

# Methods for sparse analysis of high-dimensional data, II

Rachel Ward

May 23, 2011

# High dimensional data with low-dimensional structure



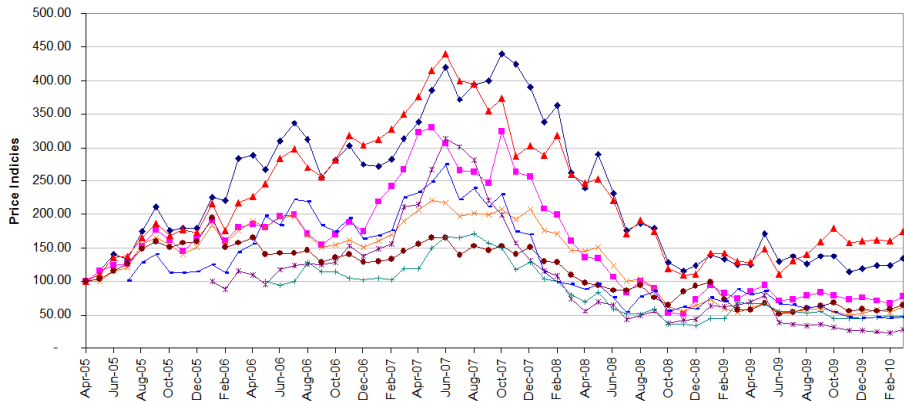
300 by 300 pixel images = 90,000 dimensions

# High dimensional data with low-dimensional structure



# High dimensional data with low-dimensional structure

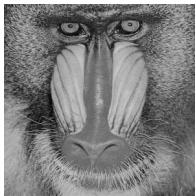
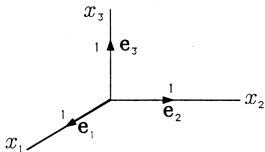
Chart 1: Monthly Stock Price Movements Over 5-Yr Period



We need to recall some ...

- Euclidean geometry
- Statistics
- Linear algebra

# Euclidean Geometry



- An element of  $\mathbb{R}^n$  is written

$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

- $\mathbb{R}^n$  is a vector space:

- $\mathbf{x} + \mathbf{y} = (x_1 + y_1, x_2 + y_2, \dots, x_n + y_n)$

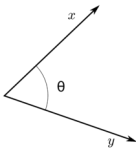
- $a\mathbf{x} = (ax_1, ax_2, \dots, ax_n)$

- $\mathbf{x} = (x_1, x_2, \dots, x_n) = \sum_{j=1}^n x_j \mathbf{e}_j$   
where

$$\mathbf{e}_1 = (1, 0, \dots, 0), \quad \mathbf{e}_2 = (0, 1, \dots, 0), \dots$$

$$\mathbf{e}_n = (0, 0, \dots, 1)$$

are the **standard basis vectors**.



- The **inner product** between  $\mathbf{x}$  and  $\mathbf{y}$  is:

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum_{j=1}^n x_j y_j$$

- $\|\mathbf{x}\| := \langle \mathbf{x}, \mathbf{x} \rangle^{1/2} = (x_1^2 + x_2^2 + \dots + x_n^2)^{1/2}$  is the Euclidean **length** of  $\mathbf{x}$ . It is a **norm**:

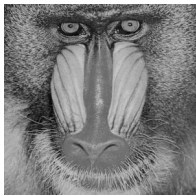
- $\|\mathbf{x}\| = 0$  if and only if  $\mathbf{x} = \mathbf{0}$ .
- $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$
- **triangle inequality**:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

- $\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} = \cos(\theta)$

- $\mathbf{x}$  and  $\mathbf{y}$  are **orthogonal** (perpendicular) if and only if  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$



# Statistics



$$\mathbf{x} = (x_1, x_2, x_3, \dots, x_n) \in \mathbb{R}^n$$

■ Sample mean:  $\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$

■ Standard deviation:

$$s = \sqrt{\frac{\sum_{j=1}^n (x_j - \bar{x})^2}{n-1}} = \frac{1}{\sqrt{n-1}} \sqrt{\langle \mathbf{x} - \bar{\mathbf{x}}, \mathbf{x} - \bar{\mathbf{x}} \rangle}$$



