#### **Dimensionality reduction**

#### Outline

- Dimensionality Reductions or data projections
- Random projections
- Singular Value Decomposition and Principal Component Analysis (PCA)

#### The curse of dimensionality

 The efficiency of many algorithms depends on the number of dimensions d

 Distance/similarity computations are at least linear to the number of dimensions

Index structures fail as the dimensionality of the data increases

#### Goals

- Reduce dimensionality of the data
- Maintain the meaningfulness of the data

#### Dimensionality reduction

- Dataset X consisting of n points in a ddimensional space
- Data point x<sub>i</sub> ∈ R<sup>d</sup> (d-dimensional real vector):
  - $x_i = [x_{i1}, x_{i2}, ..., x_{id}]$
- Dimensionality reduction methods:
  - Feature selection: choose a subset of the features
  - Feature extraction: create new features by combining new ones

#### Dimensionality reduction

- Dimensionality reduction methods:
  - Feature selection: choose a subset of the features
  - Feature extraction: create new features by combining new ones
- Both methods map vector x<sub>i</sub> ∈ R<sup>d</sup>, to vector y<sub>i</sub> ∈ R<sup>k</sup>, (k < <d)</li>

•  $F: R^d \rightarrow R^k$ 

## Linear dimensionality reduction

- Function F is a linear projection
- $\mathbf{y}_i = \mathbf{x}_i \mathbf{A}$

•  $\mathbf{Y} = \mathbf{X} \mathbf{A}$ 

• Goal: Y is as close to X as possible

#### Closeness: Pairwise distances

Johnson-Lindenstrauss lemma: Given
 ε>0, and an integer n, let k be a positive
 integer such that k≥k₀=O(ε<sup>-2</sup> logn). For
 every set X of n points in R<sup>d</sup> there exists
 F: R<sup>d</sup>→R<sup>k</sup> such that for all x<sub>i</sub>, x<sub>i</sub> ∈X

 $(1-\epsilon)||x_i - x_j||^2 \le ||F(x_i) - F(x_j)||^2 \le (1+\epsilon)||x_i - x_j||^2$ 

### What is the intuitive interpretation of this statement?

#### JL Lemma: Intuition

- Vectors x<sub>i</sub> ∈ R<sup>d</sup>, are projected onto a kdimensional space (k < <d): y<sub>i</sub> = x<sub>i</sub> A
- If ||x<sub>i</sub>||=1 for all i, then,

 $||\mathbf{x}_i - \mathbf{x}_j||^2$  is approximated by  $(\mathbf{d}/\mathbf{k})||\mathbf{y}_i - \mathbf{y}_j||^2$ 

#### • Intuition:

- The expected squared norm of a projection of a unit vector onto a random subspace through the origin is k/d
- The probability that it deviates from expectation is very small

#### Finding random projections

- Vectors x<sub>i</sub> ∈ R<sup>d</sup>, are projected onto a kdimensional space (k < < d)</li>
- Random projections can be represented by linear transformation matrix A
- $\mathbf{y}_i = \mathbf{x}_i \mathbf{A}$

• What is the matrix A?

#### Finding random projections

- Vectors x<sub>i</sub> ∈ R<sup>d</sup>, are projected onto a kdimensional space (k < < d)</li>
- Random projections can be represented by linear transformation matrix A
- $\mathbf{y}_i = \mathbf{x}_i \mathbf{A}$

• What is the matrix A?

#### Finding matrix A

- Elements A(i,j) can be Gaussian distributed
- Achlioptas\* has shown that the Gaussian distribution can be replaced by

$$A(i, j) = \begin{cases} +1 \text{ with prob } \frac{1}{6} \\ 0 \text{ with prob } \frac{2}{3} \\ -1 \text{ with prob } \frac{1}{6} \end{cases}$$

- All zero mean, unit variance distributions for A(i,j) would give a mapping that satisfies the JL lemma
- Why is Achlioptas result useful?

# Datasets in the form of matrices

Given **n** objects and **d** features describing the objects. (Each object has **d** numeric values describing it.)

#### <u>Dataset</u>

An **n-by-d** matrix **A**, **A**<sub>ij</sub> shows the "**importance**" of feature **j** for object **i**. Every row of **A** represents an object.

#### <u>Goal</u>

- 1. Understand the structure of the data, e.g., the underlying process generating the data.
- 2. Reduce the number of features representing the data

#### Market basket matrices



Find a subset of the products that characterize customer behavior



Find a subset of the groups that accurately clusters social-network users



Find a subset of the terms that accurately clusters the documents

#### **Recommendation systems**



Find a subset of the products that accurately describe the behavior or the customers

#### The Singular Value Decomposition (SVD)

Data matrices have **n** rows (one for each object) and **d** columns (one for each feature).

Rows: vectors in a Euclidean space,

Two objects are "**close**" if the angle between their corresponding vectors is small.



#### SVD: Example



Input: 2-d dimensional points

#### **Output:**

<u>1st (right) singular vector:</u> direction of maximal variance,

2nd (right) singular vector: direction of maximal variance, after removing the projection of the data along the first singular vector.

#### Singular values



 $\sigma_1$ : measures how much of the data variance is explained by the first singular vector.

