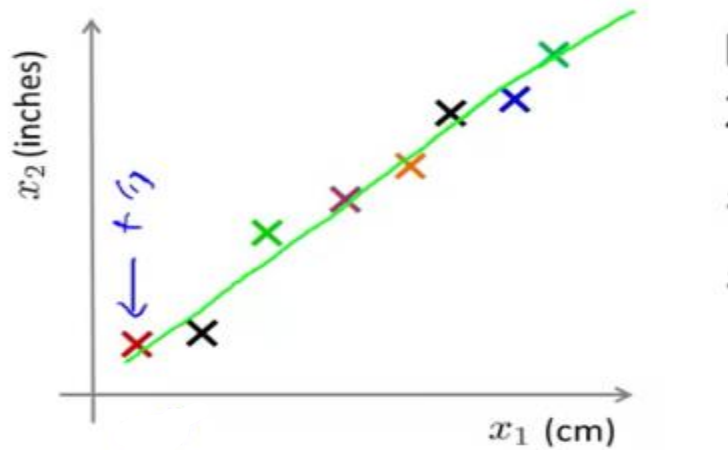


Data Compression:

Data compression is reducing data from a higher dimension to a lower dimension. Example: from 1000 dimensions to 100 dimensions or from 3D to 2D or from 2D to 1D.

Consider the data given below:

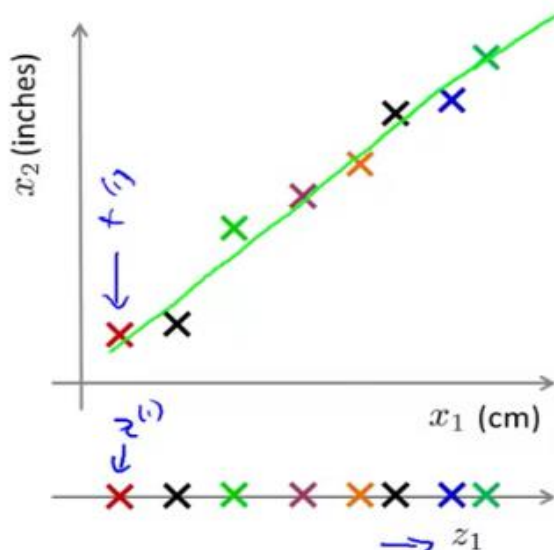
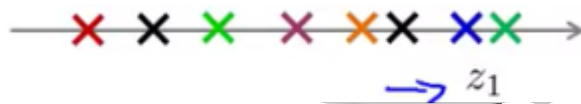


Example 1:

We have 2-dimensional data plotted on the graph. Though the data should line on a straight line it doesn't due to rounding off that has been done. Now we want to reduce this data to 1D or simply a line.

Note: x^1 represents first row of the data. There are 2 columns x_1 and x_2 in the data.

To convert this data to 1D we take projections of x^1, x^2, \dots, x^m and then obtain the line on which all the projections fall. This line is called z_1 and has only 1 dimension.



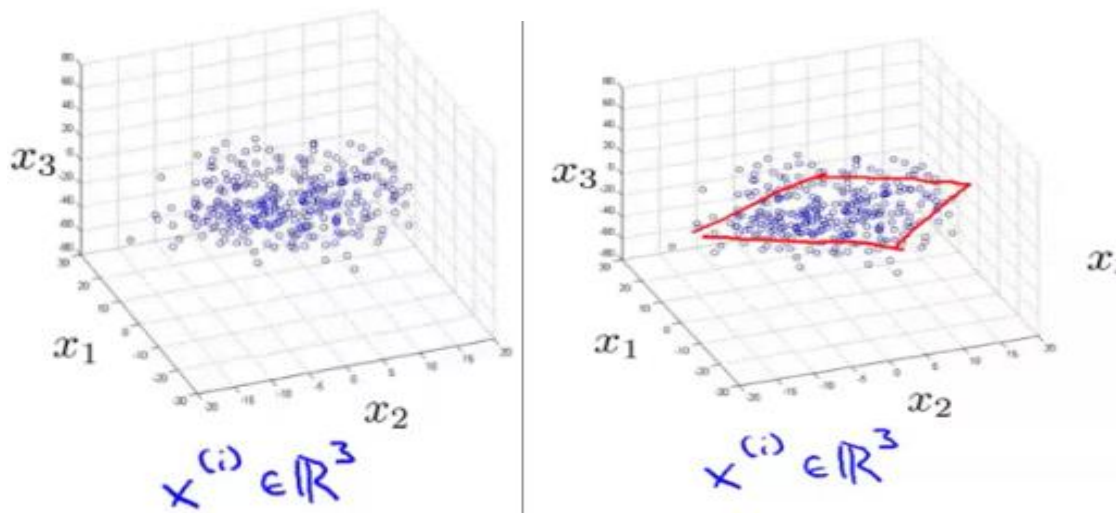
Reduce data from
2D to 1D

$$\begin{matrix} x^{(1)} \in \mathbb{R}^2 \\ x^{(2)} \end{matrix} \rightarrow z^{(1)} \in \mathbb{R}$$

