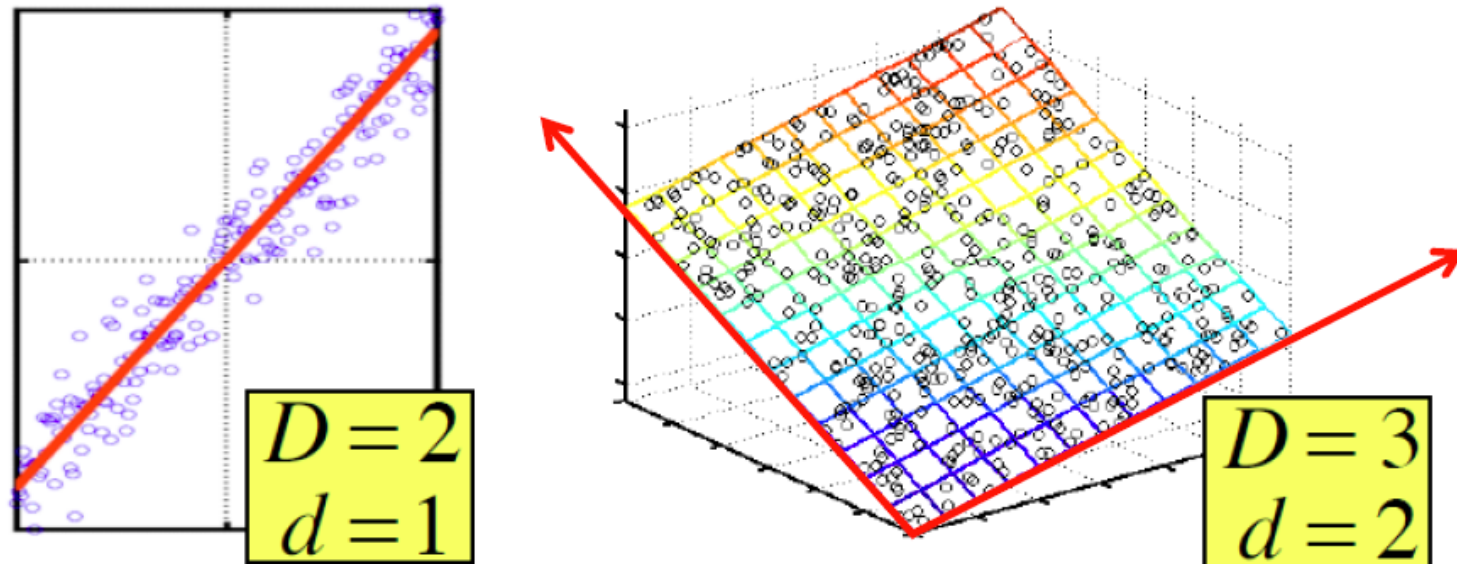


# Dimensionality Reduction: Introduction

Mining of Massive Datasets  
Leskovec, Rajaraman, and Ullman  
Stanford University



# Dimensionality Reduction



- **Assumption:** Data lies on or near a low  $d$ -dimensional subspace
- **Axes of this subspace are effective representation of the data**

# Dimensionality Reduction

- **Compress / reduce dimensionality:**
  - $10^6$  rows;  $10^3$  columns; no updates
  - Random access to any cell(s); **small error: OK**

	day	We	Th	Fr	Sa	Su
customer		7/10/96	7/11/96	7/12/96	7/13/96	7/14/96
ABC Inc.		1	1	1	0	0
DEF Ltd.		2	2	2	0	0
GHI Inc.		1	1	1	0	0
KLM Co.		5	5	5	0	0
Smith		0	0	0	2	2
Johnson		0	0	0	3	3
Thompson		0	0	0	1	1

The above matrix is really “2-dimensional.” All rows can be reconstructed by scaling  $[1 \ 1 \ 1 \ 0 \ 0]$  or  $[0 \ 0 \ 0 \ 1 \ 1]$

# Rank of a Matrix

- **Q:** What is **rank** of a matrix **A**?
- **A:** Number of **linearly independent** columns of **A**
- **For example:**
  - Matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$  has rank **r=2**
    - **Why?** The first two rows are linearly independent, so the rank is at least 2, but all three rows are linearly dependent (the first is equal to the sum of the second and third) so the rank must be less than 3.
- **Why do we care about low rank?**
  - We can write **A** as two “basis” vectors:  $[1 \ 2 \ 1] \ [-2 \ -3 \ 1]$
  - And new coordinates of :  $[1 \ 0] \ [0 \ 1] \ [1 \ 1]$

