Search

Go

Edit

Search

History

Log Out

# Recent changes Help Courses List of courses Contribute Video Subtitles Third-party Tools

**Toolbox** 

Upload file

## ML: Dimensionality Reduction

#### **Contents** [hide] 1 Motivation I: Data Compression 2 Motivation II: Visualization 3 Principal Component Analysis Problem Formulation 4 Principal Component Analysis Algorithm 5 Choosing the Number of Principal Components 6 Reconstruction from Compressed Representation

## 7 Advice for Applying PCA

■ We may want to reduce the dimension of our features if we have a lot of redundant data.

■ To do this, we find two highly correlated features, plot them, and make a new line that seems to describe both features accurately. We place all the new features on this single line. Doing dimensionality reduction will reduce the total data we have to store in computer memory and will speed up our

learning algorithm. Note: in dimensionality reduction, we are reducing our features rather than our number of examples. Our variable m

will stay the same size; n, the number of features each example from  $x^{(1)}$  to  $x^{(m)}$  carries, will be reduced.

Motivation II: Visualization

It is not easy to visualize data that is more than three dimensions. We can reduce the dimensions of our data to 3 or less in order to plot it.

Example: hundreds of features related to a country's economic system may all be combined into one feature that you

call "Economic Activity."

## The most popular dimensionality reduction algorithm is Principal Component Analysis (PCA)

Principal Component Analysis Problem Formulation

**Problem formulation** 

#### Given two features, $x_1$ and $x_2$ , we want to find a single line that effectively describes both features at once. We then

**PCA** is not linear regression

map our old features onto this new line to get a new single feature. The same can be done with three features, where we map them to a plane.

Reduce from 2d to 1d: find a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to

minimize the projection error. The more general case is as follows:

so as to minimize the projection error.

If we are converting from 3d to 2d, we will project our data onto two directions (a plane), so k will be 2.

More generally, in linear regression we are taking all our examples in x and applying the parameters in  $\Theta$  to predict y.

■ In PCA, we are minimizing the **shortest distance**, or shortest *orthogonal* distances, to our data points.

In PCA, we are taking a number of features  $x_1, x_2, \ldots, x_n$ , and finding a closest common dataset among them. We aren't trying to predict any result and we aren't applying any theta weights to the features.

## Before we can apply PCA, there is a data pre-processing step we must perform:

**Data preprocessing** 

Preprocess (feature scaling/mean normalization):

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

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.

Above, we first subtract the mean of each feature from the original feature. Then we scale all the features (

We can define specifically what it means to reduce from 2d to 1d data as follows: 
$$x^{(i)} \in \mathbb{R}^2 \longrightarrow z^{(i)} \in \mathbb{R}$$

The z values are all real numbers and are the projections of our features onto  $u^{\left(1\right)}.$ 

So, PCA has two tasks: figure out  $u^{(1)},\ldots,u^{(k)}$  and also to find  $z_1,z_2,\ldots,z_m$  . The mathematical proof for the following procedure is complicated and beyond the scope of this course.

1. Compute "covariance matrix"

This can be vectorized in Octave as:

confusingly---they represent entirely different things).

Sigma = (1/m) \* X' \* X;

Note that  $x^{(i)}$  is an  $n \times 1$  vector,  $(x^{(i)})^T$  is an  $1 \times n$  vector and X is a  $m \times n$  matrix (row-wise stored examples). The product of those will be an n imes n matrix, which are the dimensions of  $\Sigma$ . 2. Compute "eigenvectors" of covariance matrix  $\Sigma$ 

[U,S,V] = svd(Sigma);

 $u^{(1)}, \ldots, u^{(n)}$ , which is exactly what we want. 3. Take the first k columns of the U matrix and compute z

We'll assign the first k columns of U to a variable called 'Ureduce'. This will be an n imes k matrix. We compute z with:  $z^{(i)} = Ureduce^T \cdot x^{(i)}$ 

dimensions  $k \times 1$ . To summarize, the whole algorithm in octave is roughly:

Sigma = (1/m) \* X' \* X; % compute the covariance matrix [U,S,V] = svd(Sigma); % compute our projected directions
Ureduce = U(:,1:k); % take the first k directions z = Ureduce' \* x; % compute the projected data points