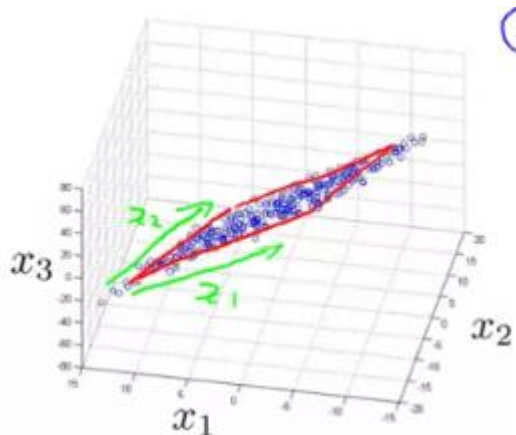
Example 2:

Reducing data from 3D to 2D

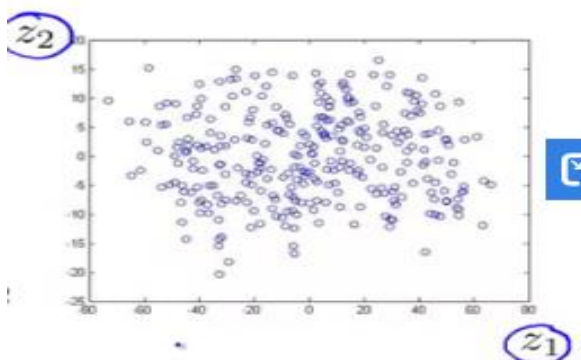
Consider the dataset given below:



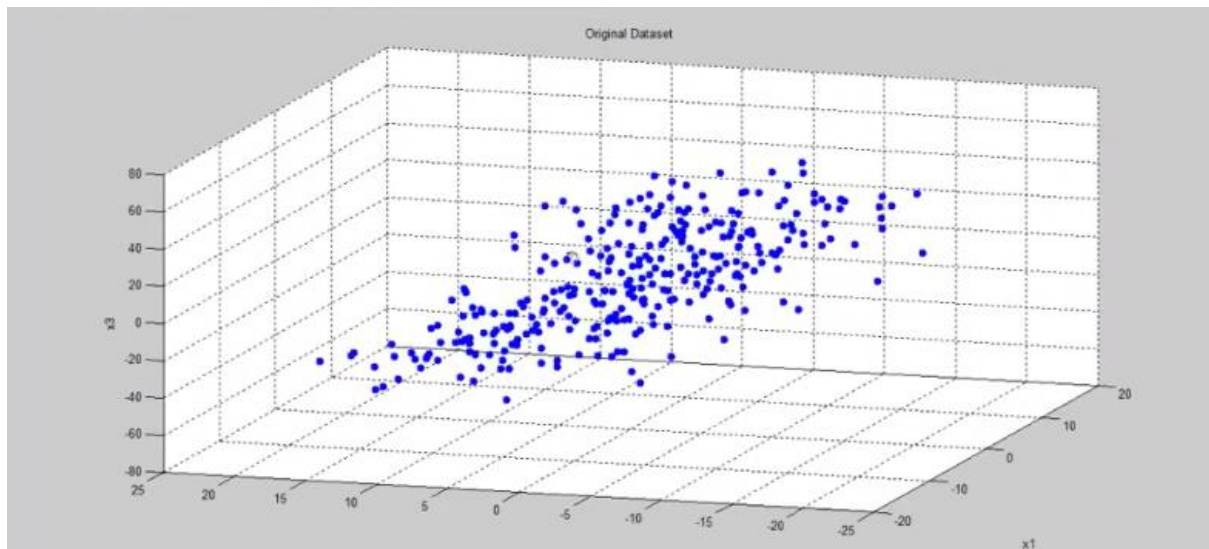
Project the dataset on a plane:



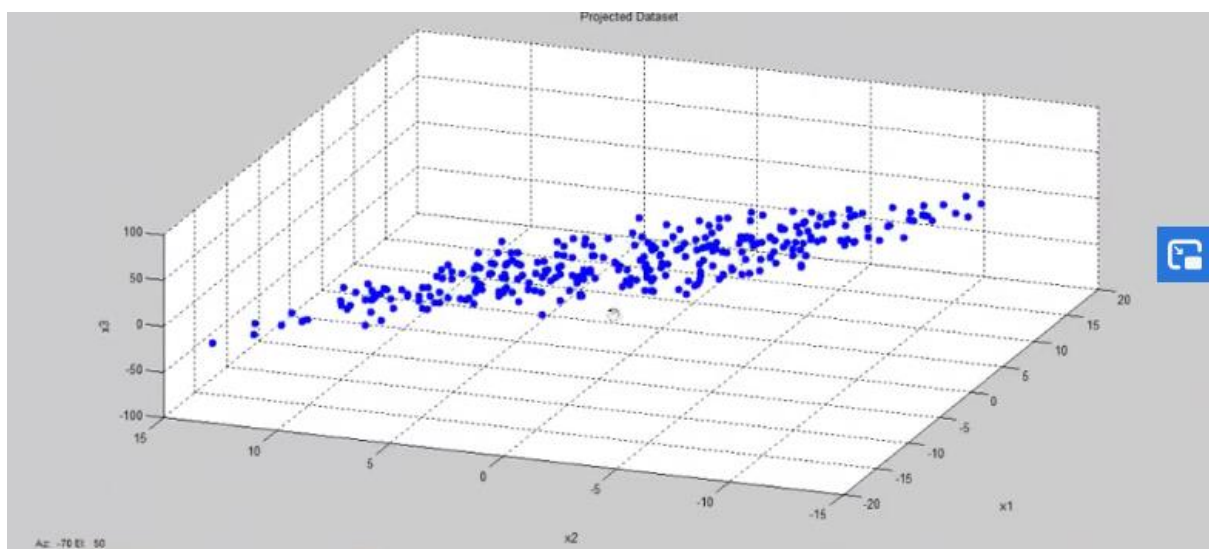
2D image of the plane:



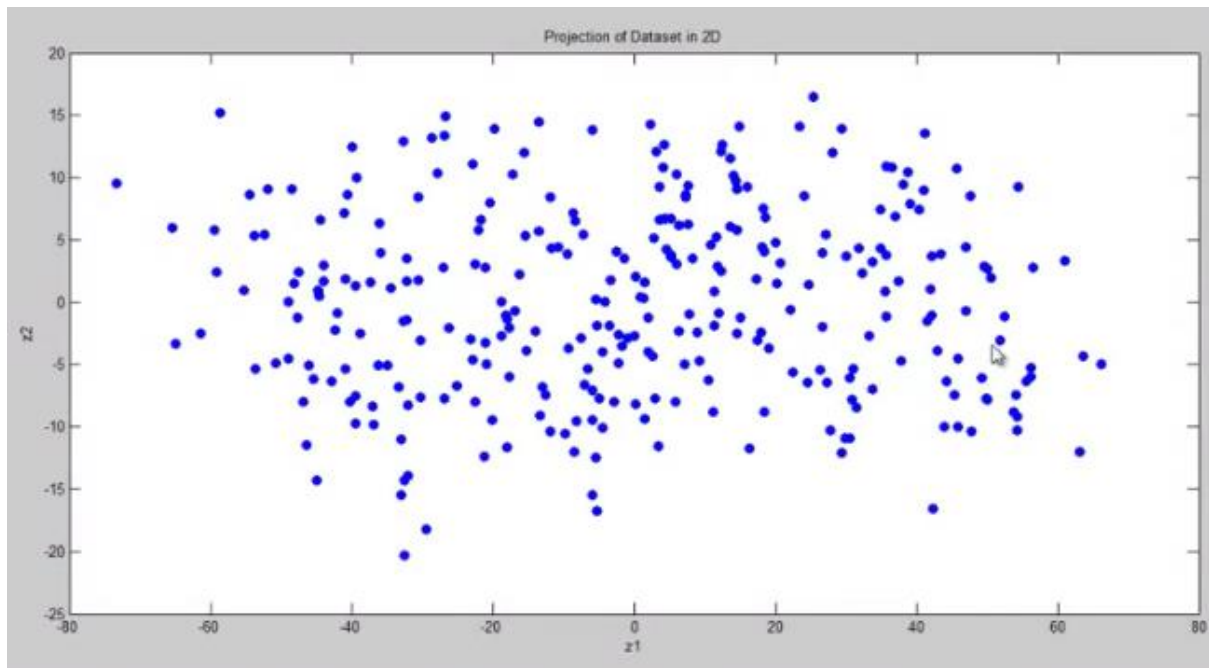
3D representation of the above example:



Projecting the data on single plane:



Projected dataset in 2D:



Questions:

Q1:

Data

Suppose we apply dimensionality reduction to a dataset of m examples $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, where $x^{(i)} \in \mathbb{R}^n$. As a result of this, we will get out:

- ☐ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}$ of k examples where $k \leq n$.
- ☐ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}$ of k examples where $k > n$.
- ☒ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ of m examples where $z^{(i)} \in \mathbb{R}^k$ for some value of k and $k \leq n$.

Correct

☐ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ of m examples where $z^{(i)} \in \mathbb{R}^k$ for some value of k and $k > n$.

