

In Search of a Matrix

CARL D. MEYER

North Carolina State University

Raleigh, NC

Canadian Mathematical Society

University of Saskatchewan

Saskatoon, Saskatchewan, June 2-4, 2001



Outline

- Background & History
- Vector Space Approach
- Frequency Weighting
- Compression & Noise Reduction
- Examples
- LSI with SVD
- LSI with URV & Centroids
- References

Background

Goal: Identify documents that best match users query

Measures

- ▶ Recall = $\frac{\#relevant\ docs\ retrieved}{\#docs\ in\ collection}$ (max # useful docs)
- ▶ Precision = $\frac{\#relevant\ docs\ retrieved}{\#docs\ retrieved}$ (min # useless docs)
- ▶ Do it *FAST!*

Methods

- ▶ Combinatorial
- ▶ Statistical
- ▶ Hashing
- ▶ Pattern matching
- ▶ Vector Space

SMART

(**S**ystem for the **M**echanical **A**nalysis and **R**etrieval of **T**ext)

Harvard 1962 – 1965

- ▶ IBM 7094 & IBM 360

Gerard Salton

- ▶ Implemented at Cornell (1965 – 1970)
- ▶ Based on matrix methods

Term–Document Matrix

Start With Dictionary of Terms

- ▶ Single words — or short phrases (e.g., *landing gear*)

Index Each Document (by human or by computer)

- ▶ Count $f_{ij} = \#$ times term i appears in document j

Unweighted Term–Document Matrix

$$\begin{array}{c} \text{TERM 1} \\ \text{TERM 2} \\ \vdots \\ \text{TERM } m \end{array} \begin{pmatrix} \text{Doc 1} & \text{Doc 2} & \cdots & \text{Doc } n \\ f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \cdots & f_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ f_{m1} & f_{m2} & \cdots & f_{mn} \end{pmatrix} = \mathbf{A}_{m \times n}$$

Features

- ▶ $\mathbf{A} \geq 0$
- ▶ \mathbf{A} can be really big! — Terms = $O(10^6)$ Docs = $O(10^7)$
- ▶ \mathbf{A} is sparse — but otherwise unstructured
- ▶ \mathbf{A} contains a lot of uncertainty (noise)

Example

(M. W. BERRY & M. BROWNE, 1999, SIAM)

Terms

T1: BAB(Y, IES, Y'S)
T2: CHILD(REN'S)
T3: GUIDE
T4: HEALTH
T5: HOME
T6: INFANT
T7: PROOFING
T8: SAFETY
T9: TODDLER

Documents

D1: INFANT AND TODDLER FIRST AID
D2: BABIES AND CHILDREN'S ROOM FOR YOUR HOME
D3: CHILD SAFETY AT HOME
D4: YOUR BABY'S HEALTH AND SAFETY: FROM INFANT TO TODDLER
D5: BABY PROOFING BASICS
D6: YOUR GUIDE TO EASY RUST PROOFING
D7: BEANIE BABIES COLLECTOR'S GUIDE

$$\mathbf{A} = \begin{matrix} & \begin{matrix} \text{D1} & \text{D2} & \text{D3} & \text{D4} & \text{D5} & \text{D6} & \text{D7} \end{matrix} \\ \begin{matrix} \text{T1} \\ \text{T2} \\ \text{T3} \\ \text{T4} \\ \text{T5} \\ \text{T6} \\ \text{T7} \\ \text{T8} \\ \text{T9} \end{matrix} & \begin{pmatrix} \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \end{matrix} \quad \mathbf{9 \times 7}$$

Example

(M. W. BERRY & M. BROWNE, 1999, SIAM)

Terms

T1: BAB(Y, IES, Y'S)
T2: CHILD(REN'S)
T3: GUIDE
T4: HEALTH
T5: HOME
T6: INFANT
T7: PROOFING
T8: SAFETY
T9: TODDLER

Documents

D1: INFANT AND TODDLER FIRST AID
D2: BABIES AND CHILDREN'S ROOM FOR YOUR HOME
D3: CHILD SAFETY AT HOME
D4: YOUR BABY'S HEALTH AND SAFETY: FROM INFANT TO TODDLER
D5: BABY PROOFING BASICS
D6: YOUR GUIDE TO EASY RUST PROOFING
D7: BEANIE BABIES COLLECTOR'S GUIDE

$$\mathbf{A} = \begin{matrix} & \begin{matrix} \text{D1} & \text{D2} & \text{D3} & \text{D4} & \text{D5} & \text{D6} & \text{D7} \end{matrix} \\ \begin{matrix} \text{T1} \\ \text{T2} \\ \text{T3} \\ \text{T4} \\ \text{T5} \\ \text{T6} \\ \text{T7} \\ \text{T8} \\ \text{T9} \end{matrix} & \begin{pmatrix} 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix} \end{matrix} \quad 9 \times 7$$