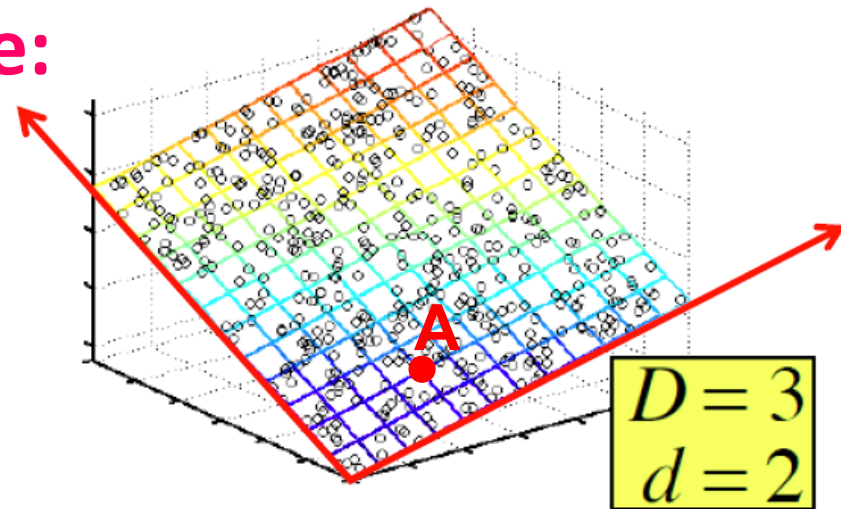
# Rank is “Dimensionality”

- **Cloud of points 3D space:**

- Think of point positions

as a matrix:  $\begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$  **A**  
**B**  
**C**

1 row per point:

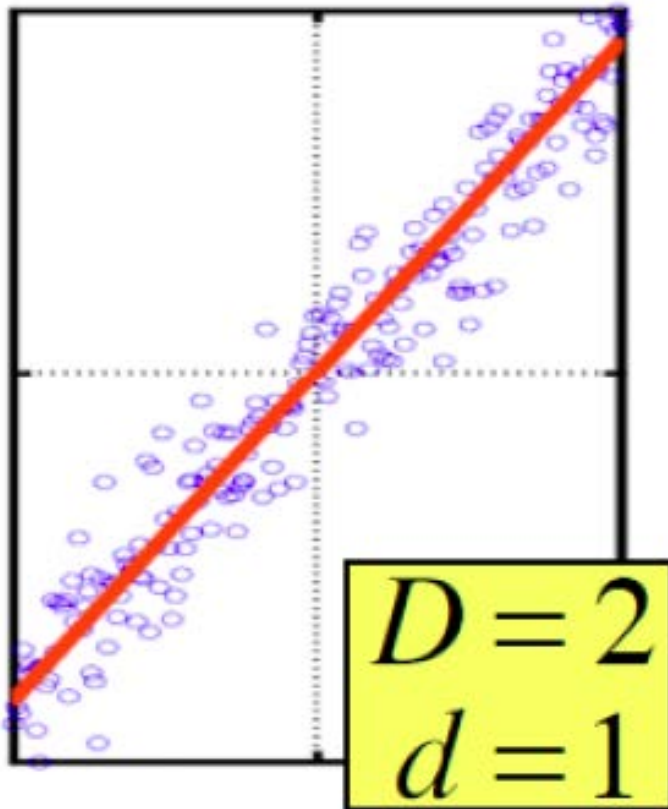


- **We can rewrite coordinates more efficiently!**

- Old coordinate system:  $[1 \ 0 \ 0] \ [0 \ 1 \ 0] \ [0 \ 0 \ 1]$
- **New coordinate system:  $[1 \ 2 \ 1] \ [-2 \ -3 \ 1]$**
- Then **A** has new coordinates:  $[1 \ 0]$ . **B**:  $[0 \ 1]$ , **C**:  $[1 \ 1]$ 
  - **Notice: We reduced the number of coordinates!**

# Dimensionality Reduction

- Goal of dimensionality reduction is to discover the axis of data!



Rather than representing every point with 2 coordinates we represent each point with 1 coordinate (corresponding to the position of the point on the red line).

By doing this we incur a bit of **error** as the points do not exactly lie on the line

# Why Reduce Dimensions?

## Why reduce dimensions?

- Discover hidden correlations/topics
  - Words that occur commonly together
- Remove redundant and noisy features
  - Not all words are useful
- Interpretation and visualization
- Easier storage and processing of the data