- Variance:  $s^2 = \frac{1}{n-1} \langle \mathbf{x} - \bar{\mathbf{x}}, \mathbf{x} - \bar{\mathbf{x}} \rangle = \frac{1}{n-1} \|\mathbf{x} - \bar{\mathbf{x}}\|^2$

- Suppose we have  $p$  data vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$

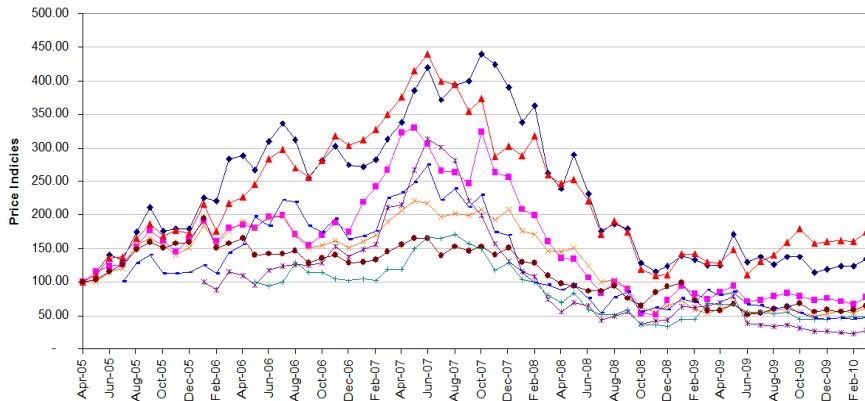
- **Covariance:**  $\text{Cov}(\mathbf{x}_j, \mathbf{x}_k) = \frac{1}{n-1} \langle \mathbf{x}_j - \bar{\mathbf{x}}_j, \mathbf{x}_k - \bar{\mathbf{x}}_k \rangle$

- **Covariance matrix** for 3 data vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ :

$$\mathcal{C} = \begin{pmatrix} \text{cov}(\mathbf{x}_1, \mathbf{x}_1) & \text{cov}(\mathbf{x}_1, \mathbf{x}_2) & \text{cov}(\mathbf{x}_1, \mathbf{x}_3) \\ \text{cov}(\mathbf{x}_2, \mathbf{x}_1) & \text{cov}(\mathbf{x}_2, \mathbf{x}_2) & \text{cov}(\mathbf{x}_2, \mathbf{x}_3) \\ \text{cov}(\mathbf{x}_3, \mathbf{x}_1) & \text{cov}(\mathbf{x}_3, \mathbf{x}_2) & \text{cov}(\mathbf{x}_3, \mathbf{x}_3) \end{pmatrix}$$

- Covariance matrix for  $p$  data vectors has  $p$  columns and  $p$  rows

Chart 1: Monthly Stock Price Movements Over 5-Yr Period



What does the covariance matrix look like?

# Linear Algebra

# Eigenvectors

Suppose  $\mathcal{A}$  is a  $p \times p$  matrix. If  $\mathcal{A}\mathbf{v} = \lambda\mathbf{v}$ , then we say  $\mathbf{v}$  is an **eigenvector** of  $\mathcal{A}$  with **eigenvalue**  $\lambda$ .

Are these eigenvectors?

$$\mathcal{A} = \begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix}, \quad \mathbf{v} = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$$

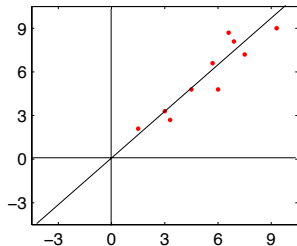
$$\mathcal{A} = \begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix}, \quad \mathbf{v} = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

- If  $\mathbf{v}$  is an eigenvector of  $\mathcal{A}$  with eigenvalue  $\lambda$ , then  $\alpha\mathbf{v}$  is also an eigenvector of  $\mathcal{A}$  with eigenvalue  $\lambda$ . **We will always use the normalized eigenvector  $\|\mathbf{v}\| = 1$ .**

