

Outline

- 1 Motivation
- 2 k -Means Clustering
- 3 Principal Component Analysis
- 4 Wrap-Up

Outline

- 1 Motivation
- 2 k -Means Clustering
- 3 Principal Component Analysis
- 4 Wrap-Up

Recap: Supervised vs. Unsupervised Learning

Supervised learning

- ▶ Machine learning task of inferring a function from **labeled training data**
- ▶ Training data includes both the input and the desired results
→ correct results (target values) are given

Unsupervised learning

- ▶ Methods try to find hidden structure in **unlabeled data**
- ▶ The model is not provided with the correct results during the training
- ▶ No error or reward signal to evaluate a potential solution
- ▶ Examples:
 - ▶ **Clustering** (e. g. by k -means algorithm)
→ group into classes only on the basis of their statistical properties
 - ▶ **Dimensionality reduction** (e. g. by principal component analysis)
 - ▶ Hidden Markov models with unsupervised learning

Unsupervised Learning

Objective

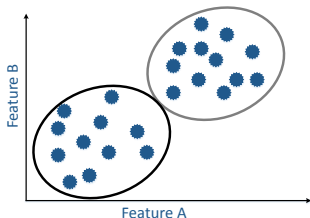
- ▶ Find interesting insights in data
- ▶ Key metrics can be relationships, main characteristics or similarity of data points
- ▶ Usually of **exploratory** nature as there are no labels

Pros and cons

- ▶ Often **easy to get unlabeled data**
 - Labels can be expensive when manual annotations are needed
- ▶ Highly **subjective** as a standardized goal is missing

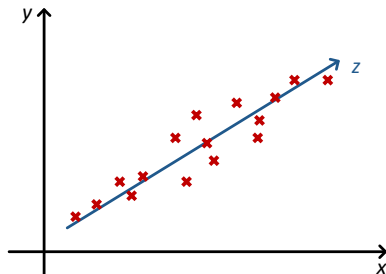
Clustering vs. Dimensionality Reduction

Clustering



- Identifies **subgroups** of data points with **homogeneous characteristics**

Dimensionality reduction



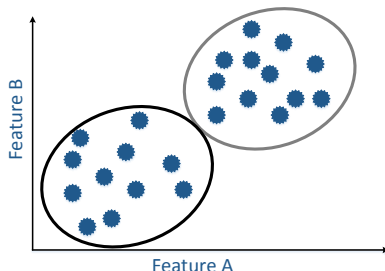
- Calculates the **main dimensions** across that data points are distributed
- Transforms data

Outline

- 1 Motivation
- 2 k -Means Clustering**
- 3 Principal Component Analysis
- 4 Wrap-Up

k -Means Clustering

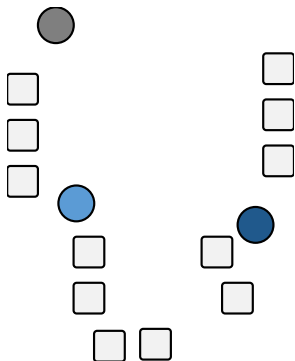
- Partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype for the cluster



- Computationally expensive; instead, we use **efficient heuristics**
- Default: Euclidean distance as metric and variance as a measure of cluster scatter

Lloyd's Algorithm: Outline

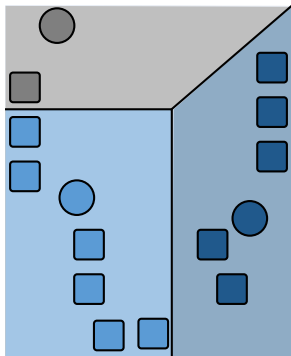
- 1 Randomly generated k initial "means" (here: $k = 3$)



- 2 Create k clusters by associating every observation with the nearest mean (colored partitions)
- 3 Centroid of each of the k clusters becomes the new mean
- 4 Repeat steps 2 and 3 until convergence

Lloyd's Algorithm: Outline

- 1 Randomly generated k initial "means" (here: $k = 3$)
- 2 Create k clusters by associating every observation with the nearest mean (colored partitions)

