

# Document Classification

$$P(\mathbf{x}|y) = N! \prod_{w \in D} \frac{p(w|y)^{h(w,D)}}{h(w,D)!}$$

- $N$  = number of words in the document  $\mathbf{x}$
- $p(w|y)$  = probability that the word  $w$  occurs in class  $y$
- $h(w,D)$  = total occurrences of the word  $w$  in document  $D$

# Classifying “burger” sentiment

- Let us assume that we would like to classify whether the following document is positive (1) or negative (-1) in sentiment.
- $D = \text{“the burger i ate was an awesome burger”}$
- Further assume that we encounter these words in positive and negative classes as follows. To make our computations easier let us assume that we removed “the”, “i”, “was”, “an” as stop words.

word	+1	-1
burger	3	2
ate	3	2
awesome	4	1

# Computing class conditional probabilities

word	+1	-1	$p(w +1)$	$p(w -1)$
burger	3	2	$3/10$	$2/5$
ate	3	2	$3/10$	$2/5$
awesome	4	1	$4/10$	$1/5$

# Quiz

- Using the probabilities in the table in slide 21 and multinomial naive Bayes formula in slide 19, compute  $P(D|+1)$  and  $P(D|-1)$

$$P(D|+1) = p(\text{burger}|+1)^2 p(\text{ate}|+1)^1 p(\text{awesome}|+1)^1 \frac{4!}{2!1!1!}$$

$$= \left(\frac{3}{10}\right)^2 \left(\frac{3}{10}\right)^1 \left(\frac{4}{10}\right)^1 \times \frac{24}{2}$$

$$= (0.3)^3 \times 0.4 \times 12 \quad \text{--- ①}$$

$$P(D|-1) = p(\text{burger}|-1)^2 p(\text{ate}|-1)^1 p(\text{awesome}|-1)^1 \frac{4!}{2!1!1!}$$

$$= \left(\frac{2}{5}\right)^2 \left(\frac{2}{5}\right)^1 \left(\frac{1}{5}\right)^1 \times 12$$

$$= (0.4)^3 \times 0.2 \times 12 \quad \text{--- ②}$$

# Quiz

- Assuming that both positive and negative classes have equal probability
  - $p(+1) = p(-1) = 0.5$
- Compute  $P(+1|D)$  and  $P(-1|D)$  for the document in slide 22.
- Is this document D positive or negative?

$$p(+1|D) = \frac{p(D|+1)p(+1)}{p(D)}$$

$$p(D) = p(\text{burger})^2 p(\text{ate})^1 p(\text{awesome})^1 \frac{4!}{2!1!1!}$$

$$= \left(\frac{5}{15}\right)^2 \left(\frac{5}{15}\right)^1 \left(\frac{5}{15}\right)^1 12$$

$$\therefore p(+1|D) = \frac{(0.3)^3 \times 0.4 \times 0.5}{(1/2)^4}$$

$$p(-1|D) = \frac{p(D|-1)p(-1)}{p(D)} = \frac{(0.4)^3 \times 0.2 \times 0.5}{(1/3)^4}$$

$$\frac{p(+1|D)}{p(-1|D)} = \left(\frac{0.3}{0.4}\right)^3 \frac{0.4}{0.2} = \frac{81}{32} > 1 \quad \therefore \underline{\underline{D \text{ is positive}}}$$