- Any **real-valued and symmetric** matrix  $\mathcal{C}$  has  $n$  eigenvectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  which form an **orthonormal basis** for  $\mathbb{R}^n$  (a.k.a. rotated coordinate view).
- Any  $\mathbf{x} \in \mathbb{R}^n$  can be expressed in this basis via  $\mathbf{x} = \sum_{j=1}^n \langle \mathbf{x}, \mathbf{v}_j \rangle \mathbf{v}_j$ .
- $\mathcal{C}\mathbf{x} = \sum_{j=1}^n \lambda_j \langle \mathbf{x}, \mathbf{v}_j \rangle \mathbf{v}_j$
- $\mathcal{C} = \mathcal{P}\mathcal{D}\mathcal{P}^{-1}$  is **diagonalizable**:

$$\mathcal{P} = \begin{bmatrix} - & - & - & \mathbf{v}_1 & - & - & - \\ - & - & - & \mathbf{v}_2 & - & - & - \\ & & & \vdots & & & \\ - & - & - & \mathbf{v}_n & - & - & - \end{bmatrix}, \quad \mathcal{D} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & & \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

## Example



$$\mathbf{x} = (7.5, 1.5, 6.6, 5.7, 9.3, 6.9, 6, 3, 4.5, 3.3),$$
$$\mathbf{y} = (7.2, 2.1, 8.7, 6.6, 9, 8.1, 4.8, 3.3, 4.8, 2.7)$$

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \langle \mathbf{x} - \bar{\mathbf{x}}, \mathbf{y} - \bar{\mathbf{y}} \rangle,$$

$$\mathcal{C} = \begin{pmatrix} \text{cov}(\mathbf{x}, \mathbf{x}) & \text{cov}(\mathbf{x}, \mathbf{y}) \\ \text{cov}(\mathbf{x}, \mathbf{y}) & \text{cov}(\mathbf{y}, \mathbf{y}) \end{pmatrix} = \begin{pmatrix} 5.549 & 5.539 \\ 5.539 & 6.449 \end{pmatrix}$$



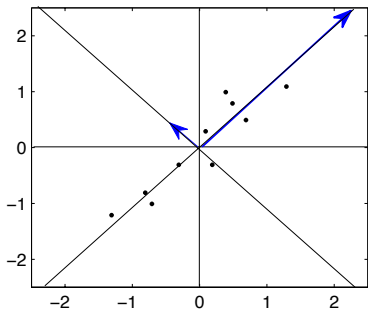


Figure:  $x - \bar{x}$  vs.  $y - \bar{y}$

Eigenvectors / values for  $C$ :

$$\blacksquare \mathbf{v}_1 = \begin{pmatrix} .6780 \\ .7352 \end{pmatrix}, \lambda_1 = 11.5562$$

$$\blacksquare \mathbf{v}_2 = \begin{pmatrix} -.7352 \\ .6780 \end{pmatrix}, \lambda_2 = .4418$$

- $\mathbf{v}_1$  the first **principal component** of the data  $(\mathbf{x}, \mathbf{y})$ , and  $\mathbf{v}_2$  the second 'principal component', and so-on ...
- **Prove:**  $\mathbf{v}_1$  is in the direction of the 'least squares fit' to the centered data  $(x_j - \bar{x}, y_j - \bar{y})$ ,  $j = 1, 2, \dots, n$ .

# Principal component analysis

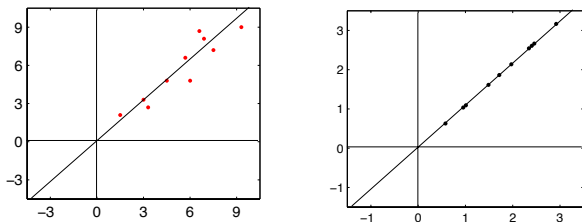


Figure: Original data and projection onto first principal component

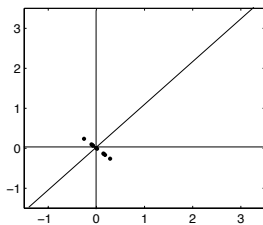
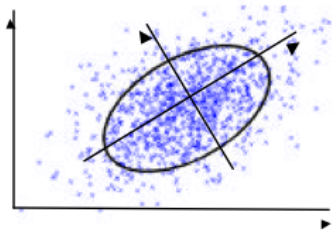