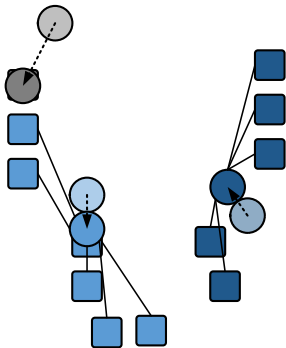


- 3 Centroid of each of the k clusters becomes the new mean

- 4 Repeat steps 2 and 3 until convergence

Lloyd's Algorithm: Outline

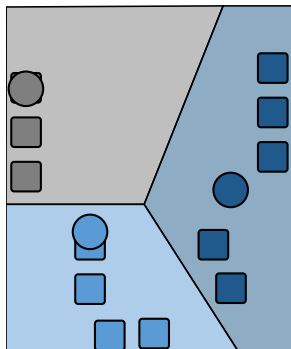
- 1 Randomly generated k initial "means" (here: $k = 3$)
- 2 Create k clusters by associating every observation with the nearest mean (colored partitions)
- 3 Centroid of each of the k clusters becomes the new mean



4 Repeat steps 2 and 3 until convergence

Lloyd's Algorithm: Outline

- 1** Randomly generated k initial "means" (here: $k = 3$)
- 2** Create k clusters by associating every observation with the nearest mean (colored partitions)
- 3** Centroid of each of the k clusters becomes the new mean
- 4** Repeat steps 2 and 3 until convergence



Lloyd's Algorithm: Pseudocode

1 Initialization

Choose a set of k means $\mathbf{m}_1^{(1)}, \dots, \mathbf{m}_k^{(1)}$ randomly

2 Assignment Step

Assign each observation to the cluster whose mean is closest to it, i.e.

$$S_i^{(t)} = \{\mathbf{x}_p : \|\mathbf{x}_p - \mathbf{m}_i^{(t)}\| \leq \|\mathbf{x}_p - \mathbf{m}_j^{(t)}\| \forall 1 \leq j \leq k\}$$

where each observation is assigned to exactly one cluster, even if it could be assigned to two or more of them

3 Update Step

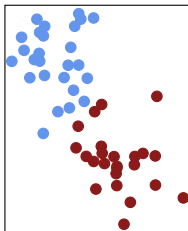
Calculate the new means to be the centroids of the observations in the new clusters

$$\mathbf{m}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x}_j \in S_i^{(t)}} \mathbf{x}_j$$

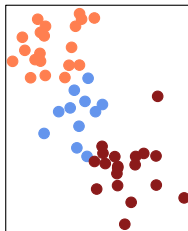
Optimal Choice of k

Example: Plots show the results of applying k -means clustering with different values of k

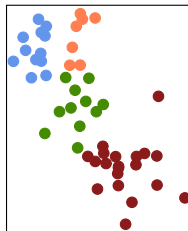
$k=2$



$k=3$



$k=4$



Note: Final results can vary according to random initial means!

→ In practice, k -means clustering will be performed using **multiple random assignments** and only the **best result is reported**

Optimal Choice of k

- ▶ **Optimal choice of k** searches for a balance between maximum compression ($k = 1$) and maximum accuracy ($k = n$)
- ▶ **Diagnostic checks** to determine the number of clusters, such as
 - 1 Simple rule of thumb sets $k \approx \sqrt{n/2}$
 - 2 Elbow Method: Plot percent of explained variance vs. number of clusters
 - 3 Usage of information criteria
 - 4 ...
- ▶ k -means minimizes the **within-cluster sum of squares (WCSS)**

$$\arg \min_S \sum_{i=1}^k \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

with clusters $S = \{S_1, \dots, S_k\}$ and mean points $\boldsymbol{\mu}_i$ in S_i

Clustering

Research Question

Group countries based on income, literacy, infant mortality and life expectancy (file: `countries.csv`) into three groups accounting for developed, emerging and undeveloped countries.

```
# Use first column as row names for each observation
countries <- read.csv("countries.csv", header=TRUE, sep=";", row.names=1)
head(countries)
```