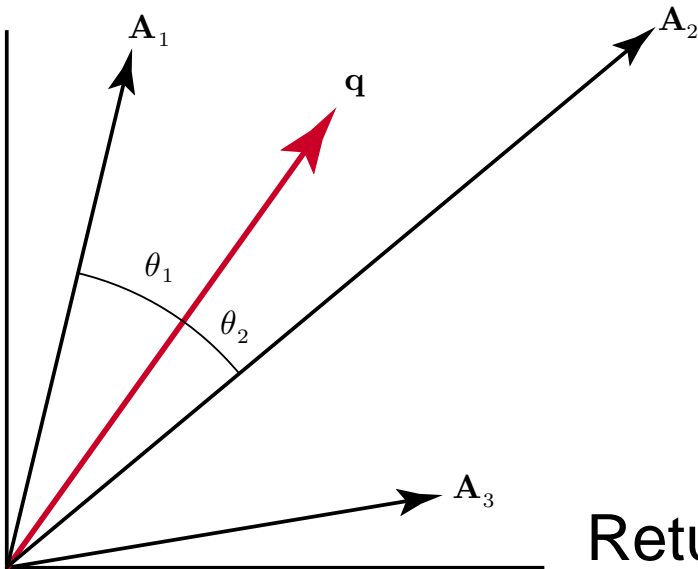
Query Matching

Query Vector

- ▶ $\mathbf{q}^T = (q_1, q_2, \dots, q_m)$ where $q_i = \begin{cases} 1 & \text{if Term } i \text{ is requested} \\ 0 & \text{if not} \end{cases}$

How Close is the Query to Each Document?

- ▶ i.e., how close is \mathbf{q} to each column \mathbf{A}_k ?



$$\|\mathbf{q} - \mathbf{A}_1\| < \|\mathbf{q} - \mathbf{A}_2\| \text{ but } \theta_2 < \theta_1$$

$$\text{Use } \delta_i = \cos \theta_i = \frac{\mathbf{q}^T \mathbf{A}_i}{\|\mathbf{q}\| \|\mathbf{A}_i\|}$$

Rank documents by size of δ_i

Return Document i to user when $\delta_i \geq tol$

Term Weighting

A Defect

- ▶ If the term *bank* occurs once in Doc 1 but twice in Doc 2, and if $\|\mathbf{A}_1\| \approx \|\mathbf{A}_2\|$, then a query containing only *bank* produces $\delta_2 \approx 2\delta_1$ (i.e., Doc 2 is rated twice as relevant as Doc 1).

To Compensate

- ▶ Set $a_{ij} = \log(1 + f_{ij})$ (other weights also possible)

Query Weights

- ▶ Terms *Boeing* and *airplanes* not equally important in query
- ▶ Importance of Term i tends to be inversely proportional to $\nu_i = \#$ Docs containing Term i

To Compensate

- ▶ Set $q_i = \begin{cases} \log(n/\nu_i) & \text{if } \nu_i \neq 0 \\ 0 & \text{if } \nu_i = 0 \end{cases}$ (other weights also possible)

Noise in A

Ambiguity in Vocabulary

- ▶ e.g., A *bank* could be
 - A financial institution
 - A river side
 - A shot in the game of pool

Variation in Writing Style

- ▶ No two authors write the same way
 - One author may write *car* and *laptop*
 - Another author may write *automobile* and *portable*

Variation in Indexing Conventions

- ▶ No two people index documents the same way
- ▶ Computer indexing is inexact and can be unpredictable

Practical Problems

Simple in Theory

- ▶ Weight terms and normalize cols — Make $\|\mathbf{A}_i\| = 1$
- ▶ For each new query, weight and normalize — Make $\|\mathbf{q}\| = 1$
- ▶ Compute $\delta_i = \cos \theta_i = (\mathbf{q}^T \mathbf{A})_i$ and return most relevant docs

Difficult in Practice

- ▶ Must be able to do it *FAST!*
 - ⇒ Somehow compress the data
- ▶ Must account for variations due to ambiguity in language and variations in writing and indexing styles.
 - ⇒ Somehow reduce “noise” or “uncertainty” in \mathbf{A}

Basis Games

Consider Vector of Data

$\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{e}_i$ with respect to o.n. basis $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$

$\implies \alpha_i = \langle \mathbf{e}_i | \mathbf{v} \rangle = \text{amount of } \mathbf{v} \text{ in direction of } \mathbf{e}_i$

Compression

- ▶ Select new o.n. basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ so that fewer vectors are needed to represent $\mathbf{v} = \sum_{i=1}^r \beta_i \mathbf{u}_i$ ($r < n$)
- ▶ Data is compressed from $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ into $\{\beta_1, \beta_2, \dots, \beta_r\}$

Even More Compression

- ▶ Eliminate data lying in insignificant directions
Arrange: $|\beta_1| \geq |\beta_2| \geq \dots \geq |\beta_k| \geq \epsilon > |\beta_{k+1}| \dots \geq |\beta_r|$
Approximate: $\mathbf{v} \approx \tilde{\mathbf{v}} = \sum_{i=1}^k \beta_i \mathbf{u}_i$ ($k < r < n$)
- ▶ Data now compressed from $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ to $\{\beta_1, \beta_2, \dots, \beta_k\}$

Added Benefit

Noise Reduction

- ▶ Assume noise (or uncertainty) is nondirectional

⇒ As much noise in one direction as in any other direction

⇒ $\mathbf{v} = \sum_{i=1}^r \beta_i \mathbf{u}_i = \sum_{i=1}^r s_i \mathbf{u}_i + \sum_{i=1}^r \epsilon \mathbf{u}_i = (\text{signal}) + (\text{noise})$

- ▶ Suppose $|s_1| \geq |s_2| \geq \dots \geq |s_{\mathbf{k}}| \geq \epsilon > |s_{\mathbf{k}+1}| \dots \geq |s_r|$

Drop $\beta_{\mathbf{k}+1}, \beta_{\mathbf{k}+2}, \dots, \beta_r$, and use $\mathbf{v} \approx \tilde{\mathbf{v}} = \sum_{i=1}^{\mathbf{k}} \beta_i \mathbf{u}_i$