Figure: Residual

# Principal component analysis



"Best fit ellipsoid" to the data

- The covariance matrix is written as  $\mathcal{C} = \mathcal{P}\mathcal{D}\mathcal{P}^{-1}$ , where

$$\mathcal{P} = \begin{bmatrix} - & - & - & \mathbf{v}_1 & - & - & - \\ - & - & - & \mathbf{v}_2 & - & - & - \\ & & & \vdots & & & \\ - & - & - & \mathbf{v}_n & - & - & - \end{bmatrix}, \quad \mathcal{D} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & & \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

Suppose that  $\mathcal{C}$  is  $n \times n$  but  $\lambda_{k+1} = \dots = \lambda_n = 0$ . Then the underlying data is **low-rank**

Suppose that  $\mathcal{C}$  is  $n \times n$  but  $\lambda_k$  through  $\lambda_n$  are **very small**. Then the underlying data is **approximately low-rank**.

# Eigenfaces



The first few principal components (a.k.a. eigenvectors of the covariance matrix) for a database of many **faces**. Different components accentuate different facial characteristics

# Eigenfaces



Top left face is projection of bottom right face onto its first principal component. Each new image from left to right corresponds to using 8 additional principal components for reconstruction

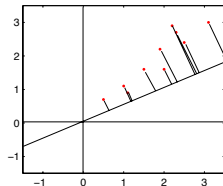
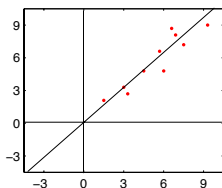
# Eigenfaces



The projections of non-face images onto first few principal components

# Reducing dimensionality using random projections





### Principal components:

Directions of projection are data-dependent

### Random projections:

Directions of projection are *independent* of the data

Why not *always* use principal components?

- 1 May not have access to all the data at once, as in **data streaming**
- 2 Computing principal components (eigenvectors) is **computationally expensive** in high dimensions:  $O(kn^2)$  'flops' to compute  $k$  principal components

# Data streaming



- Massive amounts of data arrives in small time increments
- Often past data cannot be accumulated and stored, or when they can, access is expensive.

- $\mathbf{x} = (x_1, x_2, \dots, x_n)$  at time  $(t_1, t_2, \dots, t_n)$ , and  $\tilde{\mathbf{x}} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$  at time  $(t_1 + \Delta(t), t_2 + \Delta(t), \dots, t_n + \Delta(t))$

Summary statistics that can be computed in **one pass**:

- Mean value:  $\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$
- Euclidean length:  $\|\mathbf{x}\|^2 = \sum_{j=1}^n x_j^2$
- Variance:  $\sigma^2(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2$

What about the correlation  $\langle \mathbf{x} - \bar{x}, \mathbf{y} - \bar{y} \rangle / \sigma(\mathbf{x})\sigma(\mathbf{y})$ ?

- used to assess **risk** of stock  $\mathbf{x}$  against market  $\mathbf{y}$

## Approach: introduce randomness

- Consider  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  and vector  $\mathbf{g} = (g_1, g_2, \dots, g_n)$  of **independent and identically distributed** (i.i.d.) unit normal Gaussian random variables:

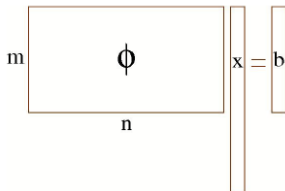
$$g_j \sim \mathcal{N}(0, 1), \quad \mathbb{P}(g_j \geq x) = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

- Consider

$$\begin{aligned} u &= \langle \mathbf{g}, \mathbf{x} \rangle - \langle \mathbf{g}, \tilde{\mathbf{x}} \rangle \\ &= (g_1 x_1 + g_2 x_2 + \dots + g_n x_n) - (g_1 \tilde{x}_1 + g_2 \tilde{x}_2 + \dots + g_n \tilde{x}_n) \\ &= \langle \mathbf{g}, \mathbf{x} - \tilde{\mathbf{x}} \rangle \end{aligned}$$

### Theorem

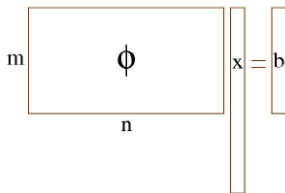
$$\mathbb{E} \langle \mathbf{g}, \mathbf{x} - \tilde{\mathbf{x}} \rangle^2 = \|\mathbf{x} - \mathbf{y}\|^2$$