	Per.capita.income	Literacy	Infant.mortality	Life.expectancy
Brazil	10326	90.0	23.60	75.4
Germany	39650	99.0	4.08	79.4
Mozambique	830	38.7	95.90	42.1
Australia	43163	99.0	4.57	81.2
China	5300	90.9	23.00	73.0
Argentina	13308	97.2	13.40	75.3

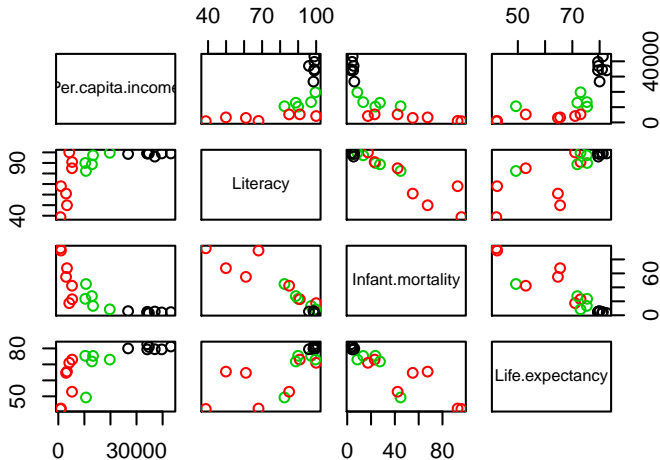
Clustering

```
km <- kmeans(countries, 3, nstart=10)
km

## K-means clustering with 3 clusters of sizes 7, 7, 5
##
## Cluster means:
##   Per.capita.income Literacy Infant.mortality Life.expectancy
## 1          35642.143    98.50          4.477143          80.42857
## 2          3267.286    70.50          56.251429          58.80000
## 3          13370.400    91.58          23.560000          68.96000
##
## Clustering vector:
##           Brazil          Germany      Mozambique      Australia          China
##           3              1              2              1              2
##   Argentina United Kingdom  South Africa      Zambia      Namibia
##           3              1              3              2              2
##           Georgia      Pakistan          India      Turkey      Sweden
##           2              2              2              3              1
##           Lithuania      Greece          Italy      Japan
##           3              1              1              1
##
## Within cluster sum of squares by cluster:
## [1] 158883600 20109876 57626083
## (between_SS / total_SS =  94.1 %)
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"
## [5] "tot.withinss" "betweenss"    "size"         "iter"
## [9] "ifault"
```


Visualizing Results of Clustering

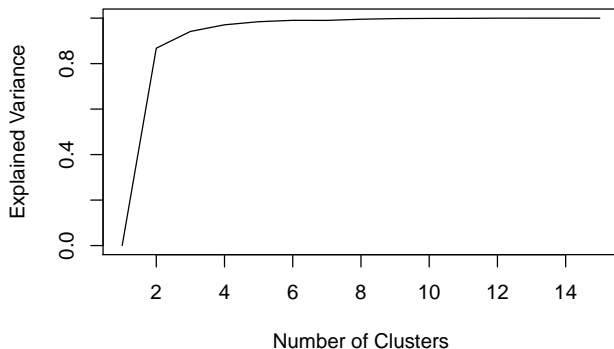
```
plot(countries, col = km$cluster)
```



Elbow Plot to Choose k

Choose k (here: $k = 3$) so that adding another cluster doesn't result in much better modeling of the data

```
ev <- c()
for (i in 1:15) {
  km <- kmeans(countries, i, nstart=10)
  ev[i] <- sum(km$betweenss)/km$totss
}
plot(1:15, ev, type="l", xlab="Number of Clusters", ylab="Explained Variance")
```



Outline

- 1 Motivation
- 2 k -Means Clustering
- 3 Principal Component Analysis**
- 4 Wrap-Up

Principal Component Analysis

Motivation

- ▶ Large datasets with **many variables** require extensive computing power
- ▶ However, only a small number of variables usually is informative
- ▶ High-dimensional data (≥ 4 dimensions) can be **difficult to visualize**

Principal component analysis (PCA)