# Dimensionality Reduction

COMP 527 Data Mining  
Danushka Bollegala



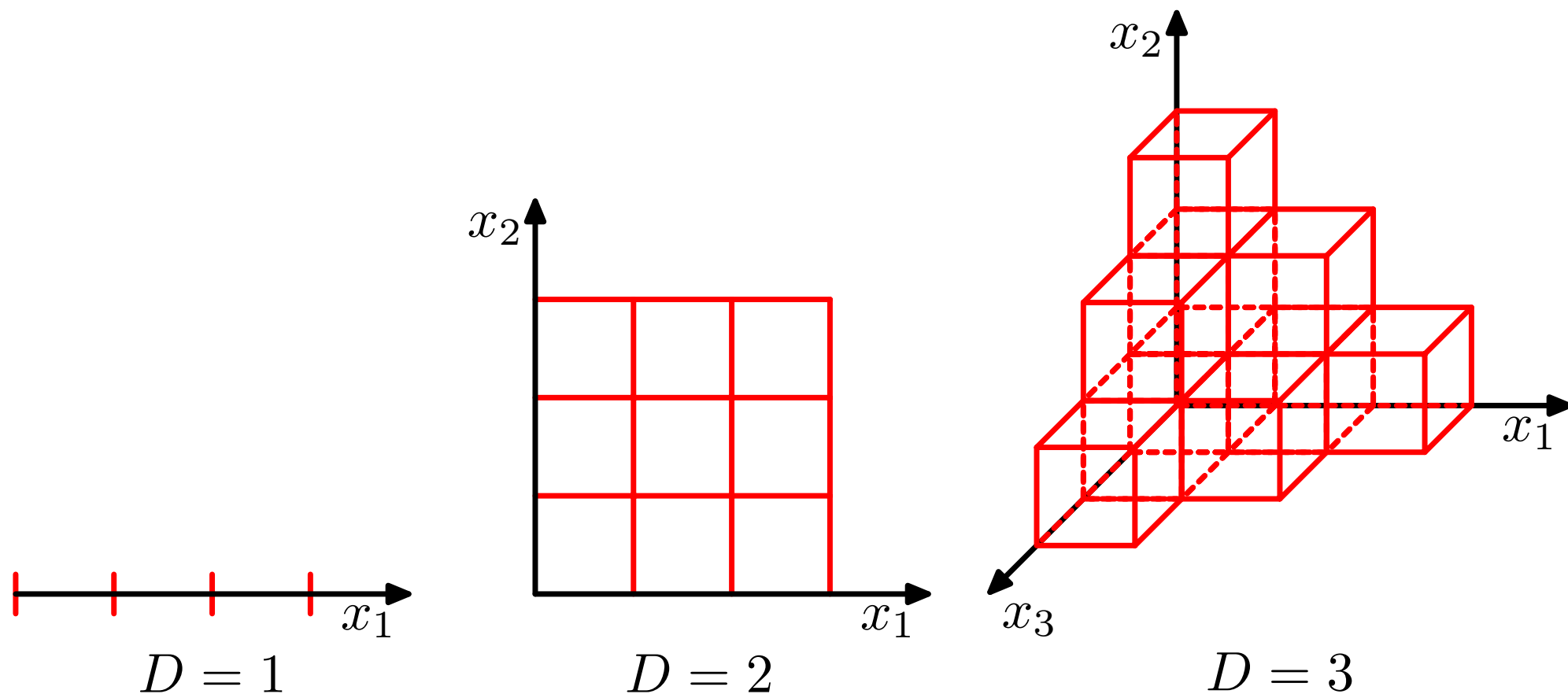
UNIVERSITY OF  
LIVERPOOL

# Outline

- Problems with high dimensional data
- Dimensionality Reduction Methods
  - Singular Value Decomposition (SVD)
  - Principal Component Analysis (PCA)

# Problems in High Dimensions

- Curse of dimensionality
- We need exponentially large number of data points to cover a high dimensional space