- For an  $m \times N$  matrix  $\Phi$  with i.i.d. Gaussian entries  $\varphi_{i,j} \sim \mathcal{N}(0, 1)$

$$\begin{aligned}
 \mathbb{E}(\| \frac{1}{\sqrt{m}} \Phi(\mathbf{x} - \mathbf{y}) \|^2) &= \frac{1}{\sqrt{m}} \mathbb{E} \left( \sum_{i=1}^m \langle \mathbf{g}_i, \mathbf{x} - \tilde{\mathbf{x}} \rangle^2 \right) \\
 &= \| \mathbf{x} - \mathbf{y} \|^2
 \end{aligned}$$

## Approach: introduce randomness



### Concentration around expectation:

- For a fixed  $\mathbf{x} \in \mathbb{R}^n$ ,

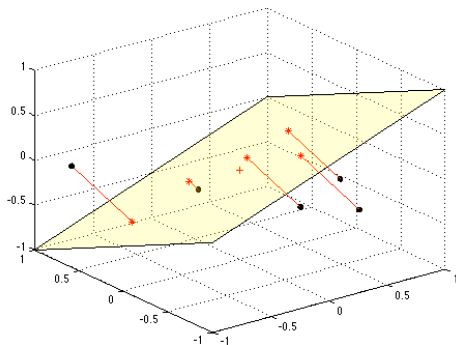
$$\mathbb{P}\left(\left\|\frac{1}{\sqrt{m}}\Phi(\mathbf{x})\right\|^2 \geq (1 + \varepsilon)\|\mathbf{x}\|^2\right) \leq \exp\left(-\frac{m}{4}\varepsilon^2\right)$$

- For  $p$  vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$  in  $\mathbb{R}^n$

$$\mathbb{P}\left(\forall \mathbf{x}_j : \left\|\frac{1}{\sqrt{m}}\Phi(\mathbf{x}_j)\right\|^2 \geq (1 + \varepsilon)\|\mathbf{x}_j\|^2\right) \leq \exp\left(\log p - \frac{m}{4}\varepsilon^2\right)$$

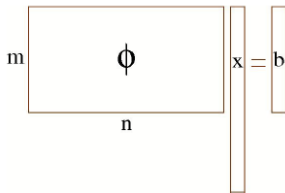
How small can  $m$  be such that this probability is still small?

# Geometric intuition



- The linear map  $\mathbf{x} \rightarrow \frac{1}{\sqrt{m}}\Phi\mathbf{x}$  is similar to a *random projection* onto an  $m$ -dimensional subspace of  $\mathbb{R}^n$
- *most* projections preserve geometry, but **not all**.

# Measure-concentration for Gaussian matrices



## Theorem (Concentration of lengths / Johnson-Lindenstrauss)

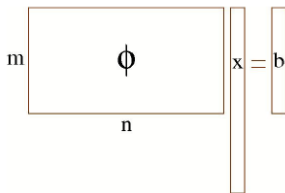
Fix an accuracy  $\varepsilon > 0$  and probability of failure  $\eta > 0$ . Fix an integer  $m \geq 10\varepsilon^{-2} \log(p)$ , and fix an  $m \times n$  Gaussian random matrix  $\Phi$ .

Then with probability greater than  $1 - \eta$ ,

$$\left| \frac{1}{\sqrt{m}} \|\Phi \mathbf{x}_j - \Phi \mathbf{x}_k\| - \|\mathbf{x}_j - \mathbf{x}_k\| \right| \leq \varepsilon \|\mathbf{x}_j - \mathbf{x}_k\|$$

for all  $j$  and  $k$ .





### Corollary (Concentration for inner products)

*Fix an accuracy  $\varepsilon > 0$  and probability of failure  $\eta > 0$ . Fix an integer  $m \geq 10\varepsilon^{-2} \log(p)$  and fix an  $m \times n$  Gaussian random matrix  $\Phi$ .*

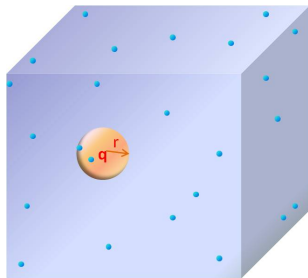
*Then with probability greater than  $1 - \eta$ ,*

$$\left\| \frac{1}{m} \langle \Phi \mathbf{x}_j, \Phi \mathbf{x}_k \rangle - \langle \mathbf{x}_j, \mathbf{x}_k \rangle \right\| \leq \frac{\varepsilon}{2} (\|\mathbf{x}_j\|^2 + \|\mathbf{x}_k\|^2)$$