⇒ Only a small proportion of the signal is lost

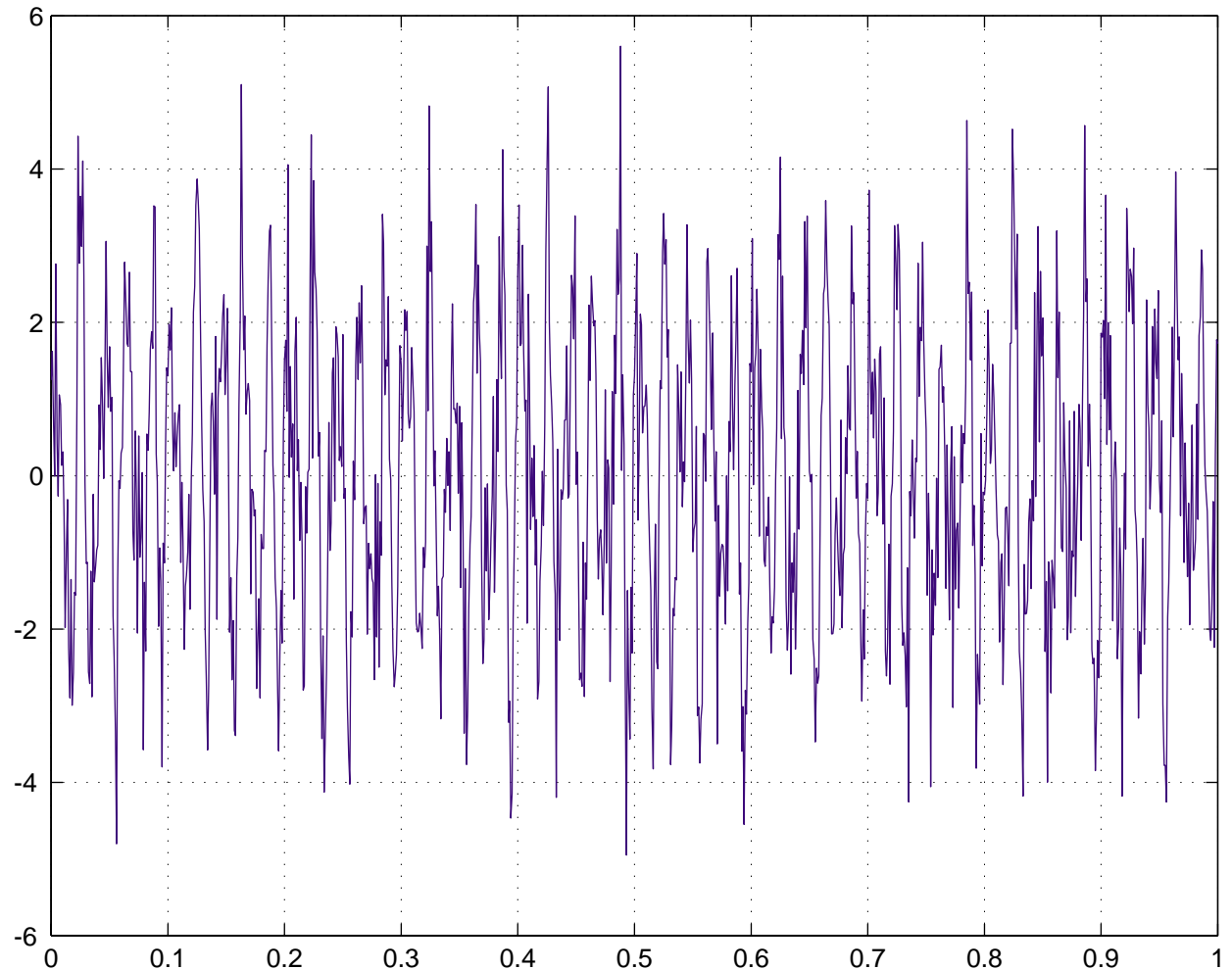
⇒ A larger proportion of the noise is lost

Example

(MATRIX ANALYSIS AND APPLIED LINEAR ALGEBRA, SIAM, 2000)

Sample Audio Signal 512 times for 1 sec.

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{510} \\ x_{511} \end{bmatrix}$$



Goal: Compress and Reduce Noise

Find A Better Basis

Oscillatory \implies Cosines & Sines

$$\mathbf{x}(t) \sim \sum_{f=1}^{\infty} \alpha_f \cos 2\pi f t + \beta_f \sin 2\pi f t \quad (\text{IF } \mathbf{x}(t) \text{ WAS A CONTINUOUS FUNCTION})$$

Discrete Time, Cosine, and Sine Vectors

$$\mathbf{t} = \begin{bmatrix} 0/n \\ 1/n \\ 2/n \\ \vdots \\ n-1/n \end{bmatrix} \quad \cos 2\pi f \mathbf{t} = \begin{bmatrix} \cos \left(2\pi f \cdot \frac{0}{n} \right) \\ \cos \left(2\pi f \cdot \frac{1}{n} \right) \\ \cos \left(2\pi f \cdot \frac{2}{n} \right) \\ \vdots \\ \cos \left(2\pi f \cdot \frac{n-1}{n} \right) \end{bmatrix} \quad \sin 2\pi f \mathbf{t} = \begin{bmatrix} \sin \left(2\pi f \cdot \frac{0}{n} \right) \\ \sin \left(2\pi f \cdot \frac{1}{n} \right) \\ \sin \left(2\pi f \cdot \frac{2}{n} \right) \\ \vdots \\ \sin \left(2\pi f \cdot \frac{n-1}{n} \right) \end{bmatrix}$$