Continue

two dimensional. So that's

Andrew Ng

Q2:

Data

Suppose you have a dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \mathbb{R}^n$. In order to visualize it, we apply dimensionality reduction and get $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ where $z^{(i)} \in \mathbb{R}^k$ is k -dimensional. In a typical setting, which of the following would you expect to be true? Check all that apply.

- ☐ $k > n$
Un-selected is correct
- ☒ $k \leq n$
Correct
- ☐ $k \geq 4$
Un-selected is correct
- ☒ $k = 2$ or $k = 3$ (since we can plot 2D or 3D data but don't have ways to visualize higher dimensional data)
Correct

Continue

per. person
GDP
(economic
activity)

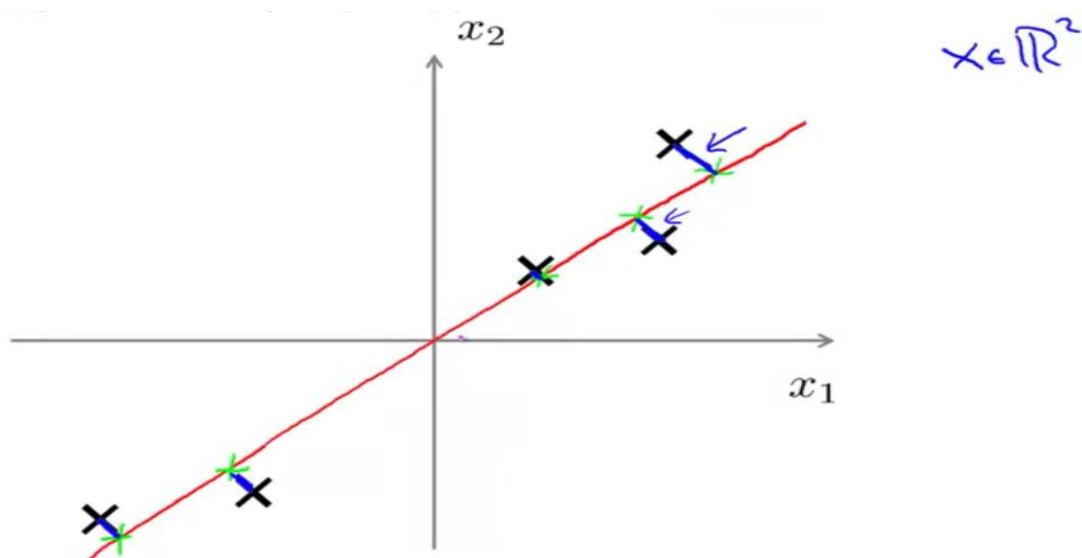
size
GDP

Andrew Ng

Principal Component Analysis:

Principal component analysis is used to find the line/plane/any other $(n-1)^{\text{th}}$ dimensional component of a n dimensional data from which the points have the minimum perpendicular distance. Using PCA, we try to minimize the square of the perpendicular distance between the points and the $(n-1)^{\text{th}}$ dimensional surface.

Consider that we have a 2D dataset and want to convert it into 1D. PCA tries to find the best possible line to do this job. Now, the best possible line is the one from which the points have the minimum perpendicular distance:

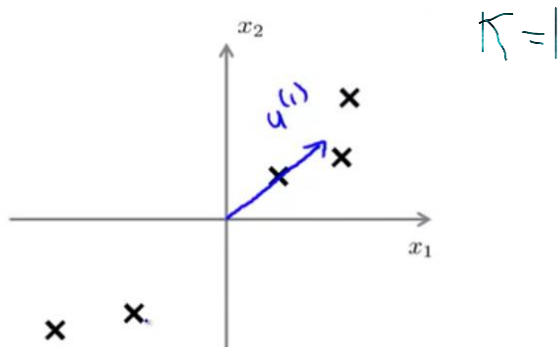


The lines projecting from the points to the red line depict the perpendicular distance of those points from the red line.

To find the best possible lower dimension surface we use vectors:

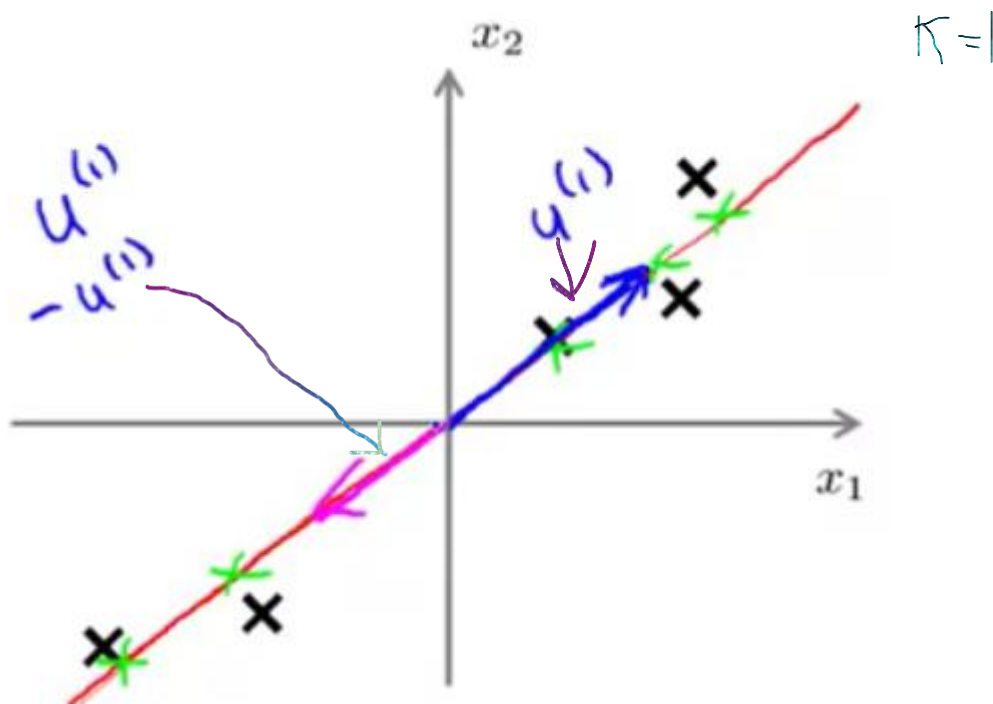
Reduce from n -dimension to k -dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

Example 1: For example, in reduction of 2D to 1D, to find the best possible line, we take a vector u^1 which when extended back and forth forms the line we are looking for:

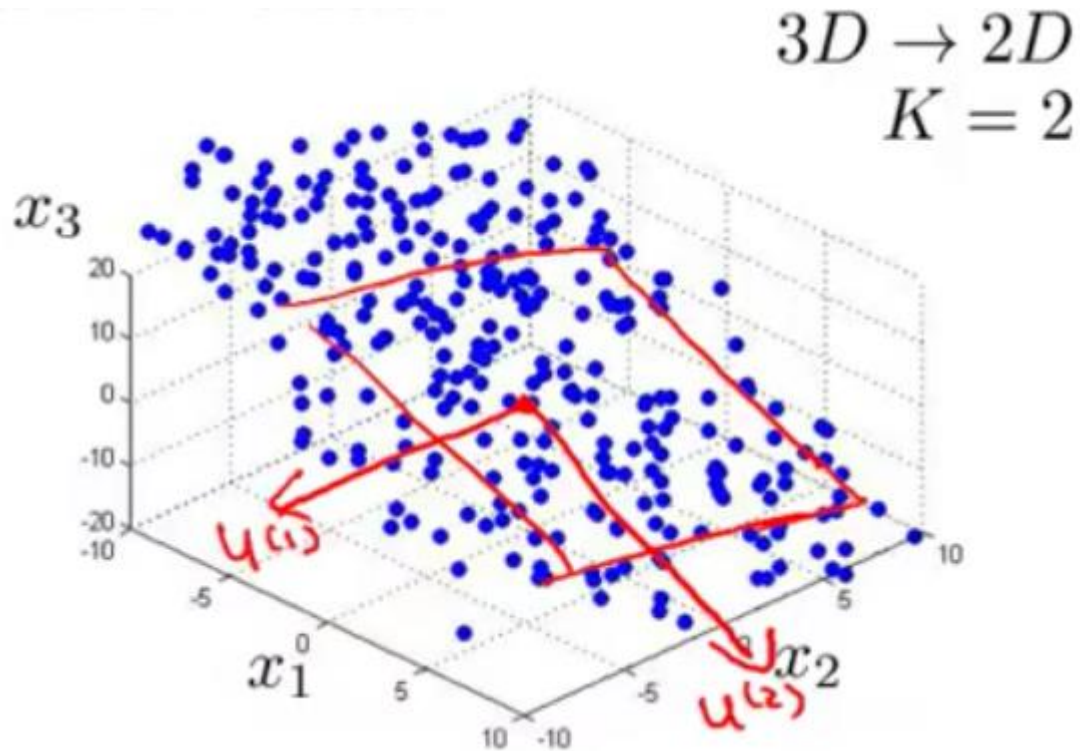


Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

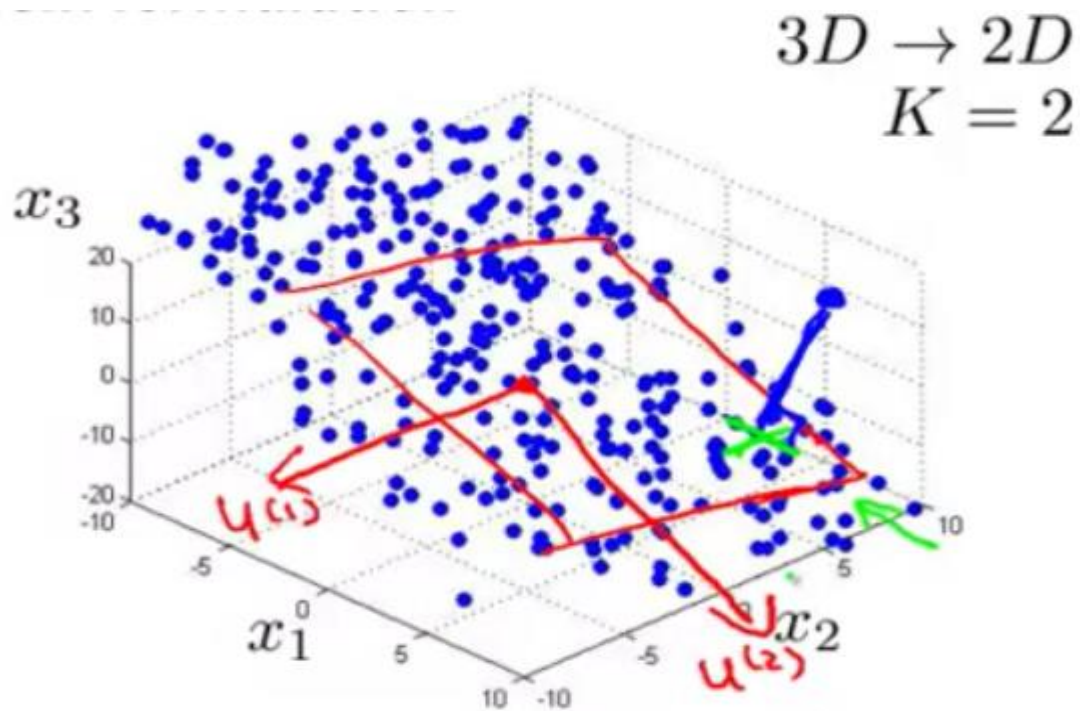
Note: Using $-u^1$ instead of u^1 does not matter since we have to extend the line in both directions.



Example 2: For example, in reduction of 3D to 2D, to find the best possible plane, we take vectors u^1 and u^2 which when extended back and forth form the plane we are looking for:



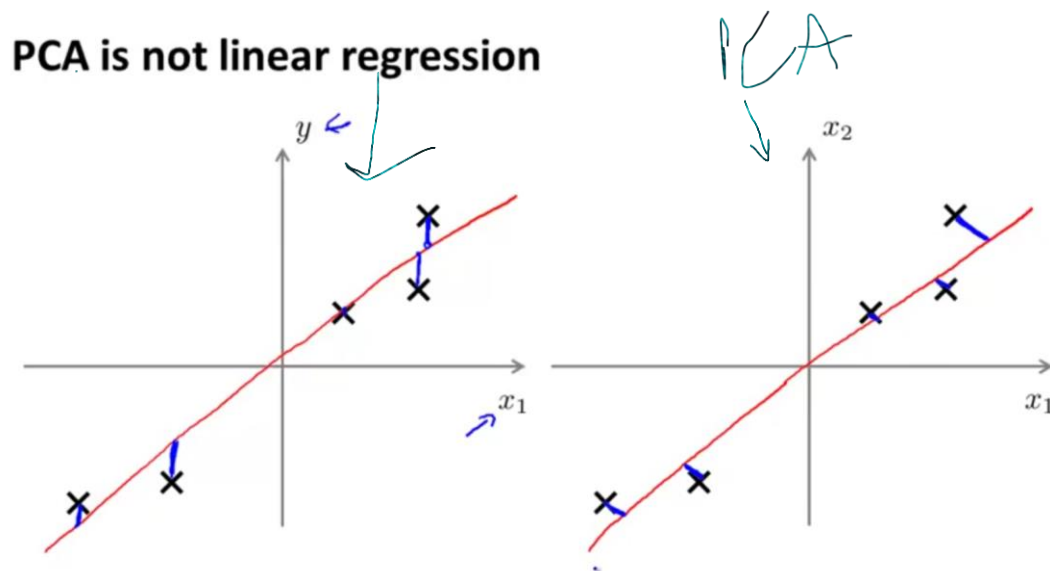
NOTE: We are trying to minimize the perpendicular distance of the points from the plane:



Difference between PCA and Linear Regression:

Though PCA seems to be the same as linear regression, the two have substantial differences:

1. In PCA, we minimize square of the perpendicular distance between the points and the lower dimension surface. Whereas in Linear Regression, we try to minimize the vertical distance between the points and the lower dimensional surface:



Andrew Ng

2. When you're doing linear regression, there is this distinguished variable y that we're trying to predict. In linear regression we take all the values of x and try to use that to predict y . Whereas in PCA, there is no distinguish, or there is no special variable y that we're trying to predict. And instead, we have a list of features, x_1, x_2, \dots and so on up to x_n , and all of these features are treated equally, so no one of them is special.

