

Machine Learning Basics: Unsupervised Learning Algorithms

Sargur N. Srihari
srihari@cedar.buffalo.edu

This is part of lecture slides on [Deep Learning](http://www.cedar.buffalo.edu/~srihari/CSE676):
<http://www.cedar.buffalo.edu/~srihari/CSE676>

Topics in Machine Learning Basics

1. Learning Algorithms
2. Capacity, Overfitting and Underfitting
3. Hyperparameters and Validation Sets
4. Estimators, Bias and Variance
5. Maximum Likelihood Estimation
6. Bayesian Statistics
7. Supervised Learning Algorithms
8. Unsupervised Learning Algorithms
9. Stochastic Gradient Descent
10. Building a Machine Learning Algorithm
11. Challenges Motivating Deep Learning

Topics in Unsupervised Learning

0. Unsupervised learning
1. Principal Components Analysis
2. K-means clustering

Unsupervised Algorithms

- They are algorithms that experience only “features” but not a supervision signal
- Refers to extract information from a distribution with no human labor to annotate examples
- Term refers to:
 - Density estimation
 - Learning to draw samples from a distribution
 - Learning to denoise data from a distribution
 - Finding a manifold that the data lies near
 - Clustering the data into groups of related examples

Finding the “best” representation

- A classic unsupervised learning task
- Generally, ‘best’ is looking for a representation:
 - that preserves as much information about x as possible while obeying some penalty or constraint
 - aimed at keeping the representation *simpler* or more accessible than x itself
- There are multiple ways of defining a *simpler* representation:
 1. Lower dimensional representations
 2. Sparse representations
 3. Independent representations

Three “simpler” representations

1. Low-dimensional representations

- Attempt to compress as much information about x as possible in a smaller representation

2. Sparse representations

- Embed the dataset into a representation whose entries are mostly zero for most inputs

3. Independent representations

- Attempt to *disentangle* sources of variation underlying data distribution so that the dimensions of the distribution are statistically independent

- Three criteria are not mutually exclusive

- Low dimensions yield weaker dependencies

Representation: central to deep learning

- Some simple examples of representation learning algorithms are:
 1. Principal components analysis
 2. k -means clustering

Principal Components Analysis

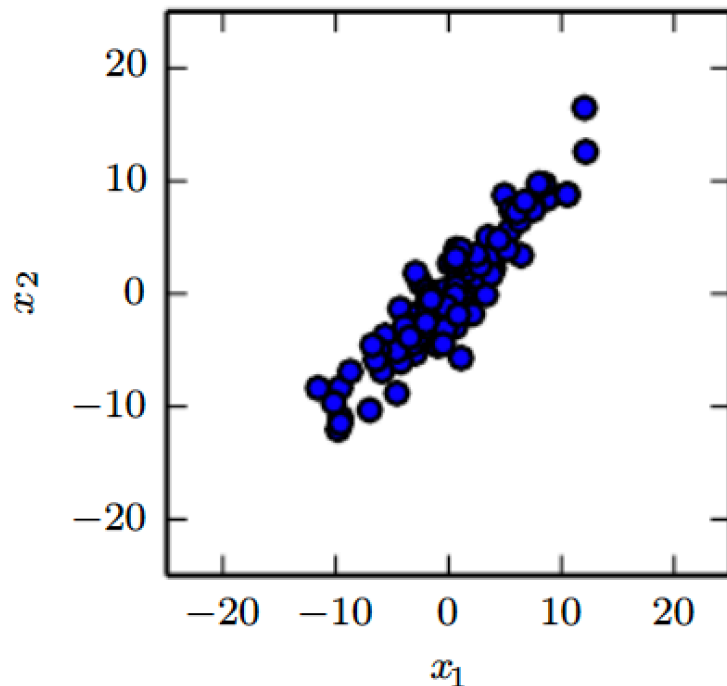
- PCA is useful for compressing data
- Can also view PCA as unsupervised learning to find a representation satisfying two conditions:
 - Lower dimensionality than original input
 - Independence (a representation with elements having no linear correlation with each other)
- This is the first step towards a representation with statistically independent elements
 - To achieve full independence, a representation learning algorithm must also remove nonlinear relationships between variables

