

# LAB #6: CLASSIFICATION

CS 109A, STAT 121A, AC 209A: Data Science

---

Fall 2016

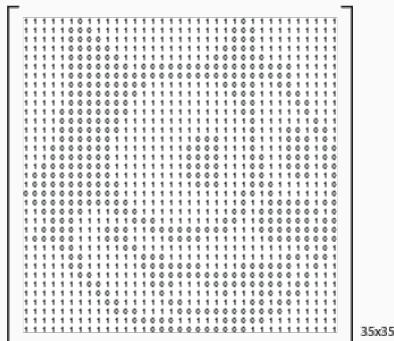
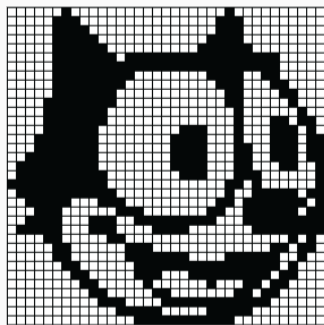
Harvard University

## IMAGE PROCESSING

---

## REPRESENTING IMAGES AS MATRICES AND ARRAYS

An image can be represented digitally as a grid of illuminated pixels (small area of illumination on display device). We can represent this as a matrix of numbers, each number encoding the intensity of the corresponding pixel.



# REPRESENTING IMAGES AS MATRICES AND ARRAYS

We can flatten the image matrix into an image vector:

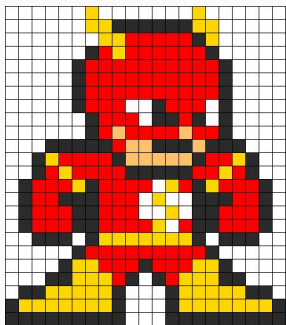


Image as matrix of pixels

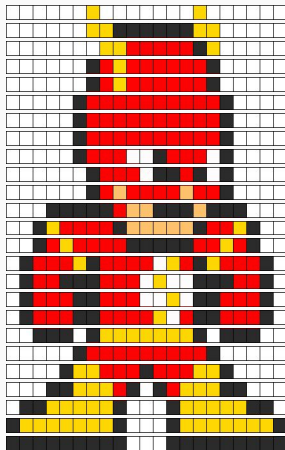


Image matrix split along rows

## REPRESENTING IMAGES AS MATRICES AND ARRAYS

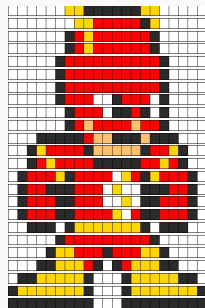


Image as matrix of pixels



Image as single vector of all the rows concatenated

## REPRESENTING IMAGES AS MATRICES AND ARRAYS

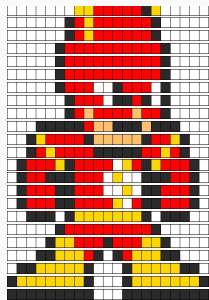


Image as matrix of pixels



Image as single vector of all the rows concatenated

## REPRESENTING IMAGES AS MATRICES AND ARRAYS

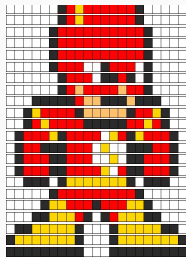


Image as matrix of pixels

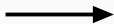
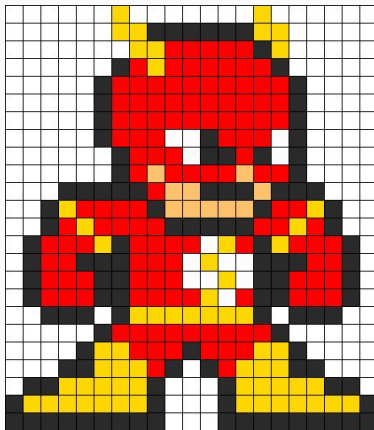


Image as single vector of all the rows concatenated

## REPRESENTING IMAGES AS MATRICES AND ARRAYS



This image, when flattened, is represented as a `numpy` array of shape `(441, )`.