Applying the PCA Algorithm:

1. First, we perform feature scaling:

Data preprocessing

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ \leftarrow

Preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each $x_j^{(i)}$ with $x_j - \mu_j$.

If different features on different scales (e.g., x_1 = size of house, x_2 = number of bedrooms), scale features to have comparable range of values.

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

2. Then we use the PCA algorithm by doing the following steps:

Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^T}_{1 \times n} \quad \text{--- } n \times n$$

Compute "eigenvectors" of matrix Σ :

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

Sigma

svd

Singular value decomposition
 $\text{eig}(\text{Sigma})$

$n \times n$ matrix.

Note: above is the code for computing eigenvectors in octave. We can also use `eig(Sigma)` instead of `svd(Sigma)`.

Note: u which represents the vector was mentioned in the above pages, is a $n \times n$ dimensional vector. To obtain our lower k dimensional vector, we choose k columns from u .

Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^T}_{1 \times n} \quad \text{Sigma } n \times n$$

Compute "eigenvectors" of matrix Σ :

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

\rightarrow Singular value decomposition
 $\text{eig}(\text{Sigma})$

$n \times n$ matrix.

$$U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \end{bmatrix}$$

$\underbrace{\hspace{10em}}_k$

$$U \in \mathbb{R}^{n \times n}$$

$$u^{(1)}, \dots, u^{(k)}$$

Andrew Ng

We obtain values of z which is a k dimensional vector but multiplying the transpose of new u with n rows and k columns with X which is a $n \times n$ matrix (Note: the new vector is names as U_{reduce}):

Principal Component Analysis (PCA) algorithm

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

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

$\underbrace{\hspace{10em}}_k$

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

$$z = \underbrace{\begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}^T}_{U_{\text{reduce}} \quad n \times k} \times \underbrace{\begin{bmatrix} \text{---} (u^{(1)})^T \text{---} \\ \vdots \\ \text{---} (u^{(k)})^T \text{---} \end{bmatrix}}_{k \times n \quad n \times 1} x$$

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

Andrew Ng

Summary of using PCA:

Principal Component Analysis (PCA) algorithm summary

→ After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

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

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

→ $\text{Ureduce} = U(:, 1:k)$;

→ $z = \text{Ureduce}' * x$;

↑

↑

$x \in \mathbb{R}^n$

~~$x_0 = 1$~~

$$X = \begin{bmatrix} - & x^{(1)T} & - \\ & \vdots & \\ - & x^{(m)T} & - \end{bmatrix}$$

→ $\text{Sigma} = (1/m) * X' * X$

Question:

In PCA, we obtain $z \in \mathbb{R}^k$ from $x \in \mathbb{R}^n$ as follows:

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

Which of the following is a correct expression for z_j ?

☐ $z_j = (u^{(k)})^T x$

☐ $z_j = (u^{(j)})^T x_j$

☐ $z_j = (u^{(j)})^T x_k$

☒ $z_j = (u^{(j)})^T x$

Correct

Choosing the value of K:

Choosing k (number of principal components)

Average squared projection error: $\frac{1}{n} \sum_{i=1}^n \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

Total variation in the data: $\frac{1}{n} \sum_{i=1}^n \|x^{(i)}\|^2$

Typically, choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq \underline{0.01} \quad \underline{(1\%)}$$

- “99% of variance is retained”

We calculate the average squared error, which is the square of the perpendicular difference between the actual x values and their projection.

Variance gives us an idea of how many features we have, more features => more variance. 99% variance retained simply means that even though we have used PCA to reduce the number of features, still 99% of the original variance is obtained.

So the 0.01 or 1% is the limit of variance sacrifice we have set ourselves meaning that we can allow maximum 1% variance to get removed in the process of data compression. Now, this limit can be 0.05 or 5% or 10% also:

Choosing k (number of principal components)

Average squared projection error: $\frac{1}{n} \sum_{i=1}^n \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

Total variation in the data: $\frac{1}{3} \sum_{i=1}^3 \|x^{(i)}\|^2$

Typically, choose k to be smallest value so that

$$\rightarrow \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq \frac{0.01}{0.05} = \frac{(1\%)}{5\%} = 10\%$$

→ "99% of variance is retained"
95 to 90%

Methods for choosing k values:

We talk about 2 methods that can be used to get the best suited value of k (the one that minimizes average squared projection error).

Method-1:

We start from k=1 and keep on increasing the value of k until we get the one that gives us average squared projection error ≤ 0.01 .

With each new value of k, we calculate the $n \times k$ matrix U_{reduce} and then calculate average squared projection error.

Algorithm:

Try PCA with $k=1$ ~~$k=2$~~ ~~$k=3$~~ $k=4$...

