

LECTURE 5: Matrix Diagonalization, Singular Value Decomposition and Principal Component Analysis

- We wish to diagonalize a square $N \times N$ matrix \mathbf{A} : $\mathbf{Ax} = \lambda \mathbf{x}$.
Here \mathbf{x} is the **eigenvector** and λ the **eigenvalue** of \mathbf{A}
- We could solve $\det|\mathbf{A} - \lambda \mathbf{I}| = 0$ and expand it terms of a polynomial of N -th order in λ
- There is a shift symmetry: $(\mathbf{A} + \tau \mathbf{I})\mathbf{x} = (\lambda + \tau)\mathbf{x}$ that does not change eigenvectors, but changes eigenvalues

Common Matrix Types

- **Symmetric**: $\mathbf{A}=\mathbf{A}^T$ or $a_{ij}=a_{ji}$
- **Hermitian** (self-adjoint): $\mathbf{A}=\mathbf{A}^+$ or $a_{ij}=a_{ji}^*$. Important concept because eigenvalues are real. (quantum mechanics!)
- **Orthogonal**: transpose equals inverse $\mathbf{A}\mathbf{A}^T=\mathbf{A}^T\mathbf{A}=\mathbf{I}$
- **Unitary**: Hermitian conjugate equals inverse $\mathbf{A}\mathbf{A}^+=\mathbf{I}$
- **Normal**: commutes with Hermitian conjugate $\mathbf{A}\mathbf{A}^+=\mathbf{A}^+\mathbf{A}$. Important because its eigenvectors are complete and orthogonal.
- For real matrices Hermitian is symmetric and unitary is orthogonal, both of which are normal.
- A normal non-symmetric matrix is hard to diagonalize.
- **Symmetric matrices are easier to diagonalize.**

Matrix Forms and Similarity Transform

- We can define the matrix of right eigenvectors as \mathbf{X}_R and the matrix of eigenvalues as \mathbf{D} , which is diagonal with λ_i on the diagonal:

$$\mathbf{A}\mathbf{X}_R = \mathbf{X}_R\mathbf{D}$$

- We can also define left eigenvectors $\mathbf{X}_L \mathbf{A} = \mathbf{D} \mathbf{X}_L$
- $\mathbf{X}_L \mathbf{X}_R \mathbf{D} = \mathbf{D} \mathbf{X}_L \mathbf{X}_R$ implies $\mathbf{X}_L \mathbf{X}_R = \mathbf{I}$ (with appropriate normalization) and so $\mathbf{A} = \mathbf{X}_R \mathbf{D} \mathbf{X}_L = \mathbf{X}_R \mathbf{D} \mathbf{X}_R^{-1}$ and $\mathbf{D} = \mathbf{X}_L \mathbf{A} \mathbf{X}_R = \mathbf{X}_R^{-1} \mathbf{A} \mathbf{X}_R$
- For a symmetric matrix $\mathbf{A} = \mathbf{X}_R \mathbf{D} \mathbf{X}_R^T$
- A transformation $\mathbf{Z}^{-1} \mathbf{A} \mathbf{Z}$ leaves eigensystem unchanged:
$$\det| \mathbf{Z}^{-1} \mathbf{A} \mathbf{Z} - \lambda \mathbf{I} | = \det| \mathbf{Z}^{-1} \mathbf{A} \mathbf{Z} - \lambda \mathbf{Z}^{-1} \lambda \mathbf{Z} | =$$
$$\det| \mathbf{Z}^{-1} | \det| \mathbf{A} - \lambda \mathbf{I} | \det| \mathbf{Z} | = \det| \mathbf{A} - \lambda \mathbf{I} |$$
- Many of the methods consist of a series of such transforms
- For real symmetric matrices $\mathbf{Z}^{-1} = \mathbf{Z}^T$

QR Decomposition

A: a set of **N** column vectors a_0, \dots, a_{N-1}

$$A = \begin{pmatrix} | & | & | & \dots & | \\ \mathbf{a}_0 & \mathbf{a}_1 & \mathbf{a}_2 & \dots & a_{N-1} \\ | & | & | & \dots & | \end{pmatrix}$$

Define: $u_0 = a_0, \quad q_0 = \frac{u_0}{|u_0|}$

$$u_1 = a_1 - (q_0 \cdot a_1)q_0, \quad q_1 = \frac{u_1}{|u_1|}$$

$$u_2 = a_2 - (q_0 \cdot a_2)q_0 - (q_1 \cdot a_2)q_1, \quad q_2 = \frac{u_2}{|u_2|}$$

:

$$u_i = a_i - \sum_{j=0}^{i-1} (q_j \cdot a_i)q_j, \quad q_i = \frac{u_i}{|u_i|}$$

We can show (by induction) that $q_i q_j = \delta_{ij}$

Note that this is closely related to **Gram-Schmidt orthogonalization**

Gram-Schmidt orthogonalization and QR/RQ Decomposition

- The RQ decomposition transforms a matrix \mathbf{A} into the product of an upper triangular matrix \mathbf{R} (also known as right-triangular) and an orthogonal matrix \mathbf{Q} . The only difference from QR decomposition is the order of these matrices.
- QR decomposition is Gram–Schmidt orthogonalization of columns of \mathbf{A} , started from the first column.
- RQ decomposition is Gram–Schmidt orthogonalization of rows of \mathbf{A} , started from the last row.

QR Decomposition

$$a_0 = |u_0|q_0, \quad a_1 = |u_1|q_1 + (q_0 \cdot a_1)q_0$$

$$a_2 = |u_2|q_2 + (q_0 \cdot a_2)q_0 + (q_1 \cdot a_2)q_1 \dots$$

$$\mathbf{A} = \begin{pmatrix} | & | & | & \dots \\ | & | & | & \dots \\ \mathbf{a}_0 & \mathbf{a}_1 & \mathbf{a}_2 & \dots \\ | & | & | & \dots \end{pmatrix} = \underbrace{\begin{pmatrix} | & | & | & \dots \\ | & | & | & \dots \\ \mathbf{q}_0 & \mathbf{q}_1 & \mathbf{q}_2 & \dots \\ | & | & | & \dots \end{pmatrix}}_{\mathbf{Q}} \cdot \underbrace{\begin{pmatrix} |u_0| & q_0 \cdot a_1 & q_0 \cdot a_2 & \dots \\ 0 & |u_1| & q_1 \cdot a_2 & \dots \\ 0 & 0 & |u_1| & \dots \\ \vdots & \vdots & \vdots & \dots \end{pmatrix}}_{\mathbf{R}}$$

- We have obtained a decomposition into an orthogonal matrix **Q** and upper diagonal matrix **R**

QR Decomposition for symmetric A

- We can repeat these QR steps

$$A = Q_1 R_1, \quad Q_1^T A = Q_1^T Q_1 R_1 = R_1$$

$$A_1 \equiv R_1 Q_1 = Q_1^T A Q_1 \quad \text{Similarity Transform}$$

Repeat

$$A_1 = Q_2 R_2, \quad A_2 \equiv R_2 Q_2 = Q_2^T A_1 Q_2 = Q_2^T Q_1^T A Q_1 Q_2$$

$$A_3 = Q_3^T Q_2^T Q_1^T A Q_1 Q_2 Q_3$$

:

A_k converges to \mathbf{D}

$$X_R = Q_1 Q_2 Q_3 \dots Q_k$$

$$D = X_R^T A X_R \rightarrow A X_R = X_R D$$

General Packaged Solutions

- QR works better if the matrix \mathbf{A} is **first transformed into tridiagonal form** (Householder or Givens method): standard method for symmetric \mathbf{A}
- **Diagonalization** methods: there are many different methods depending on what matrix we have: \mathbf{A} can be real banded symmetric, real symmetric, complex Hermitian, real non-symmetric...
- Depending of what we need: just the eigenvectors or both eigenvectors and eigenvalues
- Look at **numpy.linalg** and choose depending on what your needs are: eigh, eigvalsh...

Generalized Problems

- $\mathbf{Ax} = \lambda \mathbf{Bx}$ can be handled as $(\mathbf{B}^{-1}\mathbf{A})\mathbf{x} = \lambda \mathbf{Bx}$
- If \mathbf{A} and \mathbf{B} symmetric \mathbf{B}^{-1} is not, but we can use Cholesky decomposition $\mathbf{B} = \mathbf{LL}^T$ and by multiplying with \mathbf{L}^{-1} we get $\mathbf{C}(\mathbf{L}^T\mathbf{x}) = \lambda \mathbf{L}^T\mathbf{x}$ where $\mathbf{C} = \mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^{-1})^T$ is a symmetric matrix. This diagonalization gives the same eigenvalues with eigenvectors given by $\mathbf{L}^T\mathbf{x}$.
- Nonlinear problems: $(\mathbf{A}\lambda^2 + \mathbf{B}\lambda + \mathbf{C})\mathbf{x} = 0$ can be solved by solving $2N \times 2N$ problem

$$\begin{pmatrix} 0 & I \\ -A^{-1}C & -A^{-1}B \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \lambda \begin{pmatrix} x \\ y \end{pmatrix}$$

“Shortcomings” of Diagonalization

- Eigenvectors are not orthogonal for a non-symmetric $N \times N$ matrix
- There may not be enough eigenvectors, i.e. they may not be complete: \mathbf{A} is “defective”
- It cannot be defined as eigenvalue problem $\mathbf{Ax} = \lambda \mathbf{x}$ unless \mathbf{A} is $N \times N$
- All these “defects” can be cured by SVD
- None of these “defects” apply to a symmetric $N \times N$ matrix: SVD and diagonalization are the same in that case if \mathbf{A} is semi-positive definite

Singular Value Decomposition (SVD)

- Is a decomposition of a general $N \times M$ matrix
- Define rank r of matrix \mathbf{A} as the maximum number of linearly independent column vectors in the **matrix** or the maximum number of linearly independent row vectors in the **matrix** (both **definitions** are equivalent)
- Another way: define range as the subspace of \mathbf{b} reached by $\mathbf{Ax} = \mathbf{b}$. Rank is its dimension.
- Define two orthonormal bases, with vectors \mathbf{u} in R^M and \mathbf{v} in R^N . In matrix form \mathbf{U} is $M \times M$ and \mathbf{V} is $N \times N$

Singular Value Decomposition (SVD)

- $\mathbf{u}_1, \dots, \mathbf{u}_r$ is an orthonormal basis for the column space
- $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ is an orthonormal basis for the left null space N (defined as $\mathbf{A}^T \mathbf{x} = 0$)
- $\mathbf{v}_1, \dots, \mathbf{v}_r$ is an orthonormal basis for the row space
- $\mathbf{V}_{r+1}, \dots, \mathbf{V}_N$ is an orthonormal basis for the null space $N(\mathbf{A}\mathbf{x}=0)$.
- These bases “diagonalize” \mathbf{A} : $\mathbf{A}\mathbf{v}_1 = \sigma_1 \mathbf{u}_1$, $\mathbf{A}\mathbf{v}_2 = \sigma_2 \mathbf{u}_2 \dots$
- Also called Karhunen-Loeve (KL) Transform

SVD continued

$$(m \text{ by } n)(n \text{ by } r) \quad \mathbf{AV}_r = \mathbf{U}_r \Sigma_r \quad A \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_r \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_r \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \ddots \\ \sigma_r \end{bmatrix}$$
$$(m \text{ by } r)(r \text{ by } r)$$

- Add null space, which is already orthogonal to first r \mathbf{v} 's and \mathbf{u} 's to go from reduced to full SVD

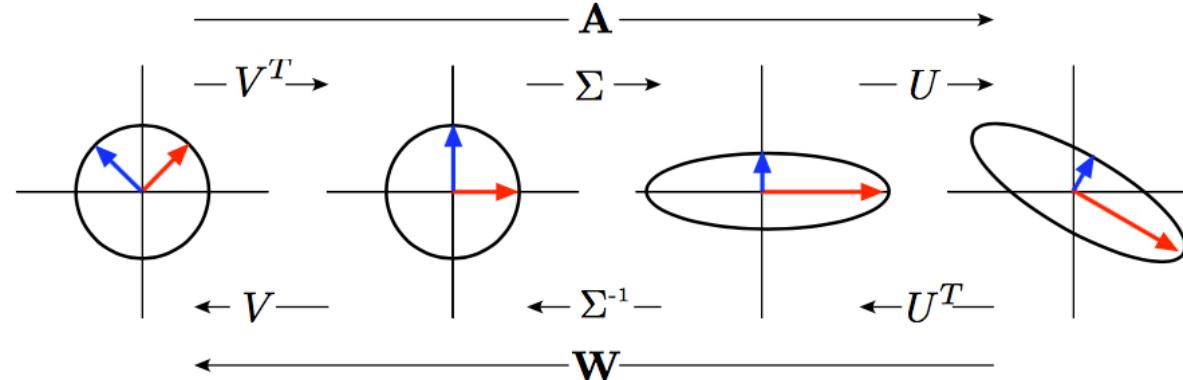
$$(m \text{ by } n)(n \text{ by } n) \quad \mathbf{Av} \text{ equals } \mathbf{U}\Sigma \quad A \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_r & \cdots & \mathbf{v}_n \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_r & \cdots & \mathbf{u}_m \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \ddots \\ \sigma_r \end{bmatrix}$$
$$(m \text{ by } m)(m \text{ by } n)$$

- Now the last matrix Σ is $M \times N$ with first r on the diagonal s_i , and the rest 0. \mathbf{U} and \mathbf{V} are now square matrices ($M \times M$ and $N \times N$) hence $\mathbf{AV} = \mathbf{U}\Sigma$ or $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T = \mathbf{u}_1\sigma_1\mathbf{v}_1 + \dots + \mathbf{u}_r\sigma_r\mathbf{v}_r$
- Rank order $\sigma_1 > \sigma_2 > \dots > \sigma_r$

SVD in Pictures

- SVD as a sequence of rotation-stretch-rotation $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$

- $\mathbf{AV} = \mathbf{U}\Sigma$



$$\Sigma \equiv \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ 0 & & \sigma_r & 0 \\ & & & \ddots & 0 \end{bmatrix} \quad \begin{pmatrix} m \\ n \end{pmatrix} \times \begin{pmatrix} m \\ m \end{pmatrix} = \begin{pmatrix} n \\ n \end{pmatrix} \times \begin{pmatrix} n \times m \\ \text{diagonal} \\ 0 \\ 0 \end{pmatrix}$$

- $\mathbf{AV}_1 = \sigma_1 \mathbf{u}_1$

$$\begin{pmatrix} m \\ n \end{pmatrix} \times \begin{pmatrix} m \\ m \end{pmatrix} = \begin{pmatrix} \text{positive number} \\ n \end{pmatrix}$$

Example

- If $\mathbf{A} = \mathbf{x}\mathbf{y}^T$ (rank 1) with unit vectors \mathbf{x} and \mathbf{y} , what is the SVD of \mathbf{A} ?

Example

- If $\mathbf{A} = \mathbf{x}\mathbf{y}^T$ (rank 1) with unit vectors \mathbf{x} and \mathbf{y} , what is the SVD of \mathbf{A} ?

Solution:

- The reduced SVD is exactly $\mathbf{x}\mathbf{y}^T$, with rank $r = 1$. It has $\mathbf{u}_1 = \mathbf{x}$ and $\mathbf{v}_1 = \mathbf{y}$ and $\sigma_1 = 1$. For the full SVD, complete $\mathbf{u}_1 = \mathbf{x}$ to an orthonormal basis of \mathbf{u} 's, and complete $\mathbf{v}_1 = \mathbf{y}$ to an orthonormal basis of \mathbf{v} 's. No new σ 's, only $\sigma_1 = 1$

How to Construct u's and v's?

- v's will be orthonormal eigenvectors of $\mathbf{A}^T\mathbf{A}$:
- $\mathbf{A}^T\mathbf{A} = (\mathbf{U}\Sigma\mathbf{V}^T)^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^T\Sigma\mathbf{V}^T$
- We see that $\sigma^2=\lambda$ of $\mathbf{A}^T\mathbf{A}$.
- u's can easily be computed from $\mathbf{A}\mathbf{v}_i = \sigma_i\mathbf{u}_i$. They are orthogonal:

$$\mathbf{u}_i^T \mathbf{u}_j = \left(\frac{\mathbf{A}\mathbf{v}_i}{\sigma_i} \right)^T \left(\frac{\mathbf{A}\mathbf{v}_j}{\sigma_j} \right) = \frac{\mathbf{v}_i^T \mathbf{A}^T \mathbf{A} \mathbf{v}_j}{\sigma_i \sigma_j} = \frac{\sigma_j^2}{\sigma_i \sigma_j} \mathbf{v}_i^T \mathbf{v}_j = \mathbf{zero}.$$

- Show that u's are eigenvectors of $\mathbf{A}\mathbf{A}^T$

How to Construct u's and v's?

- v's will be orthonormal eigenvectors of $\mathbf{A}^T\mathbf{A}$:
- $\mathbf{A}^T\mathbf{A} = (\mathbf{U}\Sigma\mathbf{V}^T)^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^T\Sigma\mathbf{V}^T$
- We see that $\sigma^2=\lambda$ of $\mathbf{A}^T\mathbf{A}$.
- u's can easily be computed from $\mathbf{Av}_i = \sigma_i \mathbf{u}_i$. They are orthogonal:

$$\mathbf{u}_i^T \mathbf{u}_j = \left(\frac{\mathbf{Av}_i}{\sigma_i} \right)^T \left(\frac{\mathbf{Av}_j}{\sigma_j} \right) = \frac{\mathbf{v}_i^T \mathbf{A}^T \mathbf{A} \mathbf{v}_j}{\sigma_i \sigma_j} = \frac{\sigma_j^2}{\sigma_i \sigma_j} \mathbf{v}_i^T \mathbf{v}_j = \mathbf{zero}.$$

- Show that u's are eigenvectors of $\mathbf{A}\mathbf{A}^T$
- Soln: $\mathbf{A}\mathbf{A}^T = \mathbf{U}\Sigma\mathbf{V}^T(\mathbf{U}\Sigma\mathbf{V}^T)^T = \mathbf{U}\Sigma\mathbf{V}^T\mathbf{V}\Sigma^T\mathbf{U}^T = \mathbf{U}\Sigma\Sigma^T\mathbf{U}^T$

Linear Algebra with SVD

- Linear algebra $\mathbf{Ax} = \mathbf{b}$, assume \mathbf{A} is $N \times N$
- Let's do matrix inversion: $\mathbf{A}^{-1} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T$
- This can be solved unless some $\sigma_i = 0$: matrix is singular, i.e. there is non-zero null space
- More generally we can define condition number σ_1/σ_N : if this is large the matrix is ill-conditioned
- SVD will diagnose the condition situation of the matrix
- If we want to solve homogeneous equation $\mathbf{Ax} = \mathbf{b}$ for $\mathbf{b} = 0$ we use the null space of \mathbf{U}/\mathbf{V} , corresponding to $\sigma_i = 0$

Solving Singular Matrix Problems

- $\mathbf{Ax} = \mathbf{b}$ can be solved if \mathbf{b} is in the range of \mathbf{A} . If so we have an infinite number of solutions, since we can add null space solutions to it in any linear combination
- We can put additional constraint, such as the norm $\|\mathbf{x}\|$ is minimized (i.e. pick the solution with the shortest length $|\mathbf{x}|^2$). Then we can set the null space solutions to 0. This is achieved by setting $\sigma_i^{-1}=0$ wherever $\sigma_i=0$ (null space) and solve
$$\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b}$$
- If \mathbf{b} is not in range we will not have an exact solution. But we can still use $\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b}$ as the solution which minimizes
$$r = |\mathbf{Ax}-\mathbf{b}|$$

Proofs that $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{b}$

- **b is in range:** suppose \mathbf{x}' is in null space, and Σ^{-1} is the modified inverse of Σ :

$$|\mathbf{x} + \mathbf{x}'| = |\mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{b} + \mathbf{x}'| = |\mathbf{V}(\Sigma^{-1}\mathbf{U}^T\mathbf{b} + \mathbf{V}^T\mathbf{x}')| = |\Sigma^{-1}\mathbf{U}^T\mathbf{b} + \mathbf{V}^T\mathbf{x}'| > |\mathbf{x}|$$

unless $\mathbf{x}' = 0$ because the two are orthogonal (\mathbf{x}' is in null space, \mathbf{x} is not).

- **b is not in range:** we add \mathbf{x}' to \mathbf{x} , $\mathbf{Ax}-\mathbf{b}$ is modified by $\mathbf{b}'=\mathbf{Ax}'$ (in range), show $r = |\mathbf{Ax}-\mathbf{b}|$ is minimized:

$$\begin{aligned} |\mathbf{Ax}-\mathbf{b}+\mathbf{b}'| &= |(\mathbf{U}\Sigma\mathbf{V}^T)\mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{b}-\mathbf{b}+\mathbf{b}'| \\ &= |(\mathbf{U}\Sigma\Sigma^{-1}\mathbf{U}^T-\mathbf{I})\mathbf{b}+\mathbf{b}'| = |\mathbf{U}[(\Sigma\Sigma^{-1}-\mathbf{I})\mathbf{U}^T\mathbf{b}+\mathbf{U}^T\mathbf{b}']| \\ &= |(\Sigma\Sigma^{-1}-\mathbf{I})\mathbf{U}^T\mathbf{b}+\mathbf{U}^T\mathbf{b}'| > |(\Sigma\Sigma^{-1}-\mathbf{I})\mathbf{U}^T\mathbf{b}| \text{ unless } \mathbf{b}'=0. \end{aligned}$$

- In practice we want to fix ill-conditioned matrices with SVD: throw away all singular values whose value is too small to be meaningful (machine precision roundoff errors etc.). So if $\sigma_i < \varepsilon$ we set $\sigma_i^{-1} = 0$, where ε is machine precision.

Non-square Matrices

- If there are fewer equations N than variables M the system is **underdetermined** and we have $M-N$ dimensional family of solutions. SVD provides the family of solutions by providing $\sigma_i < \varepsilon$.
- If there are more equations N than variables M the system is **overdetermined** and we can only obtain the least square solution, which is given by $\mathbf{x} = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^T\mathbf{b}$ (same proof). There will be no zeros of σ_i unless there are hidden degeneracies.

Key Ideas of SVD

- The SVD factors \mathbf{A} into $\mathbf{U}\Sigma\mathbf{V}^T$, with r singular values $\sigma_1 \geq \dots \geq \sigma_r > 0$.
- The numbers $\sigma_1^2, \dots, \sigma_r^2$ are the nonzero eigenvalues of $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$.
- The orthonormal columns of \mathbf{U} and \mathbf{V} are eigenvectors of $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$.
- Those columns hold orthonormal bases for the four fundamental subspaces of \mathbf{A} .
- Those bases diagonalize the matrix: $\mathbf{A}\mathbf{v}_i = \sigma_i \mathbf{u}_i$ for $i \leq r$. This is $\mathbf{AV} = \mathbf{U}\Sigma$.
- $\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$ and σ_1 is the maximum of the ratio $\|\mathbf{Ax}\|/\|\mathbf{x}\|$.

Applications of SVD

- Diagnosis of matrix condition number
- Construction of orthonormal basis (we can also do that with QR/RQ)
- Lower rank approximation of matrices

$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T = \mathbf{u}_1\sigma_1\mathbf{v}_1 + \dots + \mathbf{u}_k\sigma_k\mathbf{v}_k$, or $a_{ij} = \sum_{k=1}^N \sigma_k u_{ik} v_{jk}$ where $k < r$.

The challenge is that to construct this we need r^3 process, which can be expensive

From Lecture 3: Multivariate Linear Least Squares

- We can generalize the model to a generic functional form

$$y_i = a_0 X_0(x_i) + a_1 X_1(x_i) + \dots + a_{M-1} X_{M-1}(x_i)$$

- The problem is linear in a_j and can be nonlinear in x_i ,

e.g. $X_j(x_i) = x_i^j$

$$\chi^2 = \sum_{i=0}^{N-1} \left[\frac{y_i - \sum_{k=0}^{M-1} a_k X_k(x_i)}{\sigma_i} \right]^2$$

- We can define design matrix $A_{ij} = X_j(x_i)/\sigma_i$ and

- $b_i = y_i/\sigma_i$

$$\chi^2 = |\mathbf{A} \cdot \mathbf{a} - \mathbf{b}|^2$$

Solving Least Squares with SVD

- We want to find \mathbf{a} for which $\chi^2 = |\mathbf{A}\mathbf{a} - \mathbf{b}|^2$ is minimized.
- $\mathbf{a} = \sum_{i=1}^M (\mathbf{U}_i \mathbf{b} / \sigma_i) \mathbf{V}_i$
- The errors are $\text{Var}(a_j) = \sum_{i=1}^M (V_{ji}^2 / \sigma_i^2)$
- Note that this is no different than diagonalizing precision matrix $\mathbf{A}^T \mathbf{A} = \mathbf{V} \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} \mathbf{V}^T$, where $(\boldsymbol{\Sigma}^T \boldsymbol{\Sigma})_{ii} = \sigma_i^2$ and $(\boldsymbol{\Sigma}^T \boldsymbol{\Sigma})_{ij} = 0$ for j not equal to i .
- This is because the inverse precision matrix is covariance matrix, $(\mathbf{A}^T \mathbf{A})^{-1} = \mathbf{V} (\boldsymbol{\Sigma}^T \boldsymbol{\Sigma})^{-1} \mathbf{V}^T$,
so $(\mathbf{A}^T \mathbf{A})_{jj} = \sum_{i=1}^M (V_{ji}^2 / \sigma_i^2)$ and $(\mathbf{A}^T \mathbf{A})_{jk} = \sum_{i=1}^M (V_{ji} V_{ik} / \sigma_i^2)$

Decorrelating the Variables

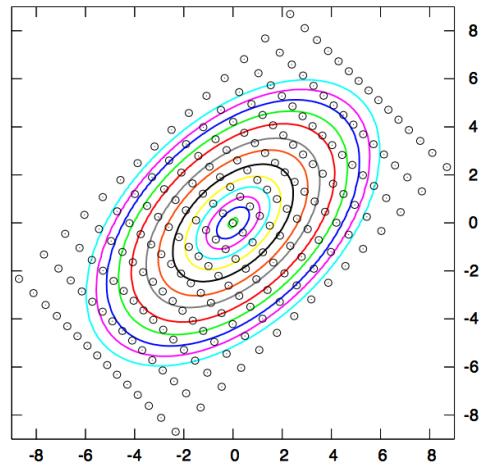
- We can decorrelate the variables using spectral principal axis rotation (diagonalization) One way to do so is to use $\alpha = \mathbf{X}_L^T \mathbf{D} \mathbf{X}_L$
- eigenvectors, columns of \mathbf{X}_L corresponding to the non-zero eigenvalues λ_i of precision matrix, to define new variables as the linear combinations of original variables: $\mathbf{u} = \mathbf{X}_L \mathbf{a}$
- the variables u_j are uncorrelated with variance λ_i^{-1}

Decorrelating the Variables

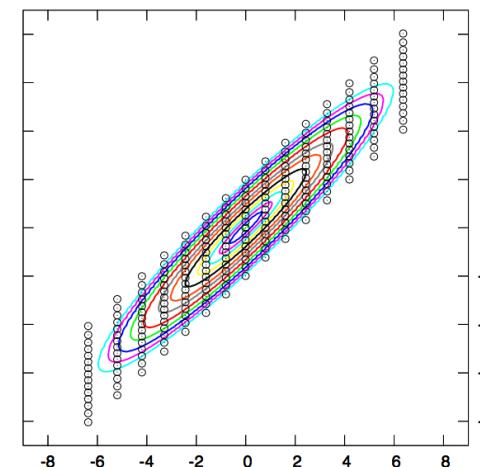
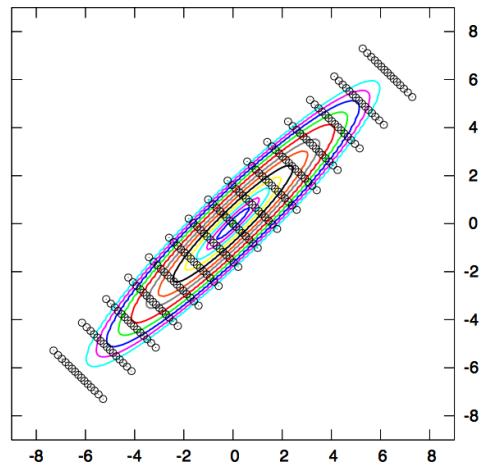
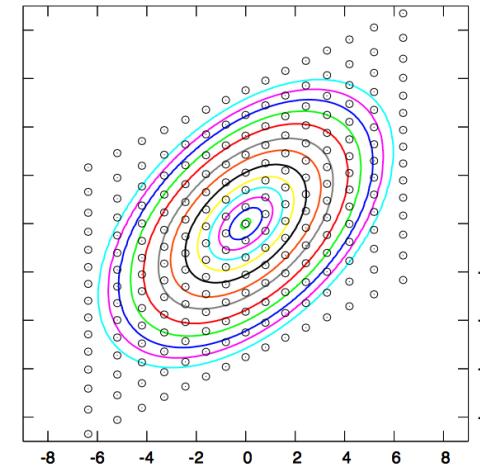
- Posterior becomes proportional to
$$\exp(-\delta \mathbf{a}^T \boldsymbol{\alpha} \delta \mathbf{a} / 2) = \exp(-\delta \mathbf{a}^T \mathbf{X}_L^T \mathbf{D} \mathbf{X}_L \delta \mathbf{a} / 2) = \exp(-\delta \mathbf{u}^T \mathbf{D} \delta \mathbf{u} / 2)$$
- This is now a product of independent variables
- $\delta \mathbf{a}^T \mathbf{X}_L^T \mathbf{D} \mathbf{X}_L \delta \mathbf{a} = \delta \mathbf{u}^T \mathbf{D} \delta \mathbf{u}$
- We can further rescale (normalize) $\mathbf{u}'_i = \mathbf{u}_i \lambda_i^{-1/2}$
- $\delta \mathbf{a}^T \mathbf{X}_L^T \mathbf{D} \mathbf{X}_L \delta \mathbf{a} = \delta \mathbf{u}'^T \delta \mathbf{u}'$
- We can also perform principal axis rotation using Cholesky
 $\mathbf{a} = \mathbf{L}^T \mathbf{L}$: $\delta \mathbf{a}^T \mathbf{L}^T \mathbf{L} \delta \mathbf{a} = \delta \mathbf{u}''^T \delta \mathbf{u}''$, where $\delta \mathbf{u}'' = \mathbf{L} \delta \mathbf{a}$
- The two are not the same

Spectral (eigen)-decomposition versus Cholesky decomposition: uniform sampling

Spectral (eigenvectors)



Cholesky



Solving Least Squares with SVD

- When normal equations are close to singular we have an under-determined system. This happens when the variables are strongly correlated.
- This is ill-conditioned problem and SVD is its cure: we want to find \mathbf{a} for which $\chi^2 = |\mathbf{A}\mathbf{a} - \mathbf{b}|^2$ is minimized.
- $\mathbf{a} = \sum_{i=1}^M (\mathbf{U}_i \mathbf{b} / \sigma_i) \mathbf{V}_i$ where we set, if $\sigma_i < \varepsilon$, $\sigma_i^{-1} = 0$.
- This is equivalent to eliminating some parameters and replacing them with a single linear combination, which is better constrained.
- Note that this is no different than diagonalizing $\mathbf{\alpha} = \mathbf{C}^{-1} = \mathbf{A}^T \mathbf{A}$ and doing the same procedure there

Singular Covariance Matrix

- What do we do if the covariance matrix is singular, i.e. some eigenvalues of precision matrix are nearly 0?
- If we invert precision matrix we will get infinity in covariance matrix. This is just a statement that we have too many variables to determine them all
- We need to “regularize” our solution by reducing the number of parameters
- One way to do so is to use eigenvectors, columns of \mathbf{X}_L corresponding to the non-zero eigenvalues λ_i of precision matrix, to define new variables as the linear combinations of original variables: $\mathbf{u} = \mathbf{X}_L \mathbf{a}$
- the variables u_j are uncorrelated with variance λ_i^{-1}
- Posterior becomes proportional to
$$\exp(-\delta \mathbf{a}^T \alpha \delta \mathbf{a} / 2) = \exp(-\delta \mathbf{a}^T \mathbf{X}_L^T D \mathbf{X}_L \delta \mathbf{a} / 2) = \exp(-\delta \mathbf{u}^T D \delta \mathbf{u} / 2)$$
- This is now a product of independent variables, but fewer than original

Principal component analysis as data analysis tool

- We do not need to remove just singular values with $\sigma_i < \varepsilon$. If we have strong ordering of singular values we can just keep the top few.
- Suppose we have some set of data measured as a function of some coordinate x_i . As an example, in HW5 we give data of quasar (QSO) light intensity as a function of wavelength for several QSOs. This is called QSO spectrum.
- We measure the data for several samples (e.g. QSOs) and we would like to describe individual variations seen in the data. For example, most QSOs look very similar in their light intensity as a function of wavelength, but there are differences in regions of strong metal line emission.
- We first normalize all the QSO spectra (so that the mean intensity is the same) and then define the mean QSO spectrum by averaging over all spectra. We subtract this mean from each spectrum. We can form $N \times M$ matrix \mathbf{A} , where N is the number of data vectors (QSOs) and M is the number of QSO pixels.
- Now we want to describe the variations from the mean. It is possible that these will be described with a set of a few simple 1-d spectra. To test this hypothesis we perform PCA. We want to minimize $\chi^2 = |\mathbf{A}\mathbf{a} - \mathbf{b}|^2$, where \mathbf{b} is the data vector for each QSO and \mathbf{a} are the coefficients. Instead of the trivial solution $a_i = 1$ we do this by reducing the rank of \mathbf{A} with SVD.

What is PCA?

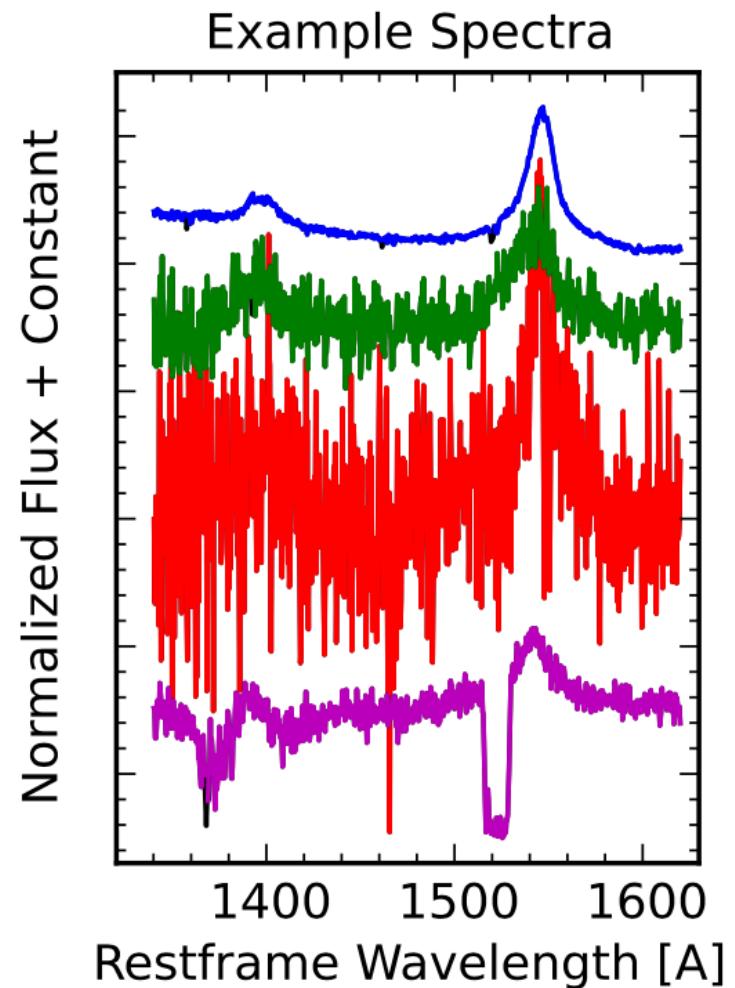
- PCA is diagonalizing $\mathbf{A}^T \mathbf{A}$ or performing SVD of \mathbf{A} and use the top few eigenvalues/principal components as a solution. We do this to reduce the rank of the matrix: **dimensionality reduction**. Note that for QSO spectra $\mathbf{A}^T \mathbf{A}$ is $M \times M$, and $(\mathbf{A}^T \mathbf{A})_{ij}$ measures the correlation between pixel i and pixel j of the spectra. If N is large SVD may be too expensive, so we do PCA by diagonalizing $\mathbf{A}^T \mathbf{A}$.
- The idea is that **if the correlation matrix is of low rank it can be approximated by a few largest eigenvectors only**. If so this will become clear by the pattern of singular values (which are sqrt of eigenvalues of $\mathbf{A}^T \mathbf{A}$): a few will be large and the rest will be small and relatively constant (because uncorrelated noise produces M eigenvalues of equal value).

What is PCA?

- With the PCA eigenvectors we can fit each QSO spectrum to these by taking a dot product between the data and the eigenvector to get the coefficient. Then we approximate the spectrum with a low rank version.
- It is non-parametric since we do not fit the data with some parametrized model. It is a linear method.

QSO Example

- We want to identify variations in QSO spectrum shapes
- Noise is a problem: for now we assume noise is the same for all spectra
- Metal absorptions are a problem: very non-Gaussian noise (outliers)
- Incomplete data are a problem
- Welcome to the world of real data analysis!

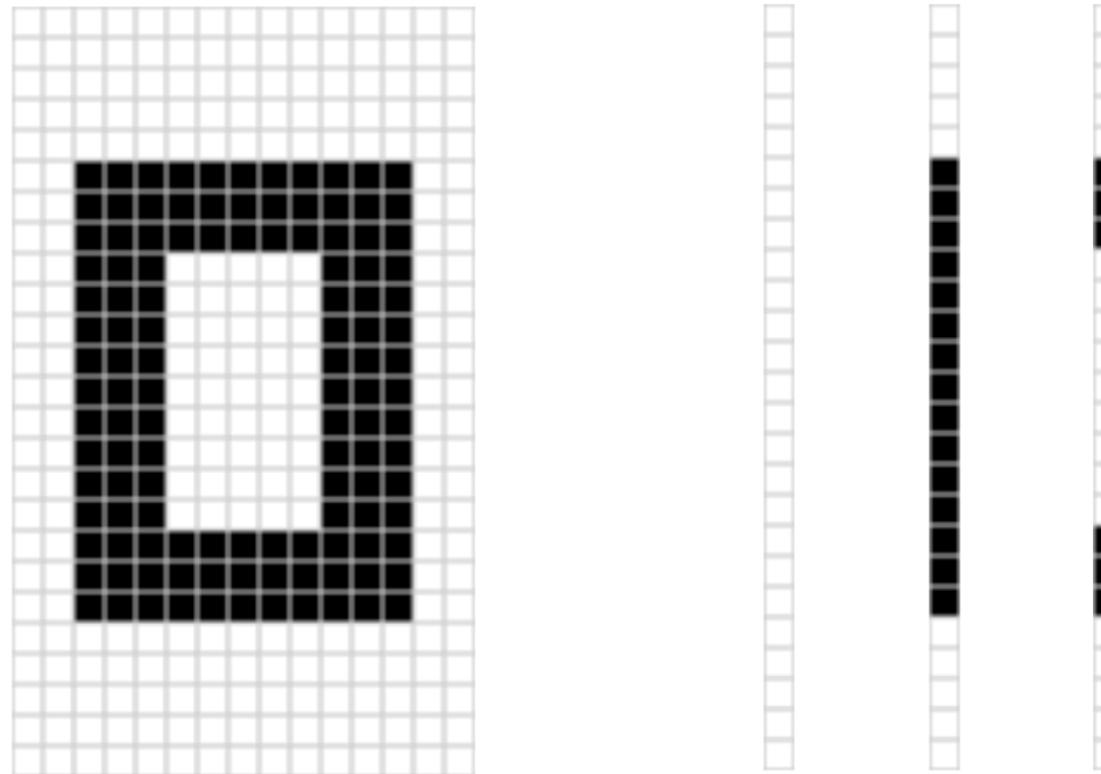


Noise in PCA

- If noise is uncorrelated it adds a diagonal to $\mathbf{A}^T \mathbf{A}$.
- If in the absence of noise the data were described with a few PCAs only, in the presence of noise all PCAs may be non-zero, but higher PCAs are meaningless since they are all noise generated
- We look at PCA values and observe where they stop decreasing significantly, or maybe use a large jump in value

An image reduction example: 15x25 block

- This block has only 3 column (row) patterns. We can do covariance analysis of columns: for any i and j column (15) we add products of ik and jk pixels over all k (25), but this can give only 3x3 possible products



We can represent it
with 0's and 1's

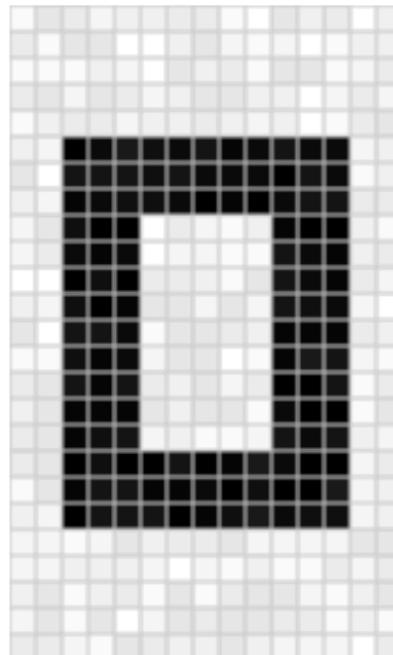
Let's perform
SVD of \mathbf{M} or
diagonalization of
 $\mathbf{M}^T \mathbf{M}$

We find 3 singular components:

$$\begin{aligned}\sigma_1 &= 14.72 \\ \sigma_2 &= 5.22 \\ \sigma_3 &= 3.31\end{aligned}$$

- We can write $M = \mathbf{u}_1 \sigma_1 \mathbf{v}_1^T + \mathbf{u}_2 \sigma_2 \mathbf{v}_2^T + \mathbf{u}_3 \sigma_3 \mathbf{v}_3^T$

- What if we have noise? We get



$$\begin{aligned}\sigma_1 &= 14.15 \\ \sigma_2 &= 4.67 \\ \sigma_3 &= 3.00 \\ \sigma_4 &= 0.21 \\ \sigma_5 &= 0.19 \\ &\dots \\ \sigma_{15} &= 0.05\end{aligned}$$

- Let us truncate at $r = 3$

$$M \approx \mathbf{u}_1\sigma_1 \mathbf{v}_1^\top + \mathbf{u}_2\sigma_2 \mathbf{v}_2^\top + \mathbf{u}_3\sigma_3 \mathbf{v}_3^\top$$

noisy

reconstructed

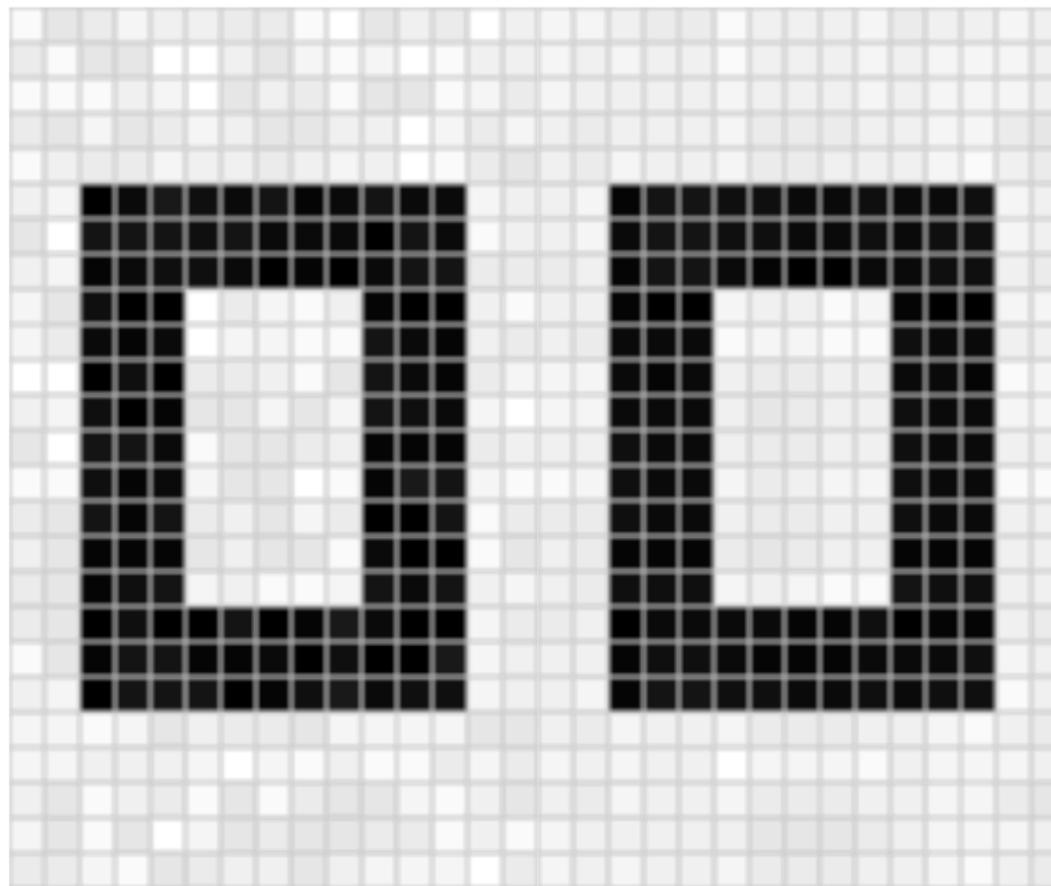


Image compression: $A=600 \times 600$ SVD (not a natural application of SVD)

Original



10 singular values



Image compression: $A=600 \times 600$ SVD (not a natural application of SVD)

50 SV

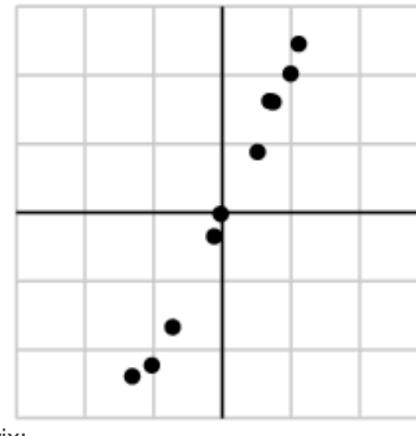


100 SV



Data Analysis

- We can do data analysis with SVD/PCA
- We have collected these data:



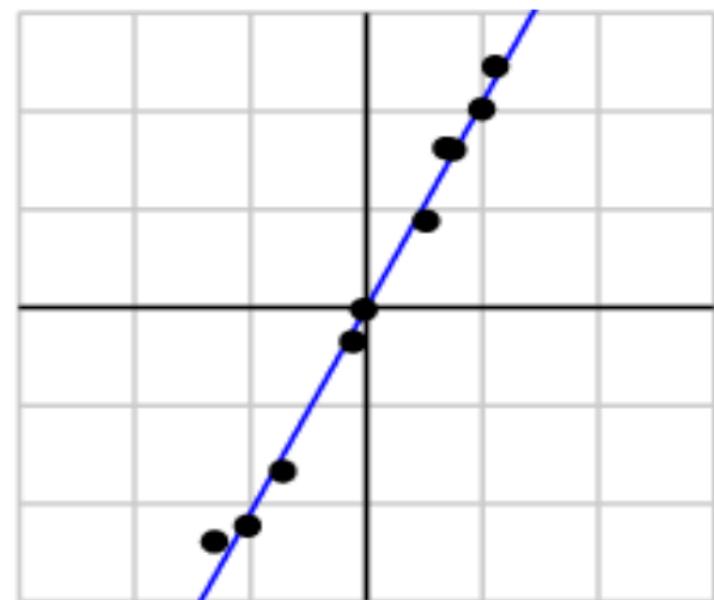
- Let us collect the data into 2x10 matrix

```
-1.03 0.74 -0.02 0.51 -1.31 0.99 0.69 -0.12 -0.72 1.11  
-2.23 1.61 -0.02 0.88 -2.39 2.02 1.62 -0.35 -1.67 2.46
```

- SVD gives $\sigma_1 = 6.04$ $\sigma_2 = 0.22$ so one value is large, second may be noise

What is SVD in this case?

- The matrix $A^T A$ is given by
$$\begin{pmatrix} \langle x^2 \rangle & \langle xy \rangle \\ \langle xy \rangle & \langle y^2 \rangle \end{pmatrix}$$
- It has 2 eigenvalues. If y is perfectly correlated with x the matrix is singular and one eigenvalue is 0.
- In this case we then have a linear relation between x and y .
- Note that PCA is a linear method

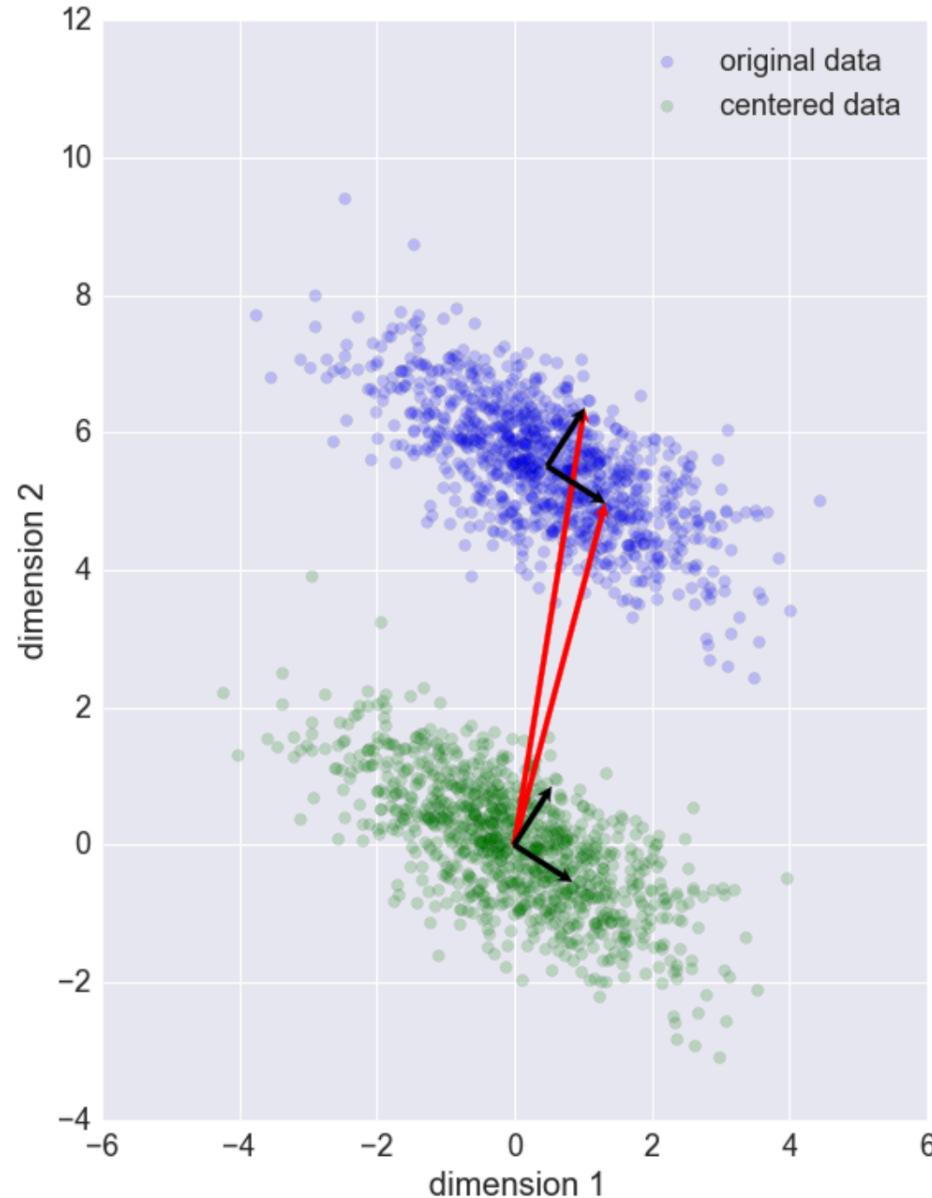


SVD Properties, Generalizations

- The method is non-parametric. Unique answers, but may not give the best answer when we have additional information
- PCA is linear. We may be able to make the data analysis more linear if we perform a non-linear transformation of variables: kernel PCA
- PCA is a dimension reduction method. One can instead require independence of components: independent component analysis (ICA)

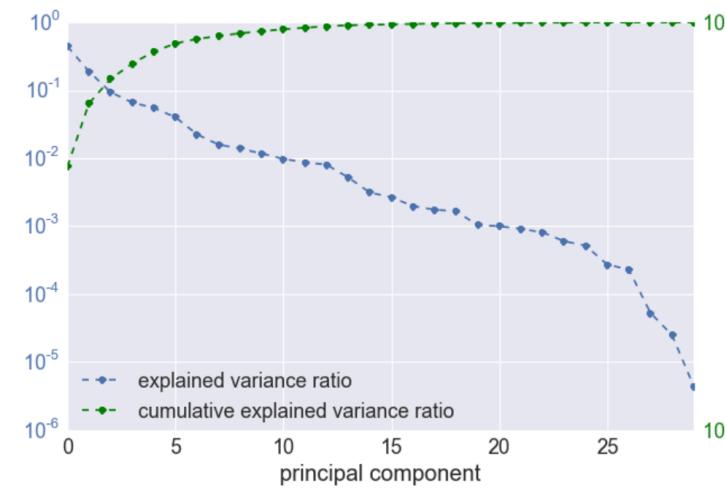
