

---

---

---

---

---

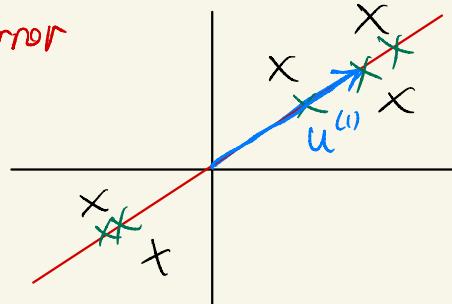


## Problem Formulation:

May 30, 2021

e.g., Task: Reduce 2D data to 1D

PCA: Find a direction (i.e. a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the square projection error



Principal Component Analysis reduces data from  $n$ -dimensions to  $k$ -dimensions ( $k \leq n$ ) by finding  $k$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  onto which to project the data, so as to minimize the projection error.

## Data Preprocessing:

Mean Normalization:

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \text{ replace each } x_j^{(i)} \text{ w/ } x_j - \mu_j$$

Feature Scaling:

Scale features to have a comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

So how do we find  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ ?

PCA Algorithm: & These are the Principal Components

To reduce from  $n$ -dimensions to  $k$ -dimensions,

1. Compute the covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^n \underset{n \times 1}{(x^{(i)})} \underset{1 \times n}{(x^{(i)})^T}, \Sigma \in \mathbb{R}^n$$

2. Compute the eigenvectors of matrix  $\Sigma$

$$[U, S, V] = \text{svd}(\Sigma);$$

↳ singular value decomposition

$$U = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(m)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$\star m=n$   
for PCA

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$\hookrightarrow U_{\text{reduce}} = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & & | \end{bmatrix}^T \in \mathbb{R}^{n \times k}$$

$$z = \begin{bmatrix} 1 & 1 & \dots & 1 \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ 1 & 1 & \dots & 1 \end{bmatrix}^T x$$

$$= \begin{bmatrix} \frac{u^{(1)T}}{\sqrt{k}} \\ \frac{u^{(2)T}}{\sqrt{k}} \\ \vdots \\ \frac{u^{(k)T}}{\sqrt{k}} \end{bmatrix} x$$

$k \times n$

$$\therefore z \in \mathbb{R}^k$$

Summary:

Mean normalization, feature scaling.

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$$

$$[U, S, V] = \text{svd}(\Sigma)$$

$$U_{\text{reduce}} = U(:, 1:k)$$

$$z = U_{\text{reduce}}^T * x$$

## Reconstruction from Compressed Data:

$$Z^{(i)} = U_{\text{reduce}}^T X^{(i)} \rightarrow \text{PCA}$$

$$X_{\text{approx}}^{(i)} = U_{\text{reduce}} Z^{(i)} \rightarrow \text{inverse operation}$$

n \times 1 \quad n \times k \quad k \times 1

Choosing  $k$  (# of Principal Components):

Avg Squared Projection Error:  $\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2$

Total Variation in Data:  $\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2$

Typically, we choose as small a " $k$ " as possible s.t.

$$\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

i.e. 99% of variance is retained.

↳ The compressed data isn't too far off from the original.

## Algorithm:

Try PCA w/  $k=1, k=2, \dots$

1. Compute  $U$  reduce,  $z^{(1)}, \dots z^{(m)}$ ,  $x_{approx}^{(1)}, \dots x_{approx}^{(m)}$

$$2. \text{Check if } \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$$

$\hookrightarrow [U, S, V] = svd(\Sigma)$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \Sigma_{nn} \end{bmatrix}$$

For a given  $k$ ,

$$1 - \frac{\sum_{i=1}^k \Sigma_{ii}}{\sum_{i=1}^n \Sigma_{ii}} \leq 0.01 ?$$

$\hookrightarrow \boxed{\frac{\sum_{i=1}^k \Sigma_{ii}}{\sum_{i=1}^n \Sigma_{ii}} \geq 0.99}$

Easy way to compute 2.

Pick smallest  $k$  for which this is satisfied.

## Supervised Learning Speedup:

e.g.  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$   
 $x^{(i)} \in {}^{10000}R^1$  for  $100 \times 100$  image data

$$x^{(1)}, x^{(2)} \dots x^{(m)} \in {}^{10000}R^1$$

$\downarrow$  PCA

$$z^{(1)}, z^{(2)} \dots z^{(m)} \in {}^{100}R^1$$

## New Training Set:

$$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}) \dots (z^{(m)}, y^{(m)})$$

For new predictions, run them through the same mapping found by PCA to get  $z$ , then use

$$h_\theta(z) = \frac{1}{1 + e^{-\theta^T z}}.$$

↳ i.e. For  $x_{cv}^{(i)}, x_{test}^{(i)}$ .

\* Only run PCA on training set to get mapping, then use mapping for validation and testing sets.

## Main Applications of PCA:

### 1. Compression

- a. Reduce memory/disk needed to store data
- b. Speed up learning algorithm

### 2. Visualization

↳  $k=2, 3$

⊗ Don't use PCA for overfitting; use regularization instead!