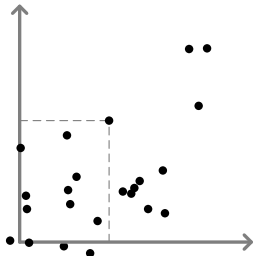
- ▶ Finds a **low-dimensional representation** of data
- ▶ Reduces n -dimensional data to k -dimensions with $k \leq n$
- ▶ Goal: keep as much of the informative value as possible

Principal Component Analysis

Intuition

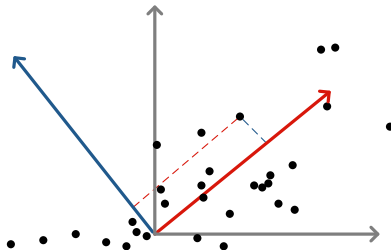
Standard basis

$$\mathbf{x} = (0.3, 0.5)^T$$



Rotated basis

$$\mathbf{z} = (0.7, 0.1)^T$$

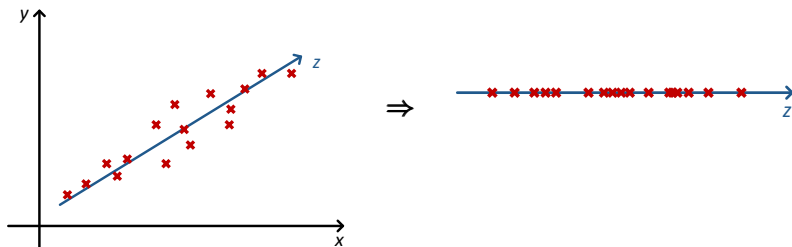


- ▶ **First principal component** is the direction with the largest variance
- ▶ **Second principal component** is orthogonal and in the direction of the largest remaining variance

Principal Component Analysis

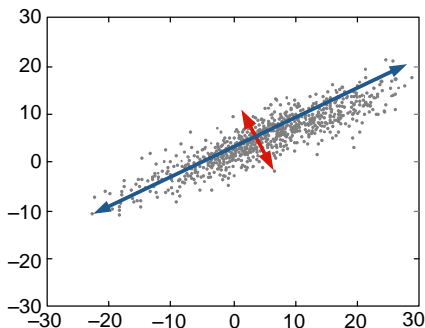
Use cases

- ▶ Principal components can work as **input for supervised learning**
→ especially suited for algorithms with super-linear time complexity in the number of dimensions
- ▶ PCA can **visualize** high-dimensional data with simple graph



Principal Component Analysis

- ▶ **Linear combination** of **uncorrelated** variables with maximal variance
→ high variance signals high information content
- ▶ Data is projected onto **orthogonal** component vectors so that the projection error is minimized
- ▶ Order of directions gives the i -th **principal component**



Standardizing

- ▶ Scaling changes results of PCA → **standardizing** is recommend
- ▶ Center variable around mean $\mu = 0$ with standard deviation $\sigma = 1$

Steps

- 1** Calculate mean and standard deviation for $\mathbf{x} = [x_1, \dots, x_N]^T$

$$\mu = \frac{1}{N-1} \sum_{i=1}^N x_i \quad \sigma = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)^2}$$

Note: R uses internally denominator $N - 1$ instead of N

- 2** **Transform** variable (built-in via `scale(x)` in R)

$$x_i \leftarrow \frac{x_i - \mu}{\sigma}$$

```
x <- scale(1:10)
c(mean(x), sd(x))
```

```
## [1] 0 1
```


Algorithm

- ▶ PCA maps \mathbf{x}_i onto a new basis via a **linear combination**

$$\mathbf{z}_i = \phi_{1,i} \mathbf{x}_1 + \phi_{2,i} \mathbf{x}_2 + \dots + \phi_{n,i} \mathbf{x}_n$$

with normalization $\sum_{j=1}^n \phi_{j,i}^2 = 1$

- ▶ \mathbf{z}_i is the i -th **principal component**
- ▶ $\phi_{1,i}, \dots, \phi_{n,i}$ are the **loadings** of the i -th principal component
- ▶ In matrix notation, this gives

$$\mathbf{Z} = \Phi \mathbf{X}$$

- ▶ Geometrically, Φ is a **rotation** with stretching
→ it also spans the **directions** of the principal components