*for all  $j$  and  $k$ .*

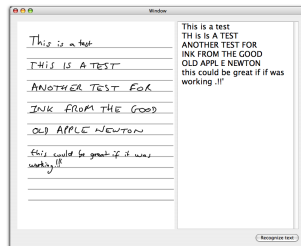
# Nearest-neighbors

# The nearest-neighbors problem



- Find the closest point to a point  $\mathbf{q}$  from among a set of points  $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$ . Originally called the “post-office problem” (1973)

# Applications



Similarity searching ...

# The nearest-neighbors problem

- Find the closest point to a point  $\mathbf{q}$  from among a set of points  $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$

■

$$\begin{aligned}\mathbf{x}^* &= \arg \min_{\mathbf{x}_j \in S} \|\mathbf{q} - \mathbf{x}_j\|^2 \\ &= \arg \min_{\mathbf{x}_j \in S} \sum_{k=1}^N (q(k) - x_j(k))^2\end{aligned}$$

- **Computational cost** (number of 'flops') per search:  $O(Np)$
- Computational cost of  $m$  searches:  $O(Nmp)$ .
- **Curse of dimensionality:** If  $N$  and  $p$  are large, this is a lot of flops!

# The $\varepsilon$ -approximate nearest-neighbors problem

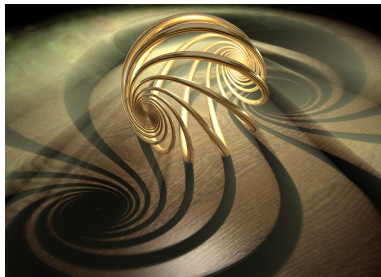
- Given a tolerance  $\varepsilon > 0$ , and a point  $\mathbf{q} \in \mathbb{R}^N$ , return a point  $\mathbf{x}_\varepsilon^*$  from the set  $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$  which is an  $\varepsilon$ -approximate nearest neighbor to  $\mathbf{q}$ :

$$\|\mathbf{q} - \mathbf{x}_\varepsilon^*\| \leq (1 + \varepsilon)\|\mathbf{q} - \mathbf{x}^*\|$$

This problem can be solved using random projections:

- Let  $\Phi$  be an  $m \times N$  Gaussian random matrix, where  $m = 10\varepsilon^{-2} \log p$ .
- Compute  $\mathbf{r} = \Phi\mathbf{q}$ . For all  $j = 1, \dots, p$ , compute  $\mathbf{x}_j \rightarrow \mathbf{u}_j = \Phi\mathbf{x}_j$ .  
**Computational cost:**  $O(Np \log(p))$ .
- Compute  $\mathbf{x}_\varepsilon^* = \arg \min_{\mathbf{x}_j \in S} \|\mathbf{r} - \mathbf{u}_j\|$ . **Computational cost:** of  $m$  searches:  $O(pm \log(p + m))$ .

**Total computation cost:**  $O((N + m)p \log(p + m)) \ll O(Np^2)$  !



# Random projections and sparse recovery

## Theorem (Subspace-preservation)

Suppose that  $\Phi$  is an  $m \times n$  random matrix with the distance-preservation property:

$$\text{For any fixed } \mathbf{x} : \quad \mathbb{P}\left(\left|\|\Phi\mathbf{x}\|^2 - \|\mathbf{x}\|^2\right| \geq \varepsilon\|\mathbf{x}\|^2\right) \leq 2e^{-c_\varepsilon m}$$

Let  $k \leq c_\varepsilon m$  and let  $T_k$  be a  $k$ -dimensional subspace of  $\mathbb{R}^n$ . Then

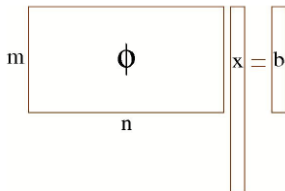
$$\mathbb{P}\left(\text{For all } \mathbf{x} \in T_k : \quad (1 - \varepsilon)\|\mathbf{x}\|^2 \leq \|\Phi\mathbf{x}\|^2 \leq (1 + \varepsilon)\|\mathbf{x}\|^2\right) \geq 1 - e^{-c'_\varepsilon m}$$

