



Dimensionality Reduction

Machine Learning Decal

Hosted by Machine Learning at Berkeley

Agenda

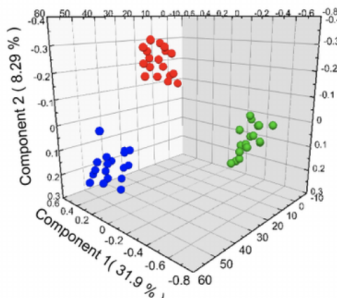
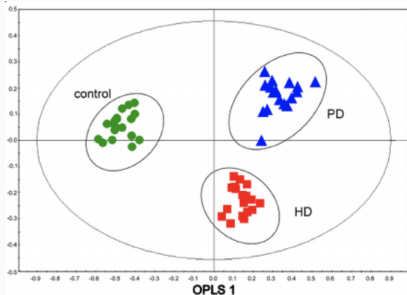
Why reduce dimensions of data?

PCA - Principle Component Analysis

LDA - Linear Discriminant Analysis

Why reduce dimensions of data?

- We can only easily visualize up to 3 dimensions
- Given a dataset with more dimensions, it is hard to understand how the points of different classes are separated

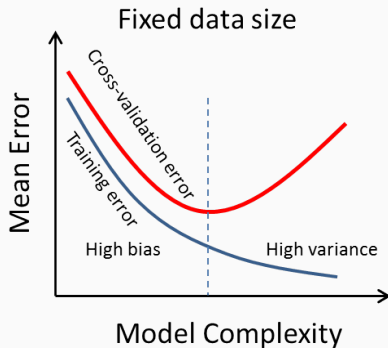


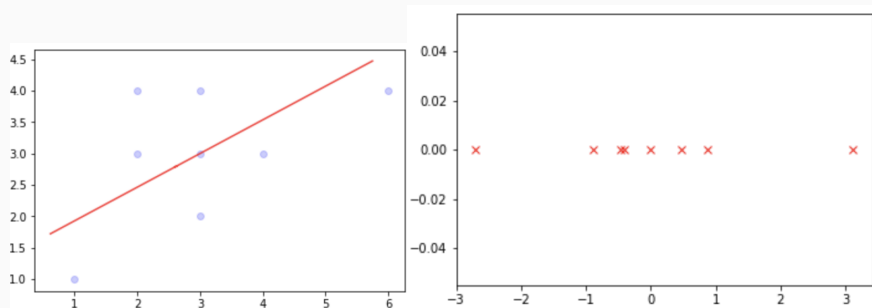
For this dataset, we don't need all 3 dimensions to visualize the clusters

Why do we want a classifier with less features?

- Recall from the bias-variance tradeoff that any feature added will increase the variance of a model
 - Models with more features are more expressive
- Ideally we want to only include features that give enough information to greatly decrease the bias to justify the increase in variance
 - Goal is to get rid of features with redundant information
- Model with fewer features can run faster

Less Complex Classifiers have Lower Variance

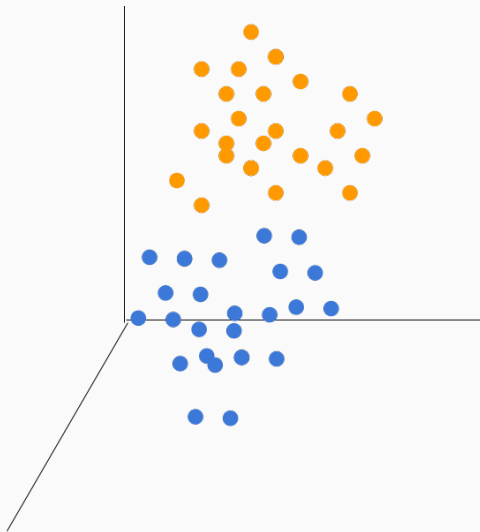


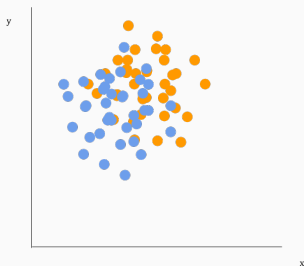


Transform $x \in R^2$ to $z \in R$

This is a form of lossy compression, but may capture enough information about data for given task

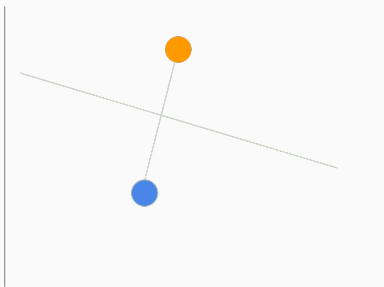
What is a Good Way to Reduce this Dataset to 2D?