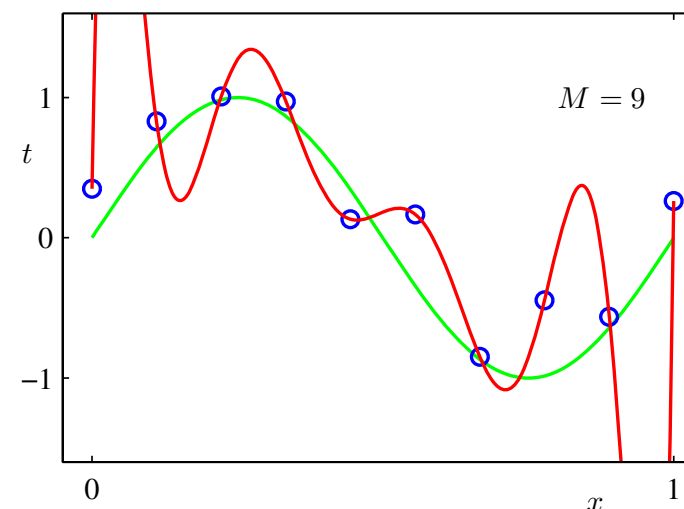
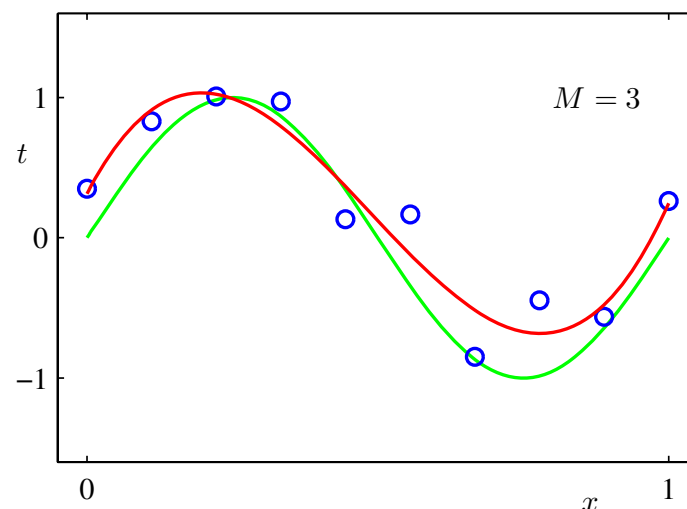
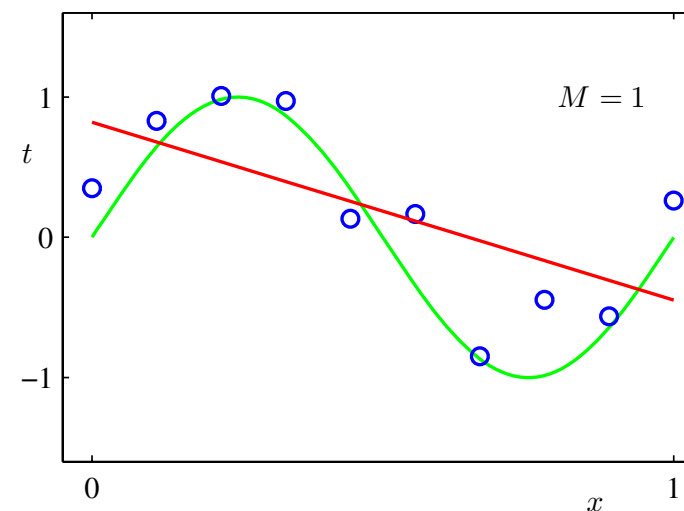
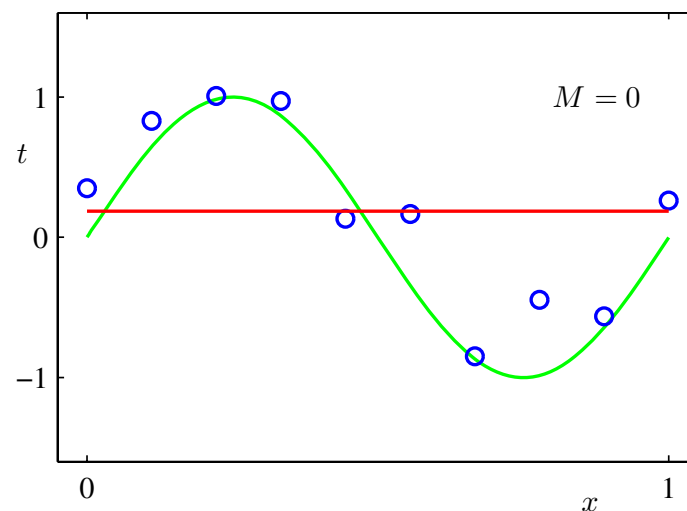