 $\sigma_2$ : measures how much of the data variance is explained by the second singular vector.



**U**(**V**): orthogonal matrix containing the left (right) singular vectors of **A**.

 $\Sigma$ : diagonal matrix containing the singular values of A: (  $\sigma_1 \geq \sigma_2 \geq ... \geq \sigma_\ell$  )

Exact computation of the SVD takes O(min{mn<sup>2</sup>, m<sup>2</sup>n}) time.

The top k left/right singular vectors/values can be **computed faster** using Lanczos/Arnoldi methods.





 $A_k$  is an approximation of A

#### SVD as an optimization problem Find C to minimize:

$$\min_{C} \left\| A - C X_{n \times k} \right\|_{F}^{2}$$
 Frobenius norm:

$$\left\|A\right\|_{F}^{2} = \sum_{i,j} A_{ij}^{2}$$

Given **C** it is easy to find **X** from standard least squares. However, the fact that we can find the optimal **C** is fascinating!

#### SVD is "the Rolls-Royce and the Swiss Army Knife of Numerical Linear Algebra."\* \*Dianne O'Leary, MMDS '06

#### Reference

Simple and Deterministic Matrix Sketching Author: Edo Liberty, Yahoo! Labs KDD 2013, Best paper award

Thanks Edo Liberty for the slides

#### Sketches of streaming matrices

- A nxd matrix
- Rows of A arrive in a stream
- Task: compute

$$AA^T = \sum_{i=1}^n A_i A_i^t$$

#### Sketches of streaming matrices

- A dxn matrix
- Rows of A arrive in a stream
- Task: compute

$$AA^T = \sum_{i=1}^n A_i A_i^t$$

- Naive solution: Compute  $AA^T$  in time  ${\cal O}(nd^2)$  and space  ${\cal O}(d^2)$
- Think of d=10^6, n = 10^6

#### Goal

• Efficiently compute a concisely representable matrix B such that

$$B \approx A \text{ or } BB^T \approx AA^T$$

woking with **B** is good enough for many tasks

• Efficiently maintain matrix B with only  $\ell=2/\epsilon$  such that

$$||AA^T - BB^T||_2 \le \epsilon ||A||_f^2$$



• obtain the frequency f(i) of each item in a stream of items



• With d counters it's easy but not good enough





• Lets keep less than a fixed number of counters



• If an item has a counter we add 1 to that counter





• Otherwise, we create a new counter for it and set it to 1





- But now we do not have less than  $\ell$  counters

# Frequent items $\delta = f_{\ell/2} = 2$

- Let  $\delta$  be the median counter value at time t





- Decrease all counters by  $\delta$  (or set to zero if less than  $\delta$ )



• And continue....





• The approximated counts are f'



• We increase the count by only 1 for each item appearance

$$f'(i) \le f(i)$$

- Because we decrease each counter by at most  $\delta_t$  at time t

$$f'(i) \ge f(i) - \sum_t \delta_t$$

• Calculating the total approximated frequencies:

$$\begin{split} 0 \leq \sum_{i} f'(i) \leq \sum_{t} \left(1 - (\ell/2)\delta_{t}\right) &= n - (\ell/2)\sum_{t} \delta_{t} \\ \sum_{t} \delta_{t} \leq 2n/\ell \\ \text{Setting } \ell &= 2/\epsilon \\ |f(i) - f'(i)| \leq \epsilon n \end{split}$$



• We keep a sketch of at most  $\ell$  columns





Maintain the invariant that some of the columns are empty (zero-valued)



• Input vectors are simply stored in empty columns





• Input vectors are simply stored in empty columns





• When the sketch is ``full" we need to zero out some columns





• Using SVD we compute  $B = USV^T$  and set  $B_{new} = US$ 





• Note that  $BB^T = B_{new}B^T_{new}$  so we don't ``lose" anything





The columns of B are now orthogonal and in decreasing magnitude order

## **Frequent directions** $\twoheadrightarrow \delta = \|B_{\ell/2}\|^2$ d

• Let  $\delta = ||B_{\ell/2}||^2$ 





- Reduce column  $\,\ell_2^2 - {
m norms}\,$  by  $\,\delta$  (or nullify if less)





• Start aggregating columns again



Input:  $\ell$ ,  $A \in \mathbb{R}^{d \times n}$   $B \leftarrow$  all zeros matrix  $\in \mathbb{R}^{d \times \ell}$ for  $i \in [n]$  do Insert  $A_i$  into a zero valued column of Bif B has no zero valued colums then  $[U, \Sigma, V] \leftarrow SVD(B)$   $\delta \leftarrow \sigma_{\ell/2}^2$   $\check{\Sigma} \leftarrow \sqrt{\max(\Sigma^2 - I_\ell \delta, 0)}$  $B \leftarrow U\check{\Sigma}$  # At least half the columns of B are zero.

**Return:** *B* 



#### Frequent directions: proof

- Step 1:  $||AA^T - BB^T|| \le \sum_{t=1}^n \delta_t$
- Step 2:  $\sum_{t=1}^n \delta_t \leq 2||A||_f^2/\ell$
- Setting  $\ell = 2/\epsilon$  yields

$$||AA^T - BB^T|| \le \epsilon ||A||_f^2$$

#### Error as a function of $\ell$