 $Ureduce^T$  will have dimensions  $k \times n$  while  $x^{(i)}$  will have dimensions  $n \times 1$ . The product  $Ureduce^T \cdot x^{(i)}$  will have

#### How do we choose k, also called the *number of principal components*? Recall that k is the dimension we are reducing to. One way to choose k is by using the following formula:

Choose 
$$k$$
 to be the smallest value such that: 
$$\frac{\frac{1}{m}\sum_{i=1}^{m}||x^{(i)}-x_{approx}^{(i)}||^2}{\frac{1}{m}\sum_{i=1}^{m}||x^{(i)}||^2}\leq 0.01 \text{ (1\%)}$$
 In other words, the squared projection error divided by the total variation should be less than one percent, so that 99% of the variance is retained. Algorithm for choosing  $k$ 

3. Check the formula given above that 99% of the variance is retained. If not, go to step one and increase k. This procedure would actually be horribly inefficient. In Octave, we will call svd:

Which gives us a matrix S. We can actually check for 99% of retained variance using the S matrix as follows:

[U,S,V] = svd(Sigma)

1. Try PCA with  $k=1,2,\ldots$ 

Note that we can only get approximations of our original data.

Advice for Applying PCA The most common use of PCA is to speed up supervised learning.

training set. **Applications** 

validation or test sets. You can apply the mapping  $z^{(i)}$  to your cross-validation and test sets after it is defined on the

Given a training set with a large number of features (e.g.  $x^{(1)},\dots,x^{(m)}\in\mathbb{R}^{10000}$ ) we can use PCA to reduce the

Reduce space of data Speed up algorithm Visualization of data

our results y. Using just regularization will be at least as effective. Don't assume you need to do PCA. Try your full machine learning algorithm without PCA first. Then use PCA if

# Motivation I: Data Compression

# We need to find new features, $z_1, z_2$ (and perhaps $z_3$ ) that can effectively **summarize** all the other features.

The goal of PCA is to reduce the average of all the distances of every feature to the projection line. This is the projection error.

Reduce from n-dimension to k-dimension: Find k vectors  $u^{(1)}, u^{(2)}, \ldots, u^{(k)}$  onto which to project the data

In linear regression, we are minimizing the squared error from every point to our predictor line. These are vertical distances.

Principal Component Analysis Algorithm

## Given training set: $x^{(1)}, x^{(2)}, \ldots, x^{(m)}$

Replace each 
$$x_j^{(i)}$$
 with  $x_j^{(i)} - \mu_j$  If different features on different s

 $x^{(i)} \in \mathbb{R}^2 \to z^{(i)} \in \mathbb{R}$ 

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

svd() is the 'singular value decomposition', a built-in Octave function. What we actually want out of svd() is the 'U' matrix of the Sigma covariance matrix:  $U \in \mathbb{R}^{n \times n}$ . U contains

Choosing the Number of Principal Components

# Given the average squared projection error: $rac{1}{m}\sum_{i=1}^{m}||x^{(i)}-x_{approx}^{(i)}||^2$ Also given the total variation in the data: $rac{1}{m}\sum_{i=1}^m ||x^{(i)}||^2$

$$rac{1}{m}\sum_{i=1}^m ||x|^i$$
ner words, the squared projection error divided by the total variation

2. Compute 
$$U_{reduce},z,x$$
3. Check the formula given above that 99%

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

To go from 1-dimension back to 2d we do:  $z \in \mathbb{R} o x \in \mathbb{R}^2$  . We can do this with the equation:  $x_{approx}^{(1)} = U_{reduce} \cdot z^{(1)}$  .

number of features in each example of the training set (e.g.  $z^{(1)},\ldots,z^{(m)}\in\mathbb{R}^{1000}$ ). Note that we should define the PCA reduction from  $x^{(i)}$  to  $z^{(i)}$  only on the training set and not on the cross-

Compressions

Bad use of PCA: trying to prevent overfitting. We might think that reducing the features with PCA would be an effective way to address overfitting. It might work, but is not recommended because it does not consider the values of

• Choose k = 2 or k = 3

you find that you need it. Next: Anomaly Detection Back to Index: Main

Category: ML:Lecture Notes