Compute $U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$

Check if

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

$k=17$

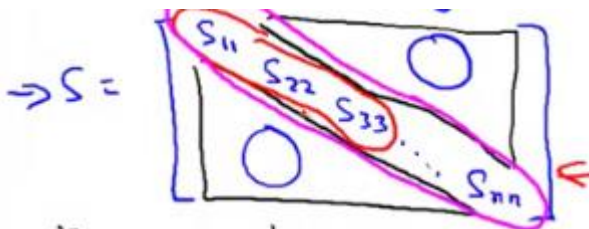
Let us assume that at $k=17$, we got average squared projection error ≤ 0.01 .

Method-2:

In this method we use the variable S which is obtained from $\text{svd}(\text{Sigma})$.

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

Note: The S matrix is a diagonal matrix and looks like this:



Now, we have to calculate the S matrix only once and its diagonal has all the values ranging from S_{11} to S_{nn} . We start from $k=1$ and keep on increasing k until for a given k :

For given k k=3

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

OR

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

Andrew Ng

Method 2 is faster than method 1 because in method 2 we don't have to run PCA and calculate things over and over again, since the diagonal of S matrix has all the values ranging from S_{11} to S_{nn} .

To summarise:

Choosing k (number of principal components)

Algorithm:

Try PCA with $k=1$ k=2 k=3 k=4
 Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

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

k=17

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

$\rightarrow S =$

For given k k=3

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

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

Andrew Ng

Choosing k (number of principal components)

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

Pick smallest value of k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

k=100

(99% of variance retained)

Even if we choose value of K manually, we can show the credibility of the chosen value by calculating:

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}}$$

And checking if:

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

Question:

The screenshot shows a web browser with multiple tabs open, including Coursera, Math behind SVM, and YouTube. The active tab is Coursera, displaying a quiz for the course "Machine Learning" (Week 8) on the topic "Choosing the Number of Principal Components". The quiz question asks for the expression that PCA tries to minimize to project data onto a set of directions $u^{(1)}, \dots, u^{(k)}$. The options are:

- ☐ $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$
- ☐ $\frac{1}{m} \sum_{i=1}^m \|x_{\text{approx}}^{(i)}\|^2$
- ☒ $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$
- ☐ $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} + x_{\text{approx}}^{(i)}\|^2$

The correct answer is the third option, which is highlighted with a green "Correct" message and a "Continue" button. The interface also includes a search bar, user profile (Ansh Tangri), and a taskbar at the bottom showing various application icons and the system clock (9:11 PM, 4/10/2020).

Applications of PCA:

g PCA

Pre

Application of PCA

- Compression

- Reduce memory/disk needed to store data
- Speed up learning algorithm ←

Choose k by % of variance retain

- Visualization

$k=2$ or $k=3$

1. Improving the speed of a supervised learning algorithm.

Supervised learning speedup

→ $(\underline{x^{(1)}}, y^{(1)}), (\underline{x^{(2)}}, y^{(2)}), \dots, (\underline{x^{(m)}}, y^{(m)})$

Extract inputs:

Unlabeled dataset: $\underline{x^{(1)}}, \underline{x^{(2)}}, \dots, \underline{x^{(m)}} \in \mathbb{R}^{10000}$ ←

↓ PCA

$\underline{z^{(1)}}, \underline{z^{(2)}}, \dots, \underline{z^{(m)}} \in \mathbb{R}^{1000}$ ←

New training set:

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

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

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

Andrew Ng

Note: The Mapping obtained from running PCA on the training set is applied on the testing and cross validation sets and not vice-versa.

2. We can reduce high dimensional data to 2D/3D, to effectively visualize it.

Bad uses of PCA:

1. Using PCA to prevent overfitting.

A network error caused the media download to fail part-way.

Bad use of PCA: To prevent overfitting

→ Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the number of features to $k < n$. — 1000

Thus, fewer features, less likely to overfit.

Bad!

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Andrew Ng

Note: Using PCA to solve the problem of overfitting is harmful since it removes some of our features and this can also lead to some important features being removed, which is detrimental for our analysis. Using regularization is much more effective.

2. Using PCA where it shouldn't be used.

Shown below is the general layout of an ML system that people often use.

Design of ML system:

- - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- - Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
- - Train logistic regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_{\theta}(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

Now, step 2 i.e. running PCA to reduce $x^{(i)}$ to get $z^{(i)}$ is not always necessary.

Usually, we should train logistic regression using $x^{(i)}$ values and if there is a lot of problem, only then should we use PCA:

PCA is sometimes used where it shouldn't be

Design of ML system:

- - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- - ~~Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$~~
- - Train logistic regression on $\{(\cancel{x^{(1)}}), y^{(1)}), \dots, (\cancel{x^{(m)}}), y^{(m)})\}$
- - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_{\theta}(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

⇒ How about doing the whole thing without using PCA?

➤ Before implementing PCA, first try running whatever you want to do with the original/raw data $x^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $z^{(i)}$.