Drawbacks:

- This method got rid of an important feature of the data
- Possible that a dataset has important information along all axes so removing any feature removes important information



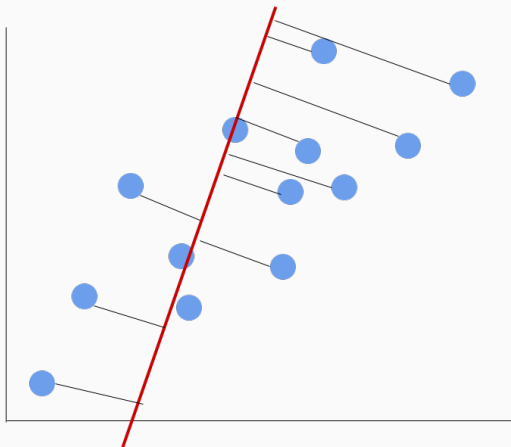
Drawbacks:

- In higher dimensions, points are more spread out, so distance between two points looks similar \Rightarrow standard clustering doesn't work well in high dimensions

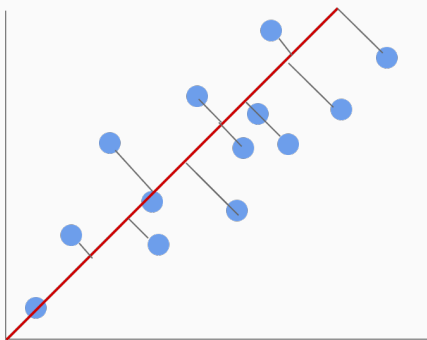
- Combination of many features rather than just eliminating a few
- Should be usable in higher dimensions

PCA - Principle Component Analysis

Which is Best Line to Project Data Onto?

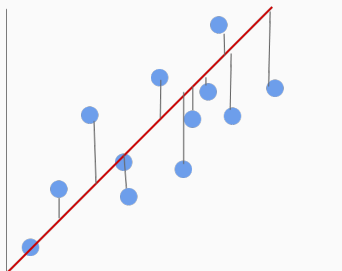
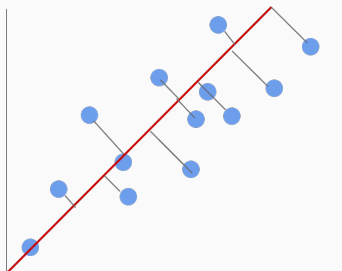


Which is Best Line to Project Data Onto?



- Goal of PCA is to find k dimensions that minimize projection error of a n dimensional dataset

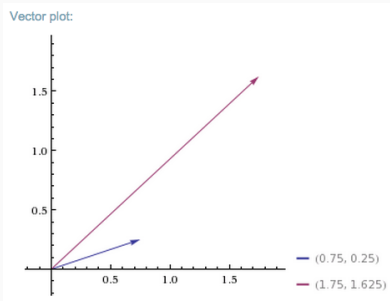
Projection Error vs Linear Regression Residuals

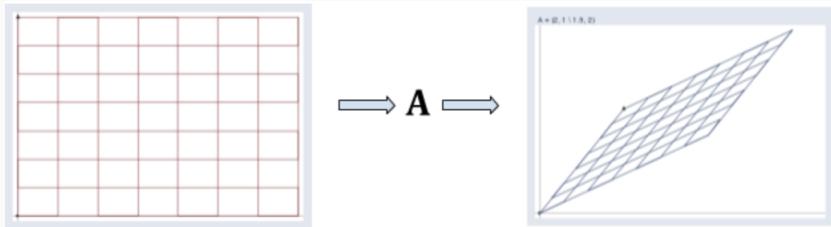


What does a matrix do to a vector?