# Problems in High Dimensions

- Data Sparseness
  - Although we have a large feature space (lots of dimensions to the data), we only observe a small number of non-zero features in any instance
  - This was the case for texts (in particular with the bag-of-words model)

# Problems in High Dimensions

- Overfitting
  - Given  $n$  train data points, we can come up with an  $n$ -dimensional ( $n$ -th order) polynomial that passes through all those data points.
  - But it is very unlikely that it will fit well for the test data points



# Problems in High Dimensions

- Time consuming (time complexity is large)
- Consider computing cosine similarity between two  $n$  dimensional vectors, when  $n$  increases.
- Memory issues (space complexity is large)
  - Storing high dimensional dense vectors can be problematic when
    - the dimensionality of the vectors is large
    - there are lots of vectors (instances) to store

# Solution

- Dimensionality Reduction
  - Try to project the original vectors to a lower dimensional space  $L$
- What constraints do we have
  - Try to minimise the error due to the projection
    - If  $X$  and  $Y$  are neighbours in the original space, then they must also be neighbours in the projected space
  - Try to retain salient/important/principal dimensions as much as possible and remove the non-salient/unimportant/auxiliary dimensions as much as possible

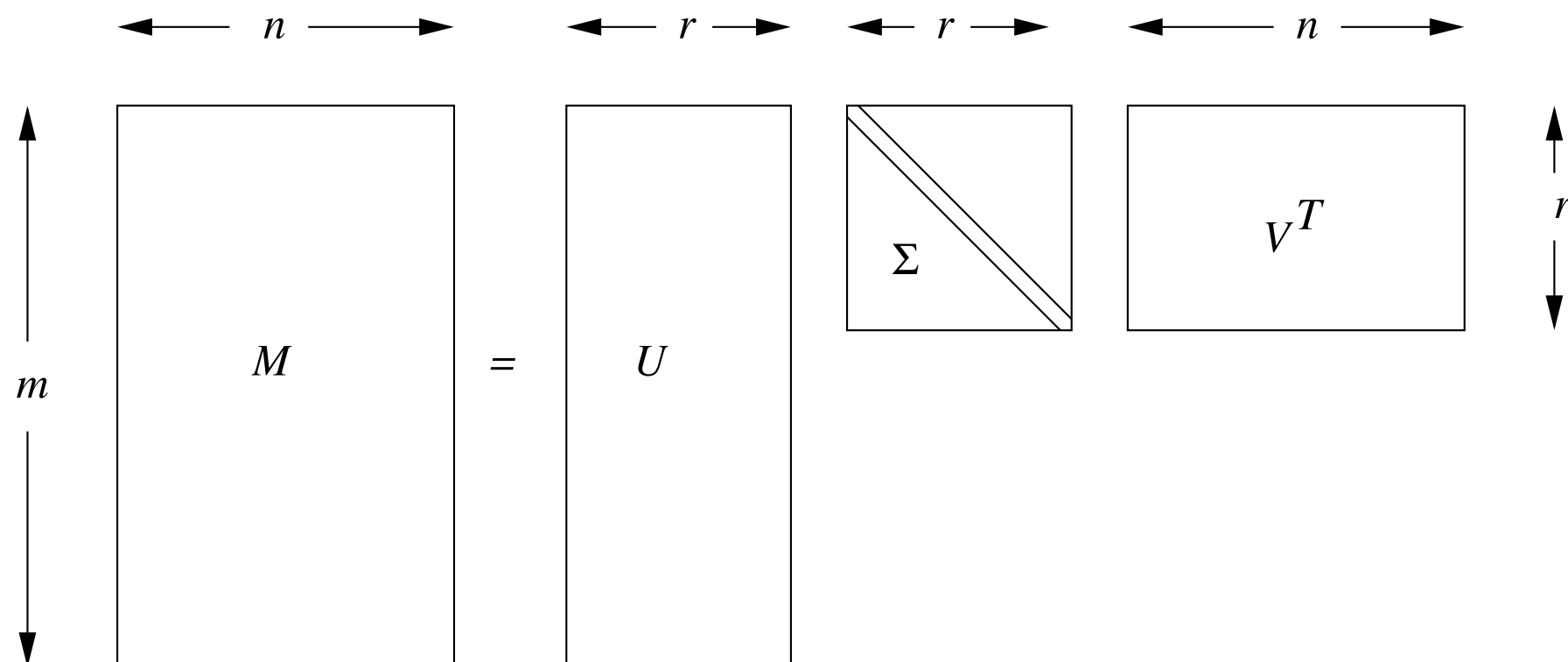
# Eigenvalue Decomposition

- Linear Algebra Revision
  - Eigenvalues and Eigenvectors of a Square Matrix
  - $Ax = \lambda x$ 
    - $x$  is the eigenvector of  $A$  corresponding to the eigenvalue  $\lambda$
- Compute the eigenvalues and eigenvectors of the following matrix

$$\begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}$$

# Singular Value Decomposition

- Eigenvalue decomposition can be performed only for square matrices.
- Singular Value Decomposition (SVD) is a operation that can be applied to any matrix
  - $M = U\Sigma V^T$