Discrete Exponential Vectors

$$e^{i2\pi f \mathbf{t}} = \cos 2\pi f \mathbf{t} + i \sin 2\pi f \mathbf{t} \quad e^{-i2\pi f \mathbf{t}} = \cos 2\pi f \mathbf{t} - i \sin 2\pi f \mathbf{t}$$

Discrete Fourier Transform

$$\omega = e^{2\pi i/n} \quad \mathbf{W} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega & \omega^2 & \cdots & \omega^{n-1} \\ 1 & \omega^2 & \omega^4 & \cdots & \omega^{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{n-1} & \omega^{n-2} & \cdots & \omega \end{bmatrix}_{n \times n} \quad \mathbf{W}_f = \frac{1}{2} e^{2\pi i f t}$$

Identities

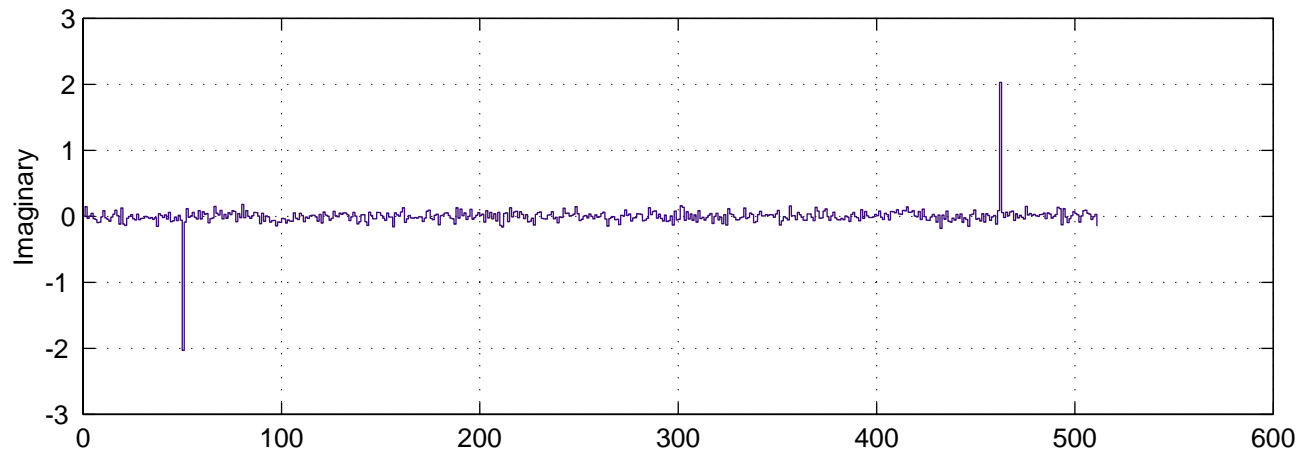
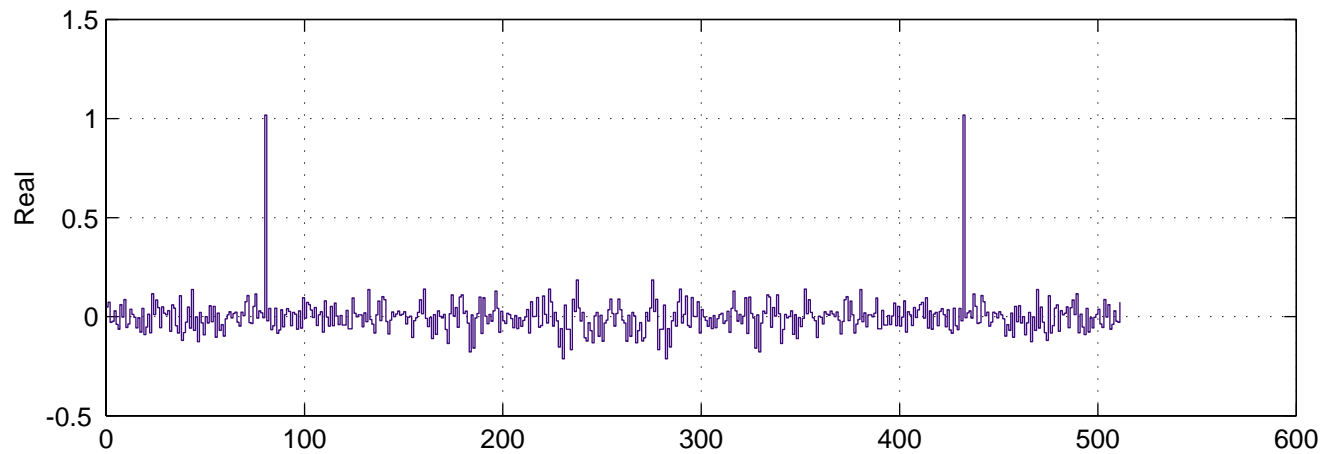
- ▶ $\cos 2\pi f t = \mathbf{W}_f + \mathbf{W}_{n-f} \quad (0 < f < n)$
- ▶ $\sin 2\pi f t = -i(\mathbf{W}_f - \mathbf{W}_{n-f})$