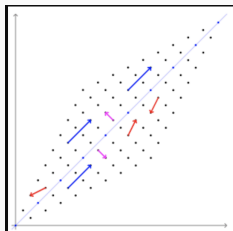
$$A = \begin{bmatrix} 2 & 1 \\ 1.5 & 2 \end{bmatrix}, v = \begin{bmatrix} 0.75 \\ 0.25 \end{bmatrix}$$

$$Av = \begin{bmatrix} 2 & 1 \\ 1.5 & 2 \end{bmatrix} \begin{bmatrix} 0.75 \\ 0.25 \end{bmatrix} = \begin{bmatrix} 1.75 \\ 1.625 \end{bmatrix}$$





- Matrix can be thought of as an operator that linearly transforms a vector



- All vectors are affected after matrix is applied, but direction of eigenvectors is preserved and only length is changed
- Eigenvectors define axis of linear transformation
 - They define directions of force matrix puts on input

"Many mathematical objects can be understood better by breaking them into constituent parts, or finding some properties of them that are universal, not caused by the way we choose to represent them.

For example, integers can be decomposed into prime factors. The way we represent the number 12 will change depending on whether we write it in base ten or in binary, but it will always be true that $12 = 2 \times 2 \times 3$."

- Yoshua Bengio

"Much as we can discover something about the true nature of an integer by decomposing it into prime factors, we can also decompose matrices in ways that show us information about their functional properties that is not obvious from the representation of the matrix as an array of elements.

One of the most widely used kinds of matrix decomposition is called eigen-decomposition, in which we decompose a matrix into a set of eigenvectors and eigenvalues."

- Yoshua Bengio

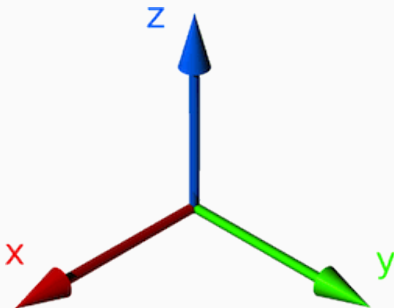
$$Av = \lambda v$$

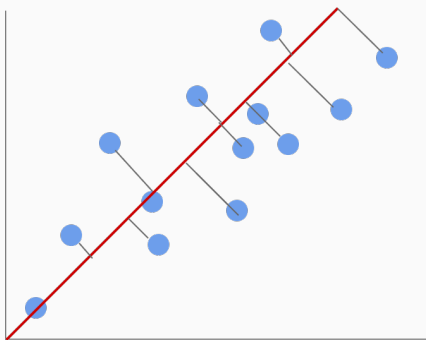
where A is a matrix, v is a eigenvector of A , λ is eigenvalue of A

- Notice that applying A to v only scales the length of v (doesn't change the direction)
- $A_{n \times n}$ can have up to n unique eigenvectors and eigenvalues

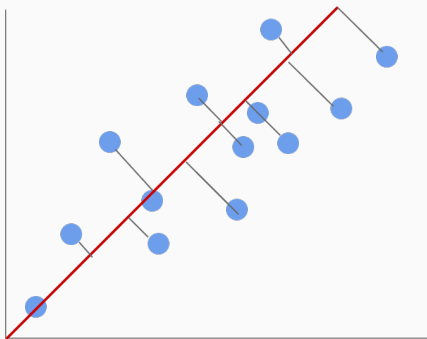
If A is symmetric ($A = A^T$) then eigenvectors are orthogonal
($v_i^T v_j = 0$)

- Orthogonal vectors define different dimension of space





- Goal of PCA is to find k dimensions that minimize projection error of a n dimensional dataset



- Minimizing projection error is equivalent to maximizing variance of points after projection (Proof)

Let there be n samples of $x_i \in \mathbb{R}^d$ so $X \in \mathbb{R}^{n \times d}$