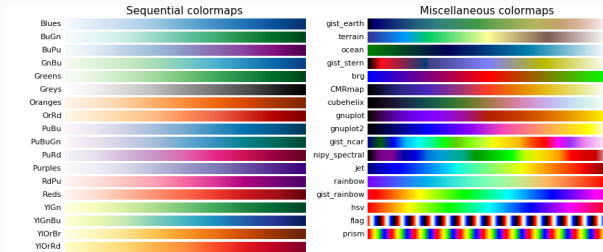


# DISPLAYING AN IMAGE ARRAY

Given an image array, we can reshape it into an image matrix; and then display the image, by mapping each number to an intensity on a colormap.

```
plt.imshow(flash_vector.reshape(21, 21), cmap=color_map)
```

where **color\_map** is an appropriately chosen **matplotlib** color map:



## DIMENSION REDUCTION

---

Now that we know how to convert digital images into numpy arrays, we can feed them as input into any of our fancy statistical models!

In particular, we can perform classification on digital images (automatically classify images as “containing you” or “containing your friends” on Facebook for example).

## THE PROBLEM OF IMAGE VECTORS

Today, you'll be working with a dataset of digital images of hand-written digits. Each data point is an 8x8 gray-scale image of a single hand-written digit, flattened into a 64-length vector.

You will perform a task that major tech companies have been trying to perfect for ages: correctly **classify** handwritten symbols.

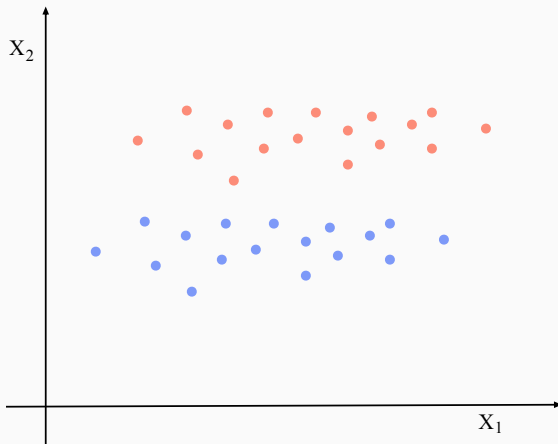
**Question:** what's a potential difficulty with working with images as arrays?

**Dimension reduction** is the task of transforming data with a large number of features (or predictors) into a data set with much fewer number of features.

Dimension reduction can be done in many ways. The most naïve approach would be to simply select a couple of features (or predictors) and chuck the rest - **variable selection**.

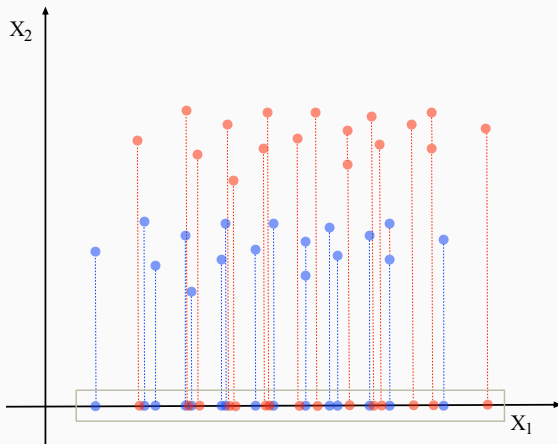
## DIMENSION REDUCTION VIA PROJECTION

A data set with two classes and two predictors,  $X_1$  and  $X_2$ .



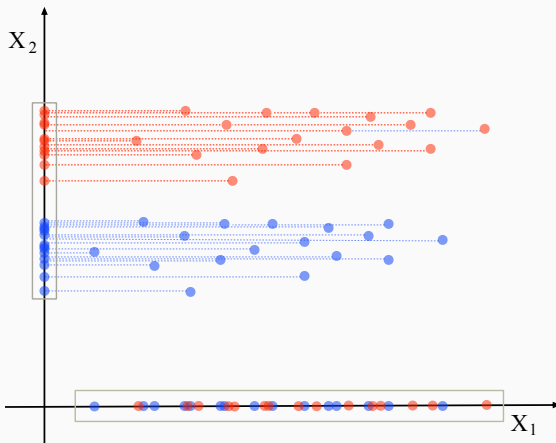
## DIMENSION REDUCTION VIA PROJECTION

We chuck  $X_2$ , and only keep  $X_1$  (projection onto  $X_1$ ).