  - $U$  and  $V$  are unitary (perpendicular) matrices,  $\Sigma$  is a diagonal matrix (singular values of  $M$  are diagonal elements of  $\Sigma$ ).
  - Columns of  $U$  and  $V$  are perpendicular.  $U^T U = I$  and  $V^T V = I$ .



# SVD?



**S**nayperskaya **V**intovka sistem'y **D**ragunova obraz'tsa  
(Dragunov Sniper Rifle or SVD)

# Example

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} .14 & 0 \\ .42 & 0 \\ .56 & 0 \\ .70 & 0 \\ 0 & .60 \\ 0 & .75 \\ 0 & .30 \end{bmatrix} \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \begin{bmatrix} .58 & .58 & .58 & 0 & 0 \\ 0 & 0 & 0 & .71 & .71 \end{bmatrix}$$

$M \qquad \qquad U \qquad \qquad \Sigma \qquad \qquad V^T$

To perform SVD in python (scipy) use `scipy.linalg.svd`  
<http://docs.scipy.org/doc/scipy/reference/generated/scipy.linalg.svd.html>



# Dimensionality Reduction with SVD

- Procedure
  - Perform SVD on  $M$ . Retain the top- $k$  (largest) singular values in  $\Sigma$  and set the remainder to zero.
  - Let us denote the diagonal matrix produced by the previous step by  $\Sigma_k$
  - The  $k$ -dimensional approximation (projection)  $M_k$  of  $M$  is then given by
    - $M_k = U \Sigma_k V^T$

# Reason

- The  $k$ -dimensional matrix  $M_k$  that minimises the Frobenius norm  $\|M - M_k\|$  is given by the matrix  $M_k$  computed as described in the previous slide
- Frobenius norm
  - Extension of the vector L2 norm to matrices
  - Frobenius norm of a matrix  $M$  is given by  $\sqrt{\sum_{ij} M_{ij}^2}$

# Proof

- By performing SVD on M, let
  - $M = PQR^T$

$$m_{ij} = \sum_k \sum_\ell p_{ik} q_{k\ell} r_{\ell j} \quad \|M\|^2 = \sum_i \sum_j (m_{ij})^2 = \sum_i \sum_j \left( \sum_k \sum_\ell p_{ik} q_{k\ell} r_{\ell j} \right)^2$$

$$\left( \sum_k \sum_\ell p_{ik} q_{k\ell} r_{\ell j} \right)^2 = \sum_k \sum_\ell \sum_m \sum_n p_{ik} q_{k\ell} r_{\ell j} p_{in} q_{nm} r_{mj}$$

$$\|M\|^2 = \sum_i \sum_j \sum_k \sum_\ell \sum_n \sum_m p_{ik} q_{k\ell} r_{\ell j} p_{in} q_{nm} r_{mj}$$

$$\|M\|^2 = \sum_i \sum_j \sum_k \sum_n p_{ik} q_{kk} r_{kj} p_{in} q_{nn} r_{nj}$$