Goal: Find the first k principle components so the variance of X projected on these components is maximized

- Must first center X around the origin so variance is not skewed
 - Variance is defined relative to the mean
- Let $\bar{X} = X - \mu$

Goal: Find the best k principle components so the variance of X projected on these components is maximized

Because covariance is given by

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T = \frac{1}{n} (X - \mu)^T (X - \mu)$$

let the sample covariance after \bar{X} is projected onto v be

$$\Sigma = \frac{1}{n} v^T \bar{X}^T \bar{X} v$$

For now, we will constrain v to have unit norm so $v^T v = 1$

Goal: Find the first principle component so the variance of X projected on these components is maximized

$$\max_{v: v^T v = 1} \frac{1}{n} v^T \bar{X}^T \bar{X} v$$

$$L(v) = \frac{1}{n} v^T \bar{X}^T \bar{X} v - \lambda(v^T v - 1)$$

$$\nabla L(v) = 2\bar{X}^T \bar{X} v - 2\lambda v = 0$$

$$\bar{X}^T \bar{X} v = \lambda v$$

v is an eigenvector of the covariance matrix

Goal: Find the first principle component so the variance of X projected on these components is maximized

$$\max_{v: v^T v = 1} \frac{1}{n} v^T \bar{X}^T \bar{X} v$$

$$\bar{X}^T \bar{X} v = \lambda v$$

Substituting this back into the original objective $\frac{1}{n} v^T \bar{X}^T \bar{X} v$ gives
 $\frac{1}{n} v^T \lambda v = \frac{1}{n} \lambda v^T v$

$$\max_{v: v^T v = 1} \frac{1}{n} v^T \bar{X}^T \bar{X} v = \max_{v: v^T v = 1} \frac{1}{n} \lambda v^T v = \max \frac{\lambda}{n}$$

v is the eigenvector corresponding to the largest eigenvalue of the covariance matrix

Given n samples $x_i \in \mathbb{R}^d$

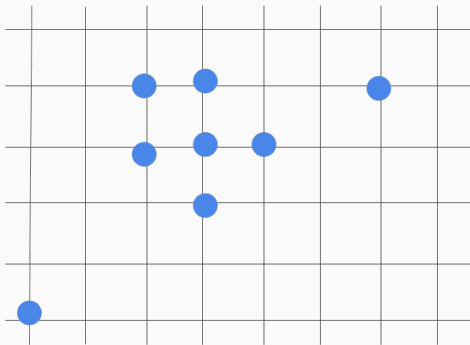
1. Compute the covariance matrix of the centered data

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

2. Find the eigenvectors v_i and eigenvalues λ_i of the covariance matrix Σ
3. Choose the k largest eigenvalues. Their corresponding eigenvectors will form the projection matrix $P_{d \times k}$
4. Project each datapoint using P to get $z_i \in \mathbb{R}^k$

Let X contain the samples $x_i \in \mathbb{R}^2$ in it's rows

$$X = \begin{bmatrix} 1 & 2 & 2 & 3 & 3 & 3 & 4 & 6 \\ 1 & 3 & 4 & 2 & 3 & 4 & 3 & 4 \end{bmatrix}^T$$



Let X contain the samples $x_i \in \mathbb{R}^2$ in it's rows

$$X = \begin{bmatrix} 1 & 2 & 2 & 3 & 3 & 3 & 4 & 6 \\ 1 & 3 & 4 & 2 & 3 & 4 & 3 & 4 \end{bmatrix}^T$$

$$\mu = \frac{1}{8} \begin{bmatrix} 1 + 2 + 2 + 3 + 3 + 3 + 4 + 6 \\ 1 + 3 + 4 + 2 + 3 + 4 + 3 + 4 \end{bmatrix}^T = \begin{bmatrix} 3 & 3 \end{bmatrix}$$

$$X - \mu = \begin{bmatrix} -2 & -1 & -1 & 0 & 0 & 0 & 1 & 3 \\ -2 & 0 & 1 & -1 & 0 & 1 & 0 & 1 \end{bmatrix}^T$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T$$

$$X - \mu = \begin{bmatrix} -2 & -1 & -1 & 0 & 0 & 0 & 1 & 3 \\ -2 & 0 & 1 & -1 & 0 & 1 & 0 & 1 \end{bmatrix}^T$$

$$\Sigma = \frac{1}{8} \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix} + \dots + \begin{bmatrix} 9 & 3 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0.75 \\ 0.75 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2 & 0.75 \\ 0.75 & 1 \end{bmatrix}$$