# DIMENSION REDUCTION VIA PROJECTION

We chuck  $X_1$ , and only keep  $X_2$  (projection onto  $X_2$ ).

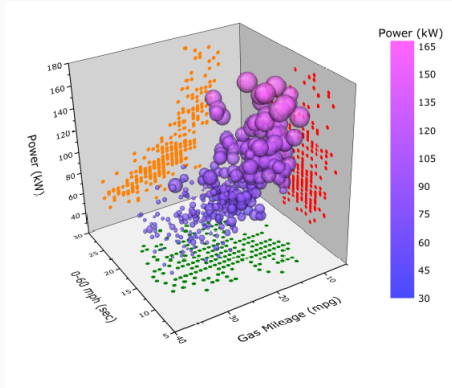


**Question:** Which projection is better? Why?



# DIMENSION REDUCTION VIA PROJECTION

We can do the same for data with three predictors.



**Question:** What is the problem with doing variable selection this way?

Recall that PCA finds the (orthogonal) directions in which the data exhibits the maximum variance.

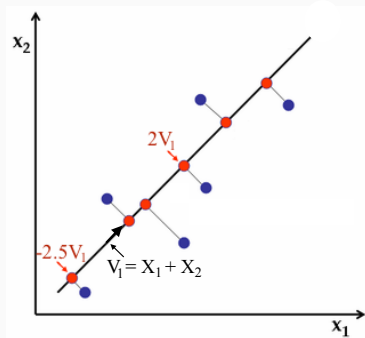
The top PCA component is the direction (given by a vector) along which the data has maximum variance; the second component of the PCA is the direction (orthogonal to the first) along which the data exhibits the “second greatest” amount of variance, etc.

Each component of the PCA is a linear combination of the original set of predictors, for example

$$\underbrace{v_1 = x_1 + x_2}_{\text{component 1}}, \quad \underbrace{v_2 = x_1 - 3x_2}_{\text{component 2}}$$

## PCA FOR DIMENSION REDUCTION

Recall that PCA finds the (orthogonal) directions in which the data exhibits the maximum variance.



When we project the data onto the axes formed by the components, each data point is now a linear combination of the components.

```
from sklearn.decomposition import PCA  
pca = PCA(n_components=2)  
pca.fit(x)  
x_reduced = pca.transform(x)
```

## WHAT DO WE WANT OUT OF THIS?

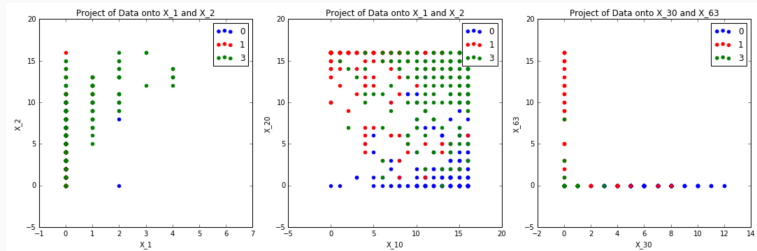
We want to reduce the dimension (the number of features) of the data.

PCA allows us to perform automated dimension reduction - we project the data onto the top PCA components.

**Question:** Is performing PCA for dimension reduction helpful for our classification task?

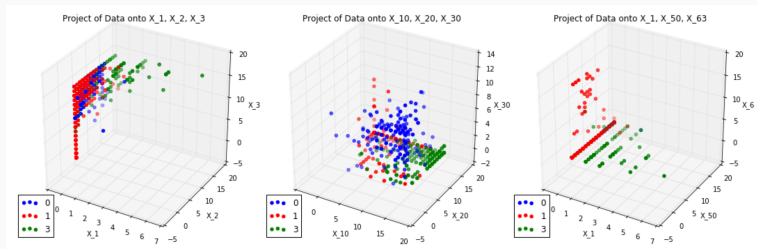
**Try it on your dataset!** Do Steps 1 and 2.