$$\|M\|^2 = \sum_j \sum_k q_{kk} r_{kj} q_{kk} r_{kj}$$

$$\|M\|^2 = \sum_k (q_{kk})^2$$

Much easier proof exists if you use the trace of a matrix and its properties

# Proof using Trace

$$\mathbf{M} = \mathbf{P}\mathbf{Q}\mathbf{R}^\top$$

$$||\mathbf{M}||_2^2 = \text{tr}(\mathbf{M}^\top \mathbf{M})$$

$$||\mathbf{M}||_2^2 = \text{tr}(\mathbf{R}\mathbf{Q}\mathbf{P}^\top \mathbf{P}\mathbf{Q}\mathbf{R}^\top)$$

$$||\mathbf{M}||_2^2 = \text{tr}(\mathbf{R}\mathbf{Q}\mathbf{Q}\mathbf{R}^\top)$$

$$||\mathbf{M}||_2^2 = \text{tr}(\mathbf{R}\mathbf{Q}^2\mathbf{R}^\top)$$

$$||\mathbf{M}||_2^2 = \text{tr}(\mathbf{R}^\top \mathbf{R}\mathbf{Q}^2)$$

$$||\mathbf{M}||_2^2 = \text{tr}(\mathbf{Q}^2)$$

# SVD and Approximation Error

$$\mathbf{M} = \mathbf{P} \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_k & & \\ & & & \sigma_{k+1} & \\ & & & & \ddots \\ & & & & & \sigma_n \end{bmatrix} \mathbf{R}^\top$$

If  $\mathbf{M}_k$  is the matrix with  $(k+1)$  and above singular values set to zero  
( $k$ -th rank approximation)

$$\mathbf{M}_k = \mathbf{P} \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_k & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{bmatrix} \mathbf{R}^\top$$

# SVD and Approximation error

- Then, the approximation error,  $\|\mathbf{M} - \mathbf{M}_k\|^2$  becomes

$$\|\mathbf{M} - \mathbf{M}_k\|^2 = \text{tr} \left( \mathbf{P} \begin{bmatrix} 0 & & & & \\ & \ddots & & & \\ & & 0 & & \\ & & & \sigma_{k+1} & \\ & & & & \ddots \\ & & & & & \sigma_n \end{bmatrix} \mathbf{R}^\top \right)$$

$$\|\mathbf{M} - \mathbf{M}_k\|^2 = \sigma_{k+1}^2 + \dots + \sigma_n^2$$

If we want to minimize this error, then we must select the largest singular values for the first  $1 \dots k$  positions!

# Applications of SVD

- Latent Semantic Analysis
  - words vs. document matrix
  - Find similar words (query expansion)
  - Find similar documents (similarity search)
- Recommendation Systems
  - users vs. items/products
  - Recommend similar products to users

# Two uses of SVD

- SVD for dimensionality reduction
  - Compute  $M = U\Sigma V^T$
  - Get the largest  $k$  singular values from  $\Sigma$  to construct a diagonal matrix  $\Sigma_k$
  - Get the corresponding left singular vectors from  $U$  to construct a matrix  $U_k$
  - Reduce the number of columns of  $M$  to  $k$  to construct the matrix  $M_k$ 
    - $M_k = U_k \Sigma_k$