Σ has eigenvectors $v_1 = \begin{bmatrix} 0.8816746 \\ 0.47185793 \end{bmatrix}$ and $v_2 = \begin{bmatrix} -0.47185793 \\ 0.8816746 \end{bmatrix}$

Eigenvalues are $\lambda_1 = 2.4013$, $\lambda_2 = 0.5986$

Suppose we are trying to reduce this dataset to 1-D:

- Because the largest eigenvalue is λ_1 , the projection matrix

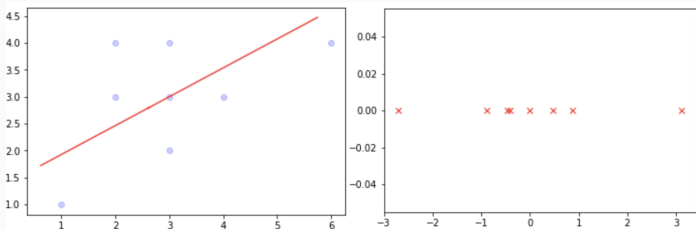
$$P = v_1 = \begin{bmatrix} 0.8816746 \\ 0.47185793 \end{bmatrix}$$

$$\text{Transform } X = \begin{bmatrix} 1 & 2 & 2 & 3 & 3 & 3 & 4 & 6 \\ 1 & 3 & 4 & 2 & 3 & 4 & 3 & 4 \end{bmatrix}^T \text{ using } P$$

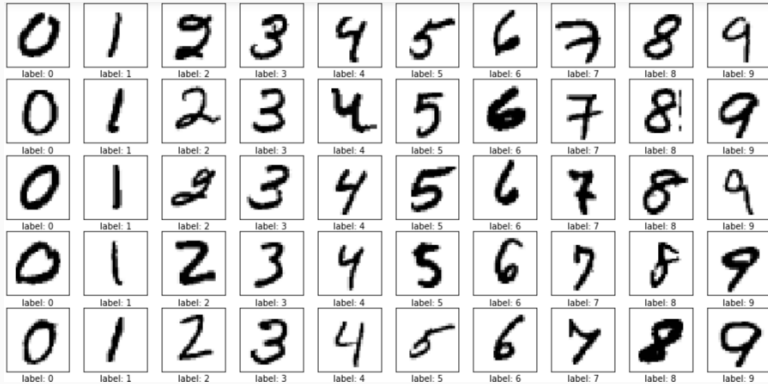
$$(X - \mu)P + \mu = \begin{bmatrix} -2.70 & -0.88 & -0.40 & -0.47 & 0.0 & 0.47 & 0.88 & 3.11 \end{bmatrix}$$

After the transformation, data in the 2-D space can be represented in the 1D space

$$XP = \begin{bmatrix} -2.70 & -0.88 & -0.40 & -0.47 & 0.0 & 0.47 & 0.88 & 3.11 \end{bmatrix}^T$$







LDA - Linear Discriminant Analysis

PCA

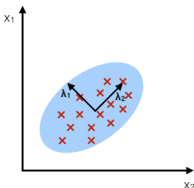
- Only uses the features of the data to find a lower dimensional subspace \Rightarrow unsupervised technique
- Maximizes the variance of the data points in the projected space

LDA

- Uses both features and data to reduce dimensions \Rightarrow supervised technique
- Maximizes the separation between two classes in projected

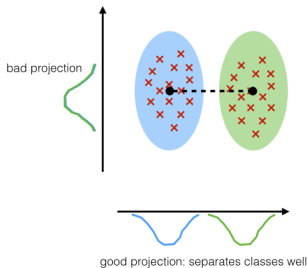
PCA:

component axes that maximize the variance



LDA:

maximizing the component axes for class-separation



Find a lower dimensional space so when the data is projected on that space

- Distance between classes is maximized
- Distance of points within a class is minimized