# WHAT DO WE WANT OUT OF THIS?



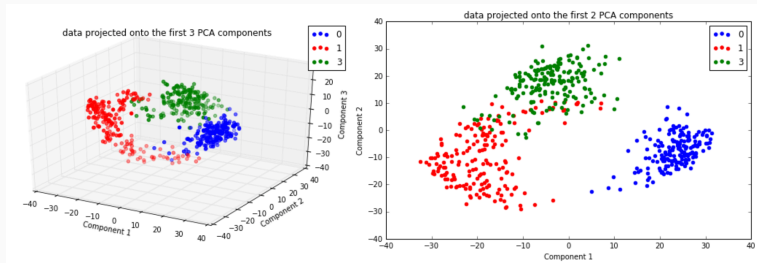
Is randomly selecting two predictors and then projecting the data yielding good results?

# WHAT DO WE WANT OUT OF THIS?



Is randomly selecting two predictors and then projecting the data yielding good results?

# WHAT DO WE WANT OUT OF THIS?

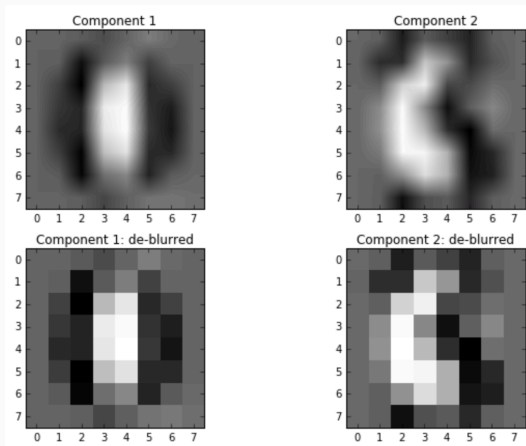


Is projecting the data onto the first 2 or 3 PCA components a good idea? How many components do we need?



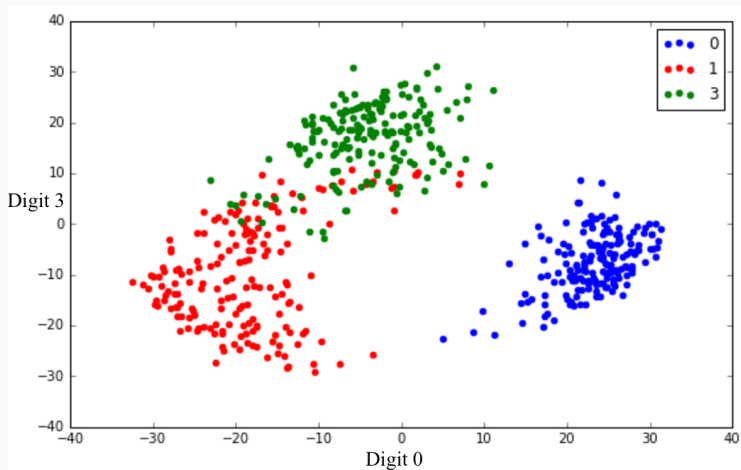
# WHAT DO WE WANT OUT OF THIS?

Each component is a combination of the original predictors (pixels),  
what do the components mean?

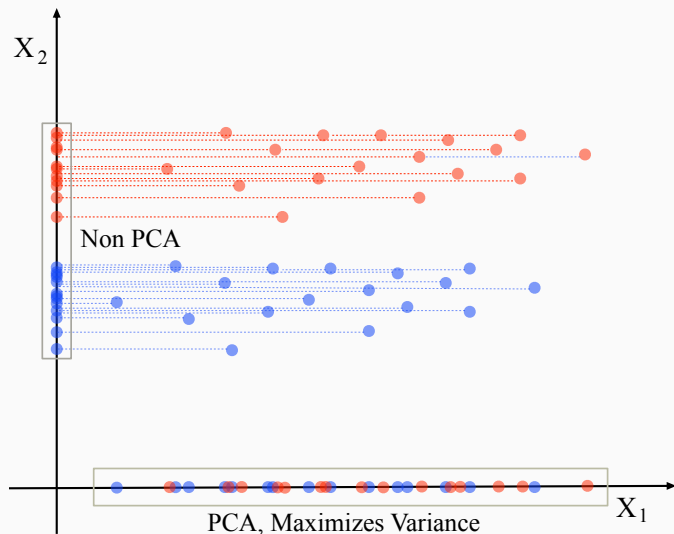


## WHAT DO WE WANT OUT OF THIS?

Each reduced data point is a linear combination of the components, what do the coefficients mean?