Algorithm

- ▶ If \mathbf{x}_i is standardized, it has mean zero and also \mathbf{z}_i
- ▶ Hence, the **variance** of \mathbf{z}_i is

$$\frac{1}{N} \sum_{j=1}^N z_{j,i}^2$$

- ▶ First loading vector searches a **direction to maximize the variance**

$$\max_{\phi_{j,1}} \frac{1}{N} \sum_{j=1}^N z_{j,i}^2 = \max_{\phi_{j,1}} \frac{1}{N} \sum_{i=1}^N \left[\sum_{j=1}^n \phi_{j,1} x_{i,j} \right]^2 \quad \text{subject to} \quad \sum_{j=1}^n \phi_{j,1}^2 = 1$$

- ▶ Numerically solved via a **singular value decomposition**

Singular Value Decomposition

Covariance matrix

- ▶ **Covariance matrix** Σ for the standardized data is given by

$$\Sigma = \frac{1}{N} X^T X \quad \Leftrightarrow \quad \Sigma_{ij} = \frac{1}{N} \mathbf{x}_i^T \mathbf{x}_j$$

- ▶ $\Sigma \in \mathbb{R}^{N \times N}$ is symmetric with diagonals being the variance
- ▶ Goal: high variance but orthogonality, i. e. zero off-diagonal elements

Singular value decomposition

- ▶ **Singular value decomposition** of square matrix X gives

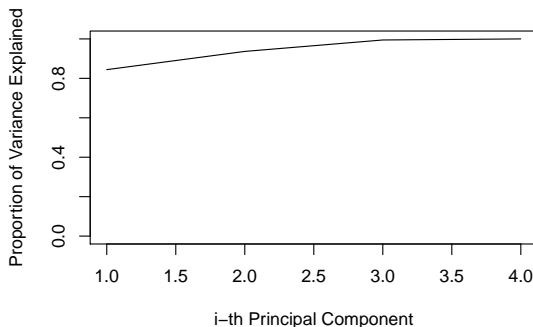
$$X = V \Sigma V^{-1}$$

- ▶ V is a matrix with the **eigenvectors** of X ($\Rightarrow VV^T = I_N$)
- ▶ Σ is a **diagonal** matrix with the corresponding **eigenvalues**
- ▶ Then $\Phi = V$

Proportion of Variance Explained

- Plot with **cumulative proportion** of variance explained

```
pve <- pca$sdev^2 / sum(pca$sdev^2)
plot(cumsum(pve), xlab="i-th Principal Component",
     ylab="Proportion of Variance Explained",
     type="l", ylim=c(0, 1))
```

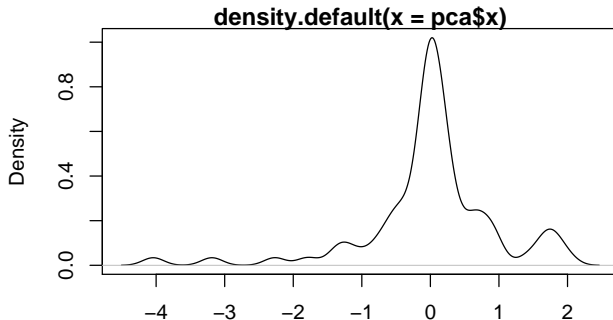


→ First principal component explains more than 80 % of the variance

PCA Example

- **Density estimation** reveals subgroups in one dimension

```
plot(density(pca$x))
```



N = 76 Bandwidth = 0.1545

→ One also observes three groups: a peak, as well as a tail and a leading group

Outline

- 1 Motivation
- 2 k -Means Clustering
- 3 Principal Component Analysis
- 4 Wrap-Up**

Summary

- ▶ Unsupervised learning usually provides **explanatory** insights
- ▶ **k-means clustering** identifies subsets of similar points
- ▶ **Elbow plot** determines a suitable number of clusters k
- ▶ **PCA reduces dimensions** with a minimal amount of information loss

Commands in R

<code>kmeans(d, k, nstart=n)</code>	<i>k</i> -means clustering
<code>prcomp(d, scale=TRUE)</code>	PCA with scaling
<code>cumsum(x)</code>	Cumulative sums
<code>apply(d, f)</code>	Apply function f to all data points in d

How about deep learning?

Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

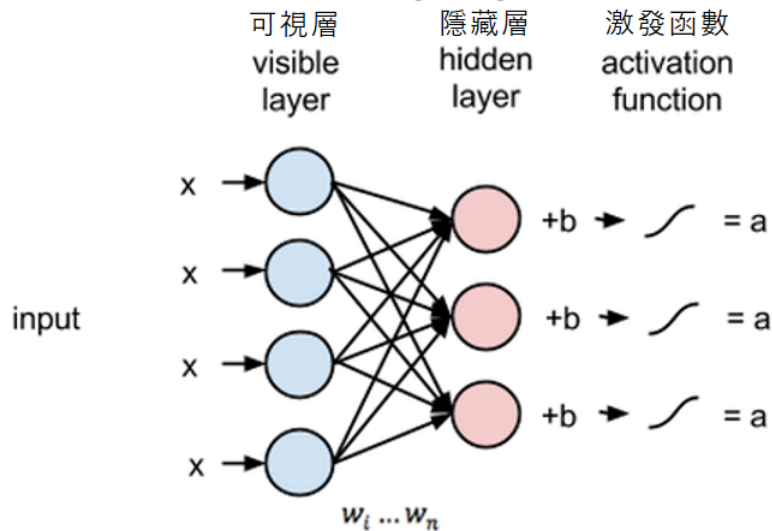
Dimensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which

finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer “encoder” network

Restricted Boltzmann Machines

前向傳導

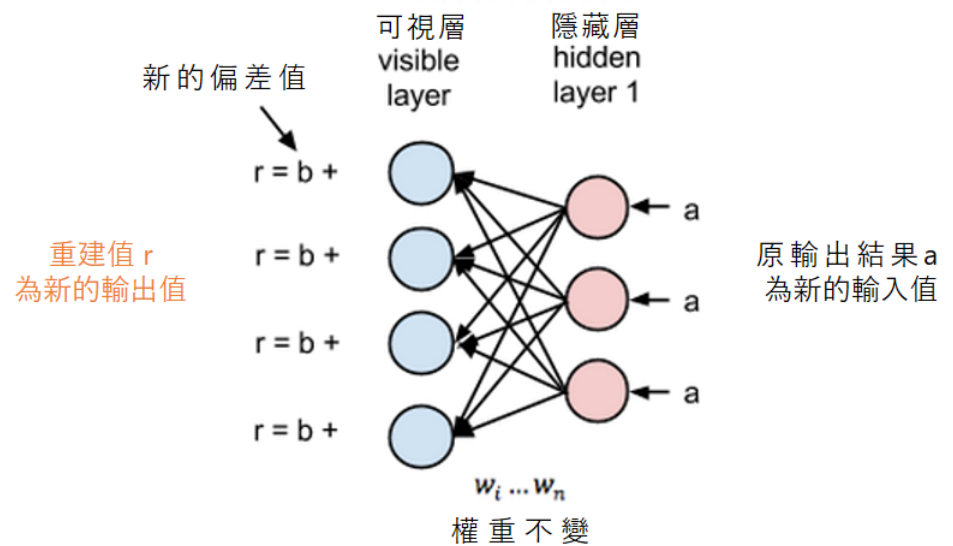
Multiple Inputs



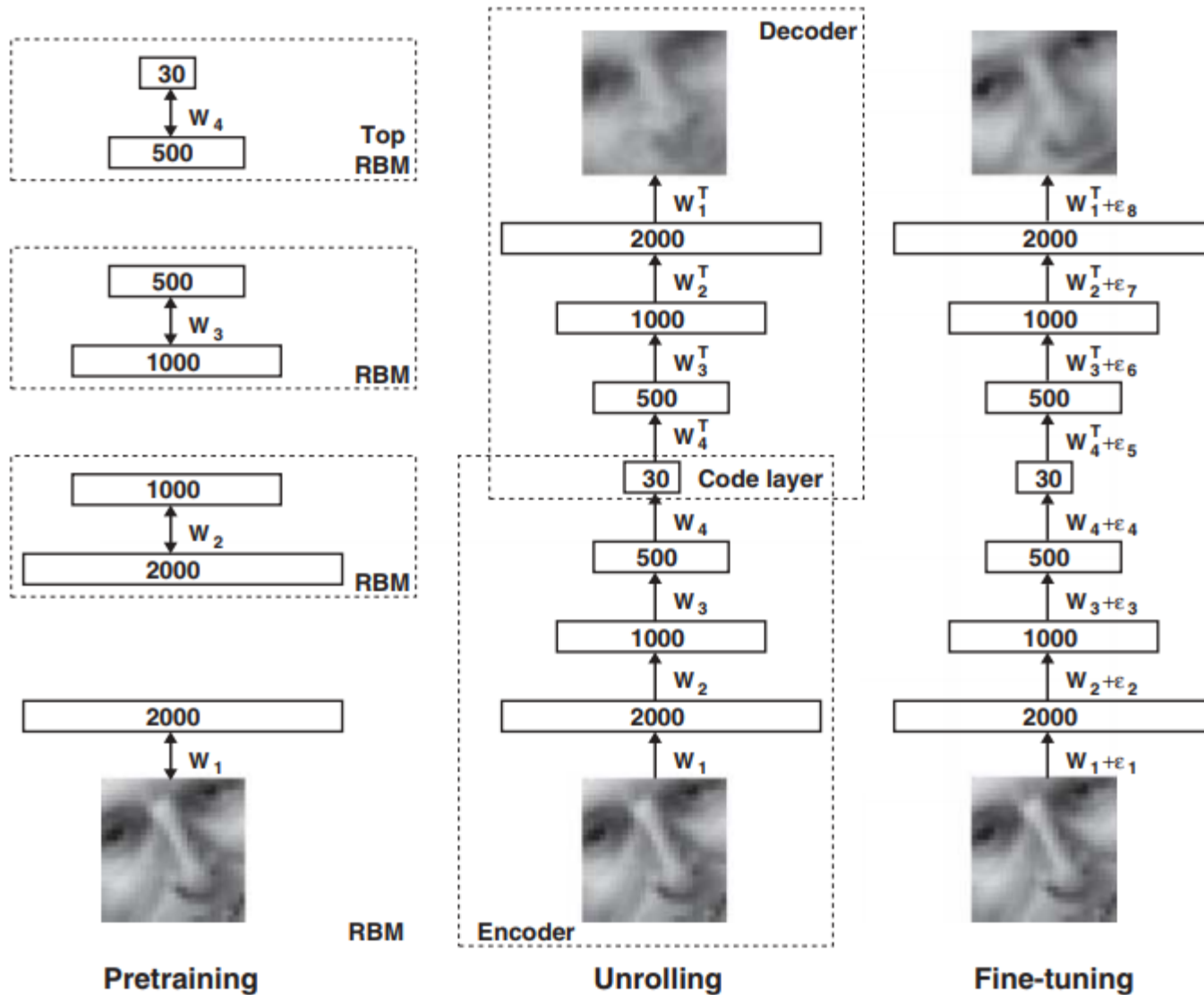
激發函數 $f(\text{權重 } w * \text{輸入值 } x) + \text{偏差值 } b) = \text{輸出結果 } a$