New Basis = $\{\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_{n-1}\}$ (Find coordinates of \mathbf{x} wrt \mathbf{W}_i 's)

$$\mathbf{x} = \mathbf{W}\mathbf{y} \implies \mathbf{y} = \mathbf{W}^{-1}\mathbf{x} = \frac{2}{n} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \xi & \xi^2 & \cdots & \xi^{n-1} \\ 1 & \xi^2 & \xi^4 & \cdots & \xi^{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \xi^{n-1} & \xi^{n-2} & \cdots & \xi \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{n-1} \end{bmatrix}$$

$$\xi = e^{-2\pi i/n} = \bar{\omega}$$

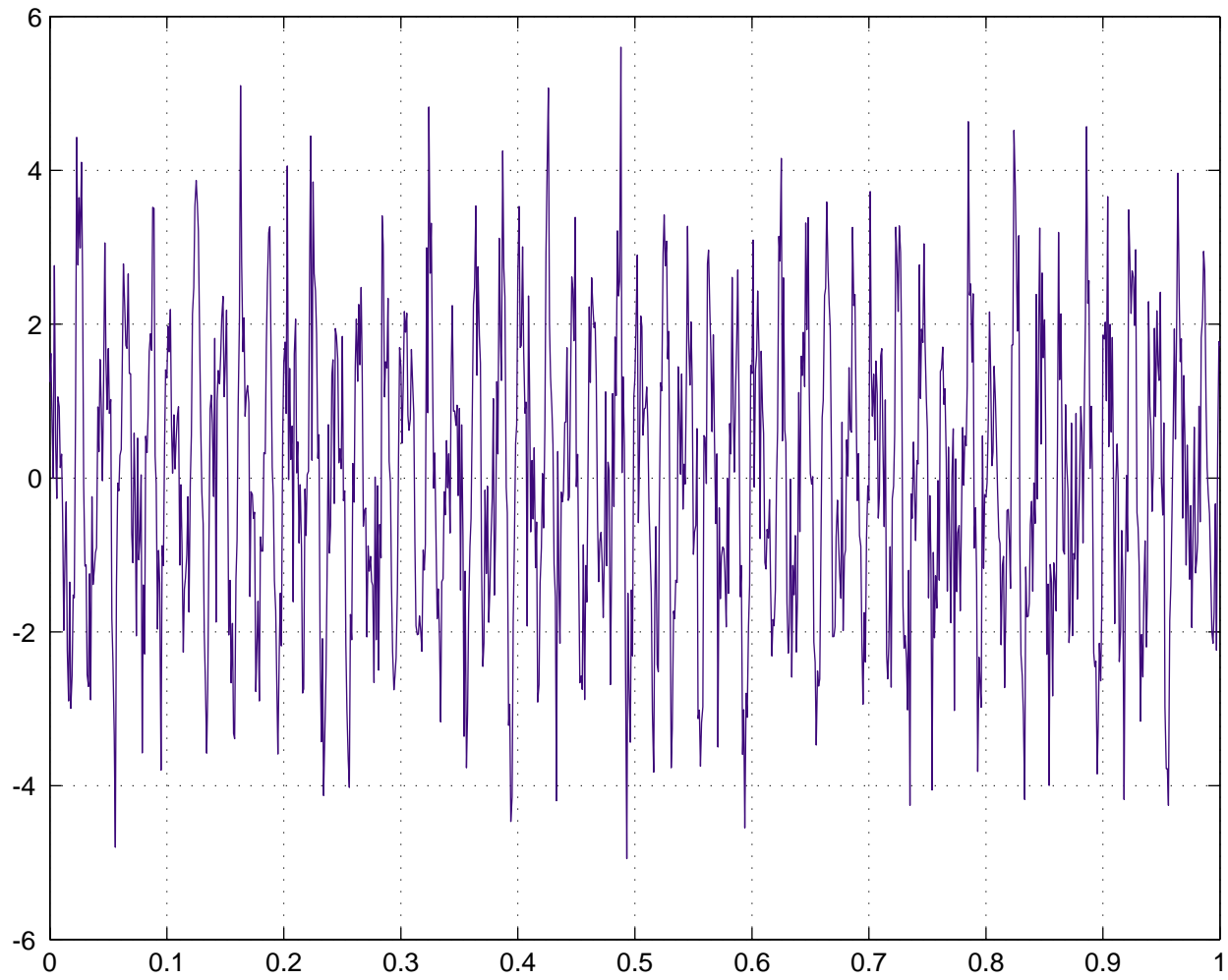
The New Coordinates



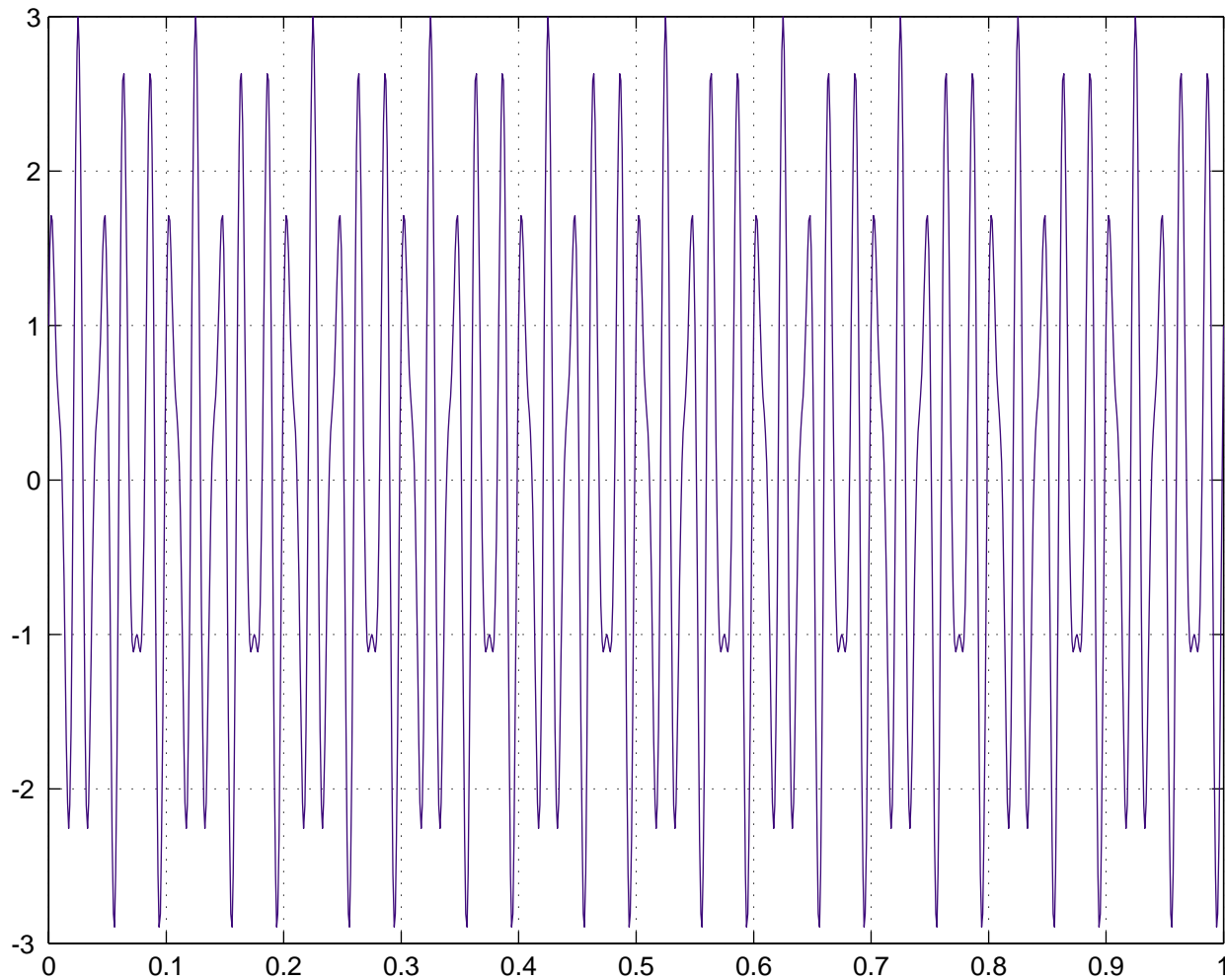
$$\begin{aligned} \mathbf{y} &= (\mathbf{e}_{80} + \mathbf{e}_{432}) - 2i(\mathbf{e}_{50} - \mathbf{e}_{462}) + \boldsymbol{\varepsilon} \\ &= (\mathbf{e}_{80} + \mathbf{e}_{n-80}) - 2i(\mathbf{e}_{50} - \mathbf{e}_{n-50}) + \boldsymbol{\varepsilon} \end{aligned} \quad (n = 512)$$

$$\begin{aligned} \mathbf{x} = \mathbf{W}\mathbf{y} &= (\mathbf{W}_{80} + \mathbf{W}_{n-80}) - 2i(\mathbf{W}_{50} - \mathbf{W}_{n-50}) + \boldsymbol{\varepsilon} \\ &= \cos 2\pi 80t + 2 \sin 2\pi 50t + \boldsymbol{\varepsilon} \end{aligned}$$

Original



Cleaned & Compressed



$$\cos 2\pi 80t + 2 \sin 2\pi 50t$$

The DFT Game

Matrix–Vector Product