## IS PCA ALWAYS GOOD FOR CLASSIFICATION?



## CLASSIFICATION

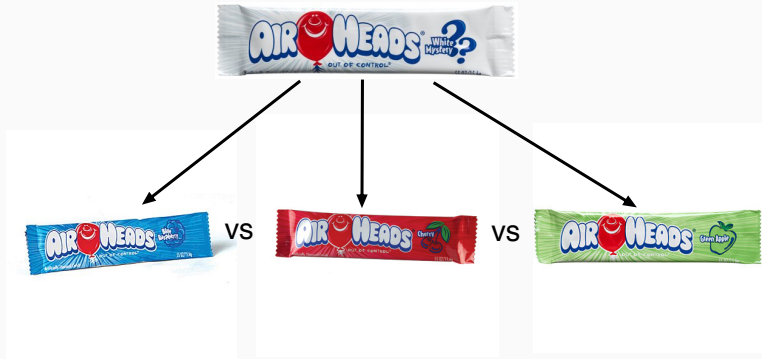
---

Implementing logistic regression using `sklearn`:

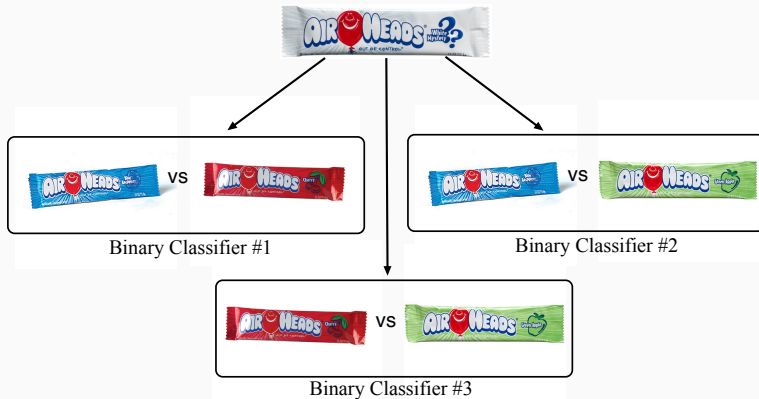
```
from sklearn.linear_model import LogisticRegression as LogReg
log = LogReg() #by default regularization param is 1
log.fit(x, y)
log.predict(x)
```

The coefficients for the predictors are given by `log.coef_`. The intercept is given by `log.intercept_`.

# NAIVE THREE-CLASS CLASSIFICATION



# NAIVE THREE-CLASS CLASSIFICATION



We tally up the predictions from all the classifiers: 2 votes red, 1 vote green, means mystery = red!

Implement a classifier for you handwriting dataset. (Step 3)

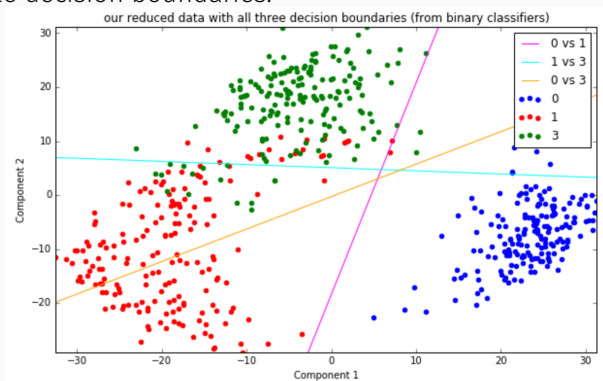
How good was your classifier?



# HOW DID WE DO?

How do we assess the quality of our model?

1. Correct Classification Rate: 0.94
2. Visualize decision boundaries:



3. Other metrics?