PCA learns a linear projection

- PCA projects an input \mathbf{x} to a representation \mathbf{z}

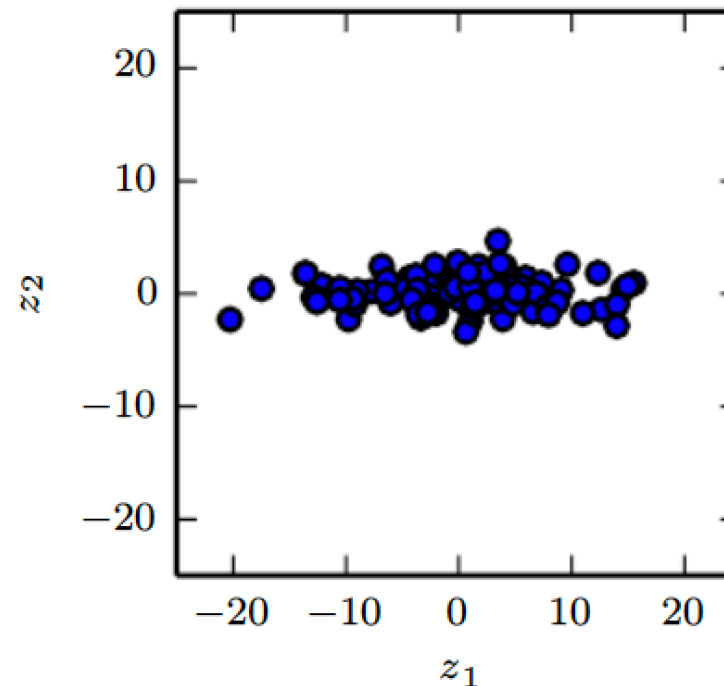
Original data
with samples in \mathbf{x}

Variance occurs along directions
not axis-aligned



Transformed data
with samples in $\mathbf{z} = \mathbf{x}^T \mathbf{W}$

Varies most along axis z_1
Direction of second most variance is z_2



Uses of PCA

- We could learn a one-dimensional representation that best reconstructs the original data
 - It corresponds to the first principal component
- PCA is a simple and effective dimensionality reduction method
 - That preserves as much of the data as possible
 - Measured by least squares reconstruction error
- In the following we study how PCA decorrelates the original data representation \mathbf{X}

How PCA uncorrelates the data

- Consider $m \times n$ dimensional design matrix \mathbf{X}
 - We assume data has mean zero, $E[\mathbf{x}] = 0$
 - Otherwise, subtract the mean from all examples
 - The unbiased sample covariance matrix associated with \mathbf{X} is $\text{Var}[\mathbf{x}] = \frac{1}{m-1} \mathbf{X}^\top \mathbf{X}$
- PCA finds a representation (through a linear transformation) $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$ where $\text{Var}[\mathbf{z}]$ is diagonal

Using SVD: $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top$

$$\begin{aligned} \text{Var}[\mathbf{x}] &= \frac{1}{m-1} \mathbf{X}^\top \mathbf{X} \\ &= \frac{1}{m-1} (\mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top)^\top \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top \\ &= \frac{1}{m-1} \mathbf{W} \mathbf{\Sigma}^\top \mathbf{U}^\top \mathbf{U} \mathbf{\Sigma} \mathbf{W}^\top \\ &= \frac{1}{m-1} \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^\top, \end{aligned}$$

$$\begin{aligned} \text{Var}[\mathbf{z}] &= \frac{1}{m-1} \mathbf{Z}^\top \mathbf{Z} \\ &= \frac{1}{m-1} \mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W} \\ &= \frac{1}{m-1} \mathbf{W}^\top \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^\top \mathbf{W} \\ &= \frac{1}{m-1} \mathbf{\Sigma}^2, \end{aligned}$$

The representation has a diagonal covariance matrix given by $\mathbf{\Sigma}^2$ which implies that the elements of \mathbf{z} are mutually uncorrelated

PCA as a disentangler

- Very important property of PCA:
 - The ability to transform data into a representation where the elements are mutually uncorrelated
 - A simple example of a representation that attempts to disentangle unknown factors of variation underlying the data
 - PCA rotates input space (by \mathbf{W}) that aligns the principal axes of variance with the basis of the new representation space associated with z
- We are also interested in more complicated forms of feature dependencies than what can be done with a simple linear transformation

K-means clustering and One-hot

- Divides the training set into k clusters
 - A cluster has examples near to each other
- An input x is represented by a *one-hot vector* h
 - If x is in cluster i then $h_i=1$ and all others are zero
- One-hot vector code is a sparse representation
 - All but one entry is zero for every input
 - One-hot : Lose benefits of distributed representation but has statistical and computational advantages
 - All examples similar to each other
 - Entire representation captured by single integer

K-means algorithm

- Initialize k centroids $\{\mu^{(1)}, \dots, \mu^{(k)}\}$ to different values
- Alternate between two steps until convergence
 - Step 1: assign sample to cluster i where i is the index of the nearest centroid $\mu^{(i)}$
 - Step 2: Each centroid $\mu^{(i)}$ is updated to the mean of all samples assigned to cluster i

One-hot vs Distributed representation

- A difficulty pertaining to clustering:
 - Suppose we run two clustering algorithms on a data set of red cars, red trucks, gray cars, grey trucks
 - Algo 1: red vehicles together, and grey vehicles together
 - Algo 2: cars together, and trucks together
 - Clustering along one feature may be relevant and irrelevant along another
 - Third Algo obtains four clusters capturing both features but loses similarity
- These illustrate as to why distributed representations may be better than one-hot
 - It contains information about both attributes