$$\mathbf{y} = \frac{2}{n} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \xi & \xi^2 & \cdots & \xi^{n-1} \\ 1 & \xi^2 & \xi^4 & \cdots & \xi^{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \xi^{n-1} & \xi^{n-2} & \cdots & \xi \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{n-1} \end{bmatrix} \quad \xi = e^{-2\pi i/n}$$

Simple in Theory but Difficult in Practice

- ▶ Must do it *FAST!*

Need For Speed \implies Matrix Factorizations \implies FFT

- ▶ $\mathbf{F}_n = \mathbf{B}_n (\mathbf{I}_2 \otimes \mathbf{F}_{n/2}) \mathbf{P}_n$ $\mathbf{B}_n = \begin{bmatrix} \mathbf{I}_{n/2} & \mathbf{D}_{n/2} \\ \mathbf{I}_{n/2} & -\mathbf{D}_{n/2} \end{bmatrix}$ $\mathbf{D}_{n/2} = \begin{bmatrix} 1 & \xi & \xi^2 & \cdots \\ & \xi & \xi^2 & \ddots \\ & & \xi^2 & \ddots \\ & & & \ddots \end{bmatrix}$

- ▶ FFT changes n^2 flop requirement into $(n/2) \log_2 n$

“THE MOST VALUABLE NUMERICAL ALGORITHM IN OUR LIFETIME.”

—G. STRANG, BULLETIN OF THE AMS, APRIL, 1993.

Things Have Changed

- ▶ “For engineers and social and physical scientists, linear algebra now fills a place that is often more important than calculus.”
- ▶ “It is partly [due to] the move from analog to digital; functions are replaced by vectors.”
- ▶ “My generation of students, and certainly my teachers, did not see this change coming.”

—Gilbert Strang

FROM THE AMERICAN SCIENTIST

APRIL, 1994.