# Two uses of SVD

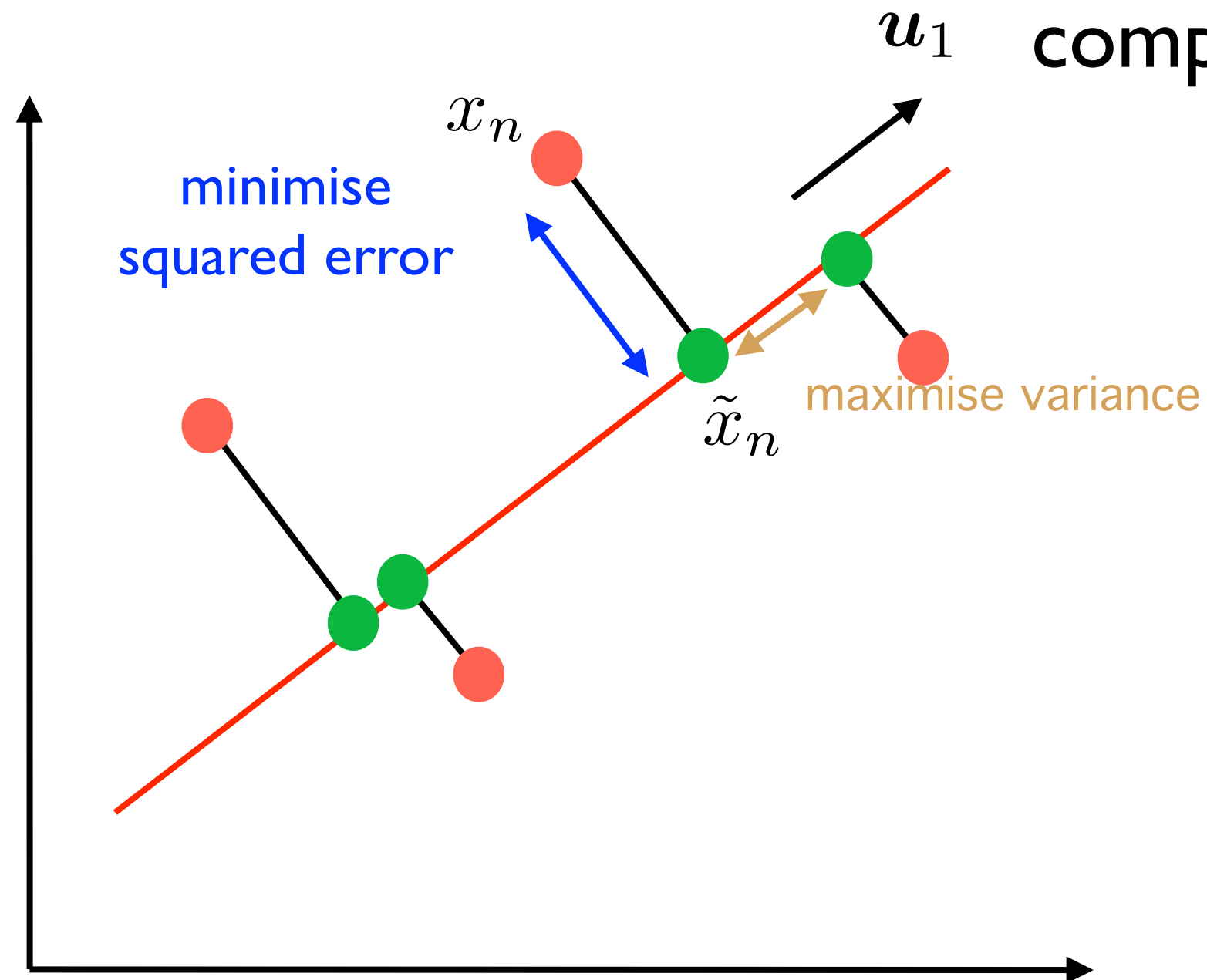
- SVD to increase the density of a matrix
  - Compute  $M = U\Sigma V^T$
  - Get the largest  $k$  singular values from  $\Sigma$  to construct a diagonal matrix  $\Sigma_k$
  - Get the corresponding left singular vectors from  $U$  to construct a matrix  $U_k$
  - Get the corresponding right singular vectors from  $U$  to construct a matrix  $V_k$
  - Reproduce a dense version of  $M$ ,  $M_k$ 
    - $M_k = U_k \Sigma_k V_k^T$
    - Lesser number of non-zero values in  $M_k$
    - However, we end up with negative values in  $M_k$  even though  $M$  is a matrix with all non-negative values!