Outline of proof:

- A  $\varepsilon$ -cover and the Vitali covering lemma
- Continuity argument



# Sparse recovery and RIP



**Restricted Isometry Property of order  $k$ :**  $\Phi$  has the RIP of order  $k$  if

$$.8\|\mathbf{x}\|^2 \leq \|\Phi\mathbf{x}\|^2 \leq 1.2\|\mathbf{x}\|^2$$

for all  $k$ -sparse vectors  $\mathbf{x} \in \mathbb{R}^n$ .

## Theorem

*If  $\Phi$  has RIP of order  $k$ , then for all  $k$ -sparse vectors  $\mathbf{x}$  such that  $\Phi\mathbf{x} = \mathbf{b}$ ,*

$$\mathbf{x} = \arg \min \left\{ \sum_{j=1}^N |z(j)| \quad : \quad \Phi\mathbf{z} = \mathbf{b}, \quad \mathbf{z} \in \mathbb{R}^n \right\}$$

## Theorem (Distance-preservation implies RIP)

Suppose that  $\Phi$  is an  $m \times N$  random matrix with the subspace-preservation property:

$$\mathbb{P}\left(\exists \mathbf{x} \in T_k : (1 - \varepsilon)\|\mathbf{x}\|^2 \leq \|\Phi\mathbf{x}\|^2 \leq (1 + \varepsilon)\|\mathbf{x}\|^2\right) \leq e^{-c'_\varepsilon m}$$

Then with probability greater than .99,

$$(1 - \varepsilon)\|\mathbf{x}\|^2 \leq \|\Phi\mathbf{x}\|^2 \leq (1 + \varepsilon)\|\mathbf{x}\|^2$$

for all  $\mathbf{x}$  of sparsity level  $k \leq c_\varepsilon m / \log(N)$ .

Outline of proof:

- Bound for a fixed subspace  $T_k$ .
- Union bound over all  $\binom{N}{k} \leq N^k$  subspaces of  $k$ -sparse vectors

# Fast principal component analysis

# Principal component analysis

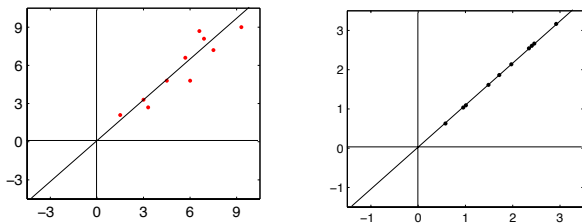


Figure: Original data and projection onto first principal component

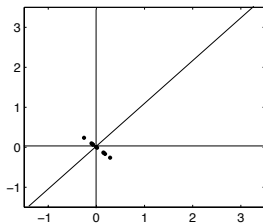


Figure: Residual

# Principal components in higher dimensions



Computing principal components is **expensive**: Use fast randomized algorithms for approximate PCA

# Randomized Principal component analysis

- First principal component is largest eigenvector  $\mathbf{v}_1 = (v_1(1), \dots, v_1(n))$  of covariance matrix  $\mathcal{C} = PDP^{-1}$ , where

$$P = \begin{bmatrix} v_1(1) & v_1(2) & \dots & v_1(n) \\ v_2(1) & v_2(2) & \dots & v_2(n) \\ \vdots & \vdots & & \vdots \\ v_n(1) & v_n(2) & \dots & v_n(n) \end{bmatrix}, \quad D = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

- 'Power method' for computing largest principal component based on observation:

If  $\mathbf{x}_0$  is a random Gaussian vector and  $\mathbf{x}_{n+1} = \mathcal{C}\mathbf{x}_n$ , then  $\mathbf{x}_n / \|\mathbf{x}_n\| \rightarrow \mathbf{v}_1$ .

# Randomized principal component analysis

- If  $\mathcal{C} = \mathcal{P}D\mathcal{P}^{-1}$  is a **rank- $k$**  (or approximately rank- $k$ ) matrix, then all principal components can be computed using  $2k$  gaussian random vectors.
- For more accurate approximate PCA, do more iterations of power method.