Back To IR

Almost the Same Problem

- ▶ Evaluate $\mathbf{q}^T \mathbf{A}$ fast

Data is Now the Term-Doc Matrix in Standard Coordinates

- ▶ $\mathbf{A} = \sum_{i,j} \langle \mathbf{E}_{ij} | \mathbf{A} \rangle \mathbf{E}_{ij}$ $\mathbf{E}_{ij} = \mathbf{e}_i \mathbf{e}_j^T$ $\langle \mathbf{E}_{ij} | \mathbf{A} \rangle = a_{ij}$
- ▶ $\langle \mathbf{X} | \mathbf{Y} \rangle = \text{trace}(\mathbf{X}^T \mathbf{Y})$ $\|\mathbf{X}\| = \langle \mathbf{X} | \mathbf{X} \rangle^{1/2}$ (Frobenius Norm)

Seek New o.n. Basis That Squeezes & Cleans

- ▶ $\mathbf{A} = \sum_{i=1}^r \langle \mathbf{Z}_i | \mathbf{A} \rangle \mathbf{Z}_i$

Think Matrix Factorizations \implies SVD \implies $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$

- ▶ $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \sum_{i=1}^r \sigma_i \mathbf{Z}_i$, $\mathbf{Z}_i = \mathbf{u}_i \mathbf{v}_i^T$, $\langle \mathbf{Z}_i | \mathbf{Z}_j \rangle = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$
- ▶ $\sigma_i = \langle \mathbf{Z}_i | \mathbf{A} \rangle =$ the amount of \mathbf{A} in direction of \mathbf{Z}_i

Same As Before

Assume Nondirectional Noise or Uncertainty

- ▶ Drop small σ_i 's and use $\mathbf{A} \approx \tilde{\mathbf{A}} = \sum_{i=1}^k \sigma_i \mathbf{Z}_i$
- ▶ Lose only small part of relevance
- ▶ Lose larger proportion of noise or uncertainty

Be Liberal in Dropping σ_i 's

- ▶ Numerical accuracy not important — 2 or 3 significant digits

New Query Matching Strategy

- ▶ Normalize
 - $\mathbf{q} \leftarrow \mathbf{q} / \|\mathbf{q}\|$
 - $\tilde{\mathbf{A}} \leftarrow \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T \mathbf{D} = \sum_{i=1}^k \sigma_i \mathbf{u}_i \tilde{\mathbf{v}}_i^T$
- ▶ Compare query to each document
 - $(\delta_1, \delta_2, \dots, \delta_n) = \mathbf{q}^T \tilde{\mathbf{A}} = \sum_{i=1}^k \sigma_i (\mathbf{q}^T \mathbf{u}_i) \tilde{\mathbf{v}}_i^T$

Pros & Cons

Advantages

- ▶ Compression
 - \mathbf{A} is replaced by only a few sing values and sing vectors
 - They are determined & normalized only once
- ▶ *SPEED!*
 - Each query requires only a few inner products

$$\mathbf{q}^T \tilde{\mathbf{A}}_{m \times n} = \sum_{i=1}^k \sigma_i (\mathbf{q}^T \mathbf{u}_i) \tilde{\mathbf{v}}_i^T$$

- ▶ Latent semantic associations are made
 - Relevant docs not found by direct matching show up

Disadvantages

- ▶ Adding & deleting docs (updating & downdating SVD) difficult
- ▶ Determining optimal k is not easy (empirical tuning required)

Variations

Projected Query

- ▶ First project the query onto the document space

$$\tilde{\mathbf{q}} = \mathbf{P}_{R(\mathbf{A})} \mathbf{q} = \sum_{j=1}^r \mathbf{u}_j \mathbf{u}_j^T \mathbf{q}$$

- ▶ Or, better yet, use truncated projection

$$\tilde{\mathbf{q}} = \mathbf{P}_{R(\tilde{\mathbf{A}})} \mathbf{q} = \sum_{j=1}^k \mathbf{u}_j \mathbf{u}_j^T \mathbf{q}$$

- ▶ Notice

- $\tilde{\mathbf{q}}^T \tilde{\mathbf{A}} = \sum_{i=1}^k \sigma_i (\tilde{\mathbf{q}}^T \mathbf{u}_i) \tilde{\mathbf{v}}_i^T = \sum_{i=1}^k \sigma_i \left(\left(\sum_{j=1}^k \mathbf{q}^T \mathbf{u}_j^T \mathbf{u}_j \right) \mathbf{u}_i \right) \tilde{\mathbf{v}}_i^T = \mathbf{q}^T \tilde{\mathbf{A}}$
- $\|\tilde{\mathbf{q}}\| \leq \|\mathbf{q}\|$
- $\cos \tilde{\theta}_i \geq \cos \theta_i$ (more documents are deemed relevant)

Other Factorizations

DFT — FFT

- ▶ No compression — no oscillatory components

Haar Transform $\mathbf{H}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ $\mathbf{H}_4 = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & -1 & 0 \\ 1 & -1 & 0 & 1 \\ 1 & -1 & 0 & -1 \end{bmatrix}$

- ▶ $\mathbf{H}_n = (\mathbf{I}_2 \otimes \mathbf{H}_{n/2}) \mathbf{P}_n \begin{bmatrix} \mathbf{H}_{n/2} & \\ & \mathbf{I}_{n/2} \end{bmatrix} \Rightarrow \mathbf{H}_n \mathbf{x}$ is *Fast* (if $n = 2^p$)
- ▶ Factor $\mathbf{A} = \mathbf{H}_m \mathbf{B} \mathbf{H}_n^T = \sum_{i,j} \beta_{ij} \mathbf{h}_i \mathbf{h}_j^T$ (\mathbf{h} 's only use $-1, 0,$ or 1)
 - More than a few β_{ij} 's may be needed
 - Needs padding if m or n not a power of 2

Semidiscrete Decomposition (T. KOLDA AND D. O'LEARY, 1998)

- ▶ Approximate $\mathbf{A} \approx \sum_{i=1}^k \alpha_i \mathbf{x}_i \mathbf{y}_j$ \mathbf{x}_i and \mathbf{y}_j only use $-1, 0,$ or 1

Other Wavelet Transforms?