# Principal Component Analysis

- We would like to **project** our **high dimensional** data points to a **low dimensional** space by **preserving the geometric properties** in the original space as much as possible
- Two ways to do this:
  - **Maximise the variance of the projected data**
  - Linear projection that minimises the average projection cost (e.g. sum of squared Euclidean distance between original and projected points)
- PCA is also known as the Karhunen-Loève Transform

# Idea

First principal component



# Maximum Variance Formulation

- Problem
  - Given  $D$  dimensional  $N$  data points  $\{x_n\}$ , where  $n=1, \dots, N$ , we must project those into an  $M < D$  dimensional space
  - $M$  is given
  - Let us consider the case  $M=1$  (one-dimensional projection)
  - The projection direction is given by the unit vector  $u_1$ 
    - $u_1^T u_1 = 1$
- $\tilde{x}_n = u_1^T x_n$

# Maximum Variance Formulation

- Mean of the data points

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- Variance of the projected data

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{u}_1^\top \mathbf{x}_n - \mathbf{u}_1^\top \bar{\mathbf{x}})^2 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

\*An unbiased estimator of variance uses (N-1) instead of N.

But this does not matter because we are only interested in the maximisation of the variance.

- $\mathbf{S}$  is the covariance matrix given by

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top$$

# Maximum Variance Formulation

- We must maximize the variance subjected to the normalization constraint on  $\mathbf{u}_1$

Lagrange multiplier method

$$L(\mathbf{u}_1, \lambda_1) = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$$

$$\frac{\partial L(\mathbf{u}_1, \lambda_1)}{\partial \mathbf{u}_1} = 0 \quad \Rightarrow \quad \mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

$$\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 = \lambda_1$$

This is variance!

$\mathbf{u}_1$  is the eigenvector of  $\mathbf{S}$  that corresponds to the largest eigenvalue of  $\mathbf{S}$

# PCA Algorithm

- INPUT

- D dimensional N data points  $\{x_n\}$ , where  $n=1, \dots, N$

- Dimensionality M

- Procedure

- Compute the covariance matrix S for the dataset

- Compute the first M eigenvectors of S

- return the computed eigenvectors

# A Word on Complexity

- Eigenvalue decomposition of a  $D \times D$  matrix is  $O(D^3)$
- However, we only need the largest  $M$  eigenvectors of  $S$
- This can be computed efficiently using truncated methods such as the power-iteration method in  $O(MD^2)$
- Reference
  - Golub & Van Loan, Matrix Computations, John Hopkins University Press, 1996.



# Computational Remarks (1/2)

- The following methods for computing PCA are equivalent (i.e. gives the same principal components)
  1. Compute the sum of squared Euclidean distance between original and projected points and **minimise** it.
  2. Compute the pairwise squared Euclidean distance between projected points and **maximise** it.
  3. Compute the variance of the projected points on the projection line and **maximise** it.

# Computational Remarks (2/2)

- Let us show that (2) and (3) are equivalent

$$\begin{aligned}(3) &= \sum_{i=1}^N (\mathbf{x}_i \mathbf{u} - \bar{\mathbf{x}}^\top \mathbf{u})^2 \\&= \sum_{i=1}^N \left( \left( \mathbf{x}_i - \frac{1}{N} \sum_{j=1}^N \mathbf{x}_j \right)^\top \mathbf{u} \right)^2 \\&= \sum_{i=1}^N \sum_{j=1}^N ((\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{u})^2 \\&= \sum_{i=1}^N \sum_{j=1}^N (\mathbf{x}_i^\top \mathbf{u} - \mathbf{x}_j^\top \mathbf{u})^2 \\&= (2)\end{aligned}$$

Using (2) for computing PCA is usually a bad idea when the dataset is large because the pairwise combinations grow quadratically with the number of data points in the dataset. But could be helpful for small examples because we do not have to compute the mean.

# References

- Mining of Massive Datasets
  - <http://infolab.stanford.edu/~ullman/mmds.html#original>
  - Chapter 11 on Dimensionality Reduction (SVD)
- Pattern Recognition and Machine Learning
  - Section 1.1 (overfitting)
  - Section 1.4 (curse of dimensionality)
  - Page 561 onwards: PCA