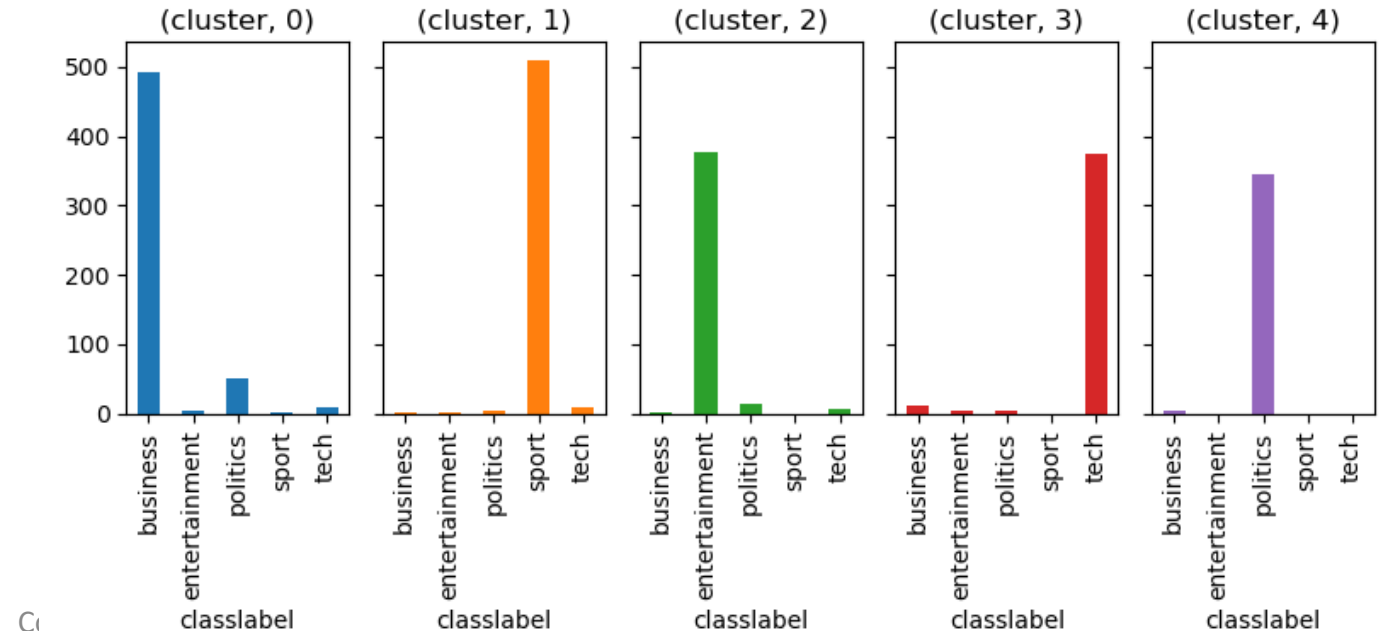
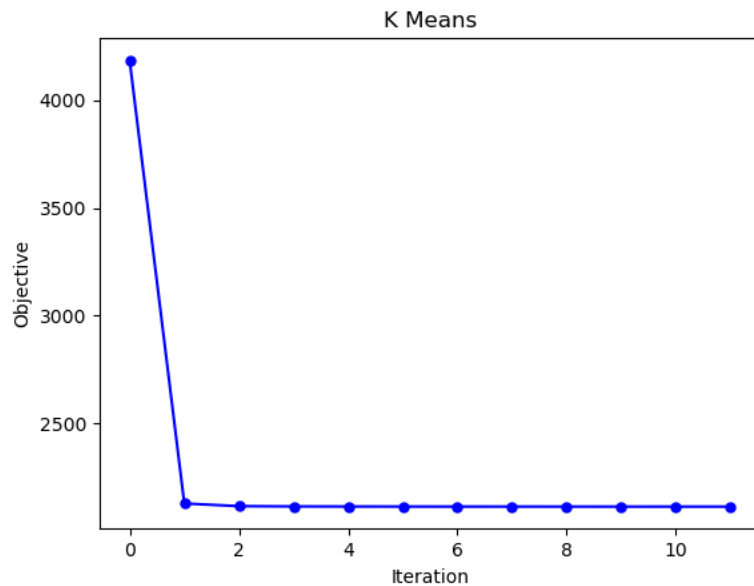
# PCA for BBC Text Dataset

- Perform PCA for BBC Text Dataset
- Since number of dimensions (12915) > number of data points (2225), can reduce number of dimensions and still retain 100% of variance
- Actually, since some singular values are 0, can retain 100% of variance with 2116 dimensions



# Example 2: K Means Clustering with PCA

- Dataset: use PCA to reduce dimension and still capture 100% of variance – results in feature matrix R ( 2116 dimensions x 2225 data points)
- Algorithm: K Means with 5 clusters, random initialization, 50 iterations maximum, tolerance of  $10^{-4}$
- Metrics: (clusters are exactly the same as in Example 1)
  - Purity: 0.943
  - Davies-Bouldin: 8.06
  - PCA Time: 13.6 seconds + Clustering Time: 2.7 seconds



# Comments

- K Means Clustering algorithm is able to achieve greater than 94% purity measure for grouping articles in BBC Text dataset
- Can use PCA to reduce dimension and maintain 100% variance capture
  - Dimension reduced from 12915 to 2116
  - Clustering results exactly the same as in case of no dimension reduction
  - Clustering time reduced by a factor of 5

# bbctext class Code Design

method	Input	Description
<code>__init__</code>		Constructor for bbctext class – saves directory and TFIDF vectorizer Return: nothing
<code>load</code>	<code>nsample (integer)</code>	Loads bbc text dataset for specified number of samples and applies TFIDF vectorization to create feature matrix  Return: X (2d numpy array), class_label (1d numpy array) -UnsupervisedML/Examples/Section02/Pandas.ipynb -UnsupervisedML/Examples/Section03/SklearnText.ipynb
<code>create_wordcloud</code>	<code>X_tfidf (2d numpy array)</code> <code>cluster_assignment (1d numpy array)</code> <code>ncluster (integer)</code> <code>nword (integer)</code>	Creates wordcloud plot for specified X_tfidf matrix, cluster assignments, and number of words  Return: nothing -UnsupervisedML/Examples/Section03/SklearnText.ipynb

# Text Clustering Code Walkthrough

Code and data located at:

- UnsupervisedML/Code/Programs
- UnsupervisedML/Code/Data\_BBCText

Files to Review	Description
data_bbctext.py	Code for loading and processing BBC text data
casestudy_bbctext.py	Driver for bbc text cluster
bbc-text.csv	BBC text data file

Course Resources at:

- <https://github.com/satishchandrareddy/UnsupervisedML/>
- Stop video if you would like to implement code yourself first