重建

Reconstruction

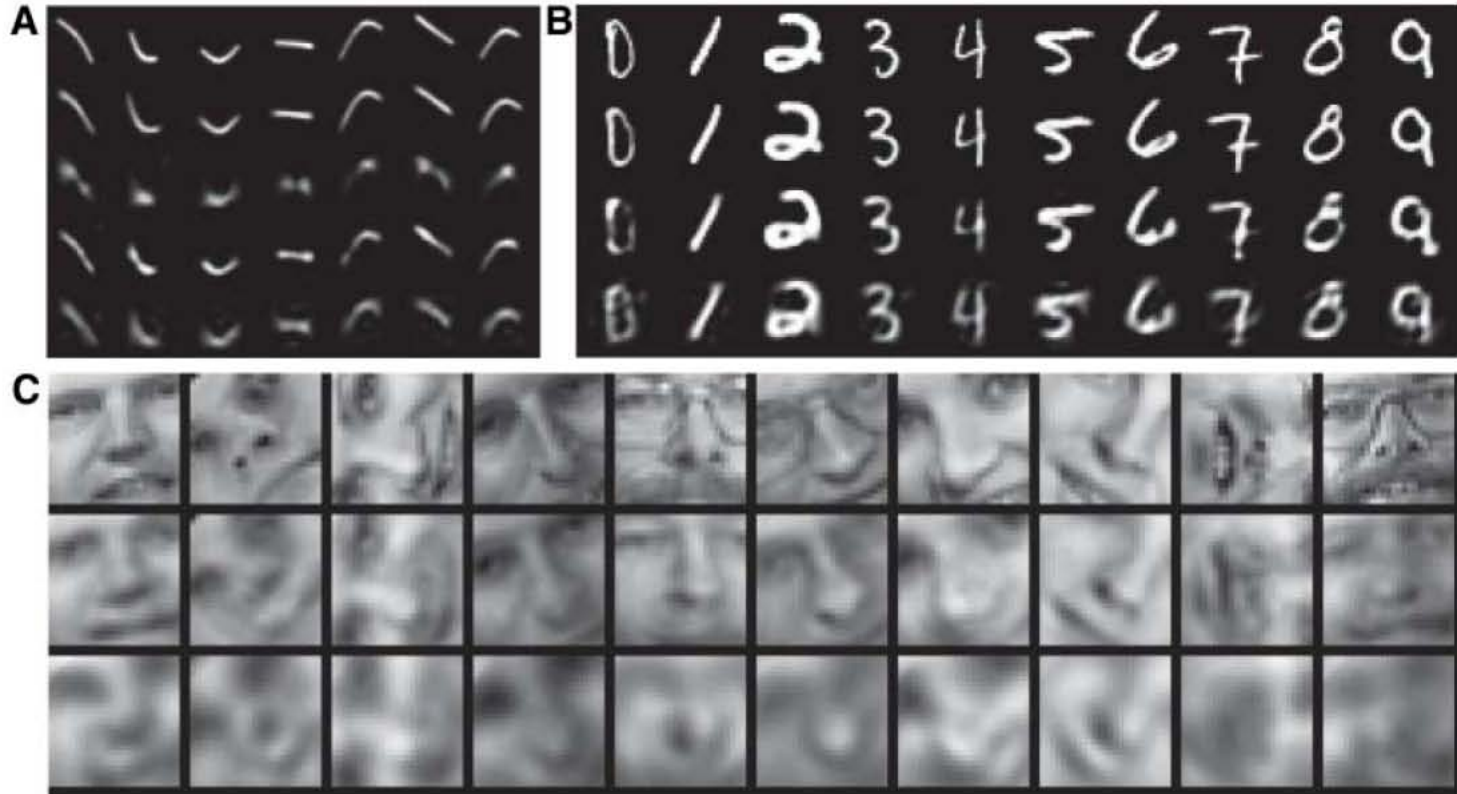


Autoencoder Network



Comparison with PCA

Fig. 2. (A) Top to bottom: Random samples of curves from the test data set; reconstructions produced by the six-dimensional deep autoencoder; reconstructions by “logistic PCA” (8) using six components; reconstructions by logistic PCA and standard PCA using 18 components. The average squared error per image for the last four rows is 1.44, 7.64, 2.45, 5.90. (B) Top to bottom: A random test image from each class; reconstructions by the 30-dimensional autoencoder; reconstructions by 30-dimensional logistic PCA and standard PCA. The average squared errors for the last three rows are 3.00, 8.01, and 13.87. (C) Top to bottom: Random samples from the test data set; reconstructions by the 30-dimensional autoencoder; reconstructions by 30-dimensional PCA. The average squared errors are 126 and 135.



Comparison with PCA

Fig. 3. (A) The two-dimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoencoder. For an alternative visualization, see (8).