Distance between classes is the distance between the **means of each class**

Suppose there are 2 classes:

$$\mu_1 = \frac{1}{N_1} \sum_{x \in C_1} x, \mu_2 = \frac{1}{N_2} \sum_{x \in C_2} x$$

After projecting onto a vector w :

$$\mu'_1 = \frac{1}{N_1} \sum_{x \in C_1} w^T x, \mu'_2 = \frac{1}{N_2} \sum_{x \in C_2} w^T x$$

The distance between means

$$\|\mu'_1 - \mu'_2\|_2^2 = \|w^T \mu_1 - w^T \mu_2\|_2^2$$

Distance between classes is the distance between the **means of each class**

$$\|\mu'_1 - \mu'_2\|_2^2 = \|w^T \mu_1 - w^T \mu_2\|_2^2$$

This can be put in matrix form as

$$(w^T \mu_1 - w^T \mu_2)^2 = w^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T w = w^T \Sigma_m w$$

where

$$\Sigma_m = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

We can check points in a class are close to each other by measuring the in-class variance

$$\sigma_i^2 = \sum_{x \in C_i} (w^T x - \mu'_i)^2$$

Then the total in-class variance is

$$\sum_{i=1}^c \sigma_i^2$$

We can check points in a class are close to each other by measuring the in-class variance

$$\sigma_i^2 = \sum_{x \in C_i} (w^T x - \mu'_i)^2$$

This can be put in matrix form as

$$\begin{aligned}\sigma_i^2 &= \sum_{x \in C_i} (w^T x - w^T \mu_i)^2 = \sum_{x \in C_i} w^T (x - \mu_i)(x - \mu_i)^T w \\ &= w^T (\sum_{x \in C_i} (x - \mu_i)(x - \mu_i)^T) w = w^T \Sigma_c w\end{aligned}$$

where

$$\Sigma_c = \sum_{x \in C_i} (x - \mu_i)(x - \mu_i)^T$$

Suppose there are two classes c_1 and c_2

We can maximize the distance between classes and minimize the in-class variance with:

$$\begin{aligned}\operatorname{argmax}_w \frac{\|\mu'_1 - \mu_2\|_2^2}{\sigma_1 + \sigma_2} &= \operatorname{argmax}_w \frac{w^T \Sigma_m w}{w^T \Sigma_{c1} w + w^T \Sigma_{c2} w} \\ &= \operatorname{argmax}_w \frac{w^T \Sigma_m w}{w^T (\Sigma_{c1} + \Sigma_{c2}) w}\end{aligned}$$

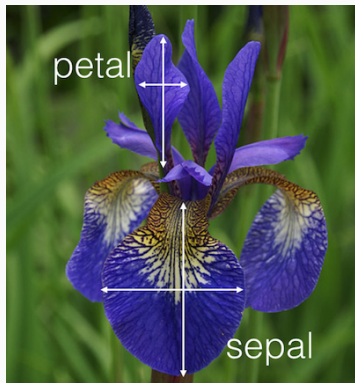
$$\operatorname{argmax}_w \frac{w^T \Sigma_m w}{w^T (\Sigma_{c1} + \Sigma_{c2}) w} = \operatorname{argmax}_w \frac{w^T \Sigma_m w}{w^T (\Sigma_c) w}$$

This takes the form of a generalized Rayleigh Quotient, so any w that maximizes this equation must satisfy

$$\Sigma_m w = \lambda \Sigma_c w = \Sigma_c \lambda w$$

$$\Sigma_c^{-1} \Sigma_m w = \lambda w$$

To transform the data into k dimensions, select the eigenvectors of $\Sigma_c^{-1} \Sigma_m$ that correspond to the first k dimensions as the projection basis



Both PCA and LDA are linear methods of dimensionality reduction

- The resulting axes after projection are linear combinations of the initial axes
- Can kernelize PCA and LDA for a nonlinear transformation to lower dimensional subspace