Using Centroids

Document Clusters

- ▶ Assume $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$ (normalized) represents a cluster
- ▶ Mean: $\mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{a}_i = \frac{\mathbf{A}\mathbf{e}}{n}$ where $\mathbf{e}^T = (1, 1, \dots, 1)$
- ▶ Centroid vector: $\mathbf{c} = \mathbf{m} / \|\mathbf{m}\| = \mathbf{A}\mathbf{e} / \|\mathbf{A}\mathbf{e}\|$

Facts

- ▶ $\min_{\mathbf{p} \geq 0, \|\mathbf{p}\|=1} \sum_{i=1}^n \cos \theta(\mathbf{a}_i, \mathbf{p}) = \sum_{i=1}^n \cos \theta(\mathbf{a}_i, \mathbf{c}).$
- ▶ i.e., \mathbf{c} is the closest vector to the cluster
- ▶ $\mathbf{c}^T \mathbf{A} \approx \mathbf{e}^T \quad (\mathbf{c}^T \mathbf{a}_i = \cos \theta_i \approx 1)$

Reflectors With Specified Unit Column — Say \mathbf{c}

- ▶ $\mathbf{L} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^T / (\mathbf{u}^T \mathbf{u}) = [\mathbf{c} \mid \mathbf{X}]$ where $\mathbf{u} = \mathbf{c} \pm \mathbf{e}_1$
- ▶ $\mathbf{L} = \mathbf{L}^T = \mathbf{L}^{-1}$

Orthogonal Factorization

Specify \mathbf{c} and \mathbf{e}/\sqrt{n} as First Columns

$$\blacktriangleright \mathbf{L} = [\mathbf{c} \mid \mathbf{X}] \quad \text{and} \quad \mathbf{R} = [\mathbf{e}/\sqrt{n} \mid \mathbf{Y}]$$

$$\begin{aligned} \mathbf{LAR} = \mathbf{L}^T \mathbf{AR} &= \begin{bmatrix} \mathbf{c}^T \mathbf{Ae}/\sqrt{n} & \mathbf{c}^T \mathbf{AY} \\ \mathbf{X}^T \mathbf{Ae}/\sqrt{n} & \mathbf{X}^T \mathbf{AY} \end{bmatrix} \approx \begin{bmatrix} \mathbf{e}^T \mathbf{e}/\sqrt{n} & \mathbf{e}^T \mathbf{Y} \\ \mathbf{X}^T \mathbf{c} \|\mathbf{Ae}\|/\sqrt{n} & \mathbf{X}^T \mathbf{AY} \end{bmatrix} \\ &= \begin{bmatrix} \sqrt{n} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^T \mathbf{AY} \end{bmatrix} = \begin{bmatrix} \sqrt{n} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \end{aligned}$$

A New Factorization

$$\blacktriangleright \mathbf{A} = \mathbf{LDR} = \mathbf{LDR}^T = \mathbf{ce}^T + \sum_{i,j \neq 1} \beta_{ij} \mathbf{x}_i \mathbf{y}_j^T$$

Truncate To Compress & Clean

Break Entire Collection Into Clusters

To Learn More

Books

- MATRIX ANALYSIS AND APPLIED LINEAR ALGEBRA, C. D. MEYER, SIAM, 2000.
- UNDERSTANDING SEARCH ENGINES; MATHEMATICAL MODELING AND TEXT RETRIEVAL, M. W. BERRY AND M. BROWNE, SIAM, 1999.
- INTRODUCTION TO MODERN INFORMATION RETRIEVAL, G. SALTON AND M. MCGILL, MCGRAW-HILL, 1983.

Papers

- M. W. BERRY, Z. DRMAC, AND E. R. JESSUP, MATRICES, VECTOR SPACES, AND INFORMATION RETRIEVAL, SIAM REV., 41(1999), PP.335-362.
- M. W. BERRY, S. T. DUMAS, AND G. W. O'BRIEN, USING LINEAR ALGEBRA FOR INTELLIGENT INFORMATION RETRIEVAL, U. TENN. COMP. SCI. REPORT CS-94-270, DEC, 1994.
- I. S. DHILLON AND D. S. MODHA, CONCEPT DECOMPOSITIONS FOR LARGE SPARSE TEXT DATA USING CLUSTERING, IBM RESEARCH REPORT RJ 10147 (95022), JULY 8, 1999-DECLASSIFIED ON MARCH 13, 2000, TO APPEAR IN MACHINE LEARNING.

URLs

- [HTTP://WWW.CS.UTK.EDU/~ LSI/](http://www.cs.utk.edu/~lsi/)
- [HTTP://LSI.RESEARCH.TELCORDIA.COM/](http://lsi.research.telcordia.com/)
- [HTTP://LSA.COLORADO.EDU/](http://lsa.colorado.edu/)
- [HTTP://WWW.SEARCHENGINEWATCH.COM/RESOURCES/INDEX.HTML](http://www.searchenginewatch.com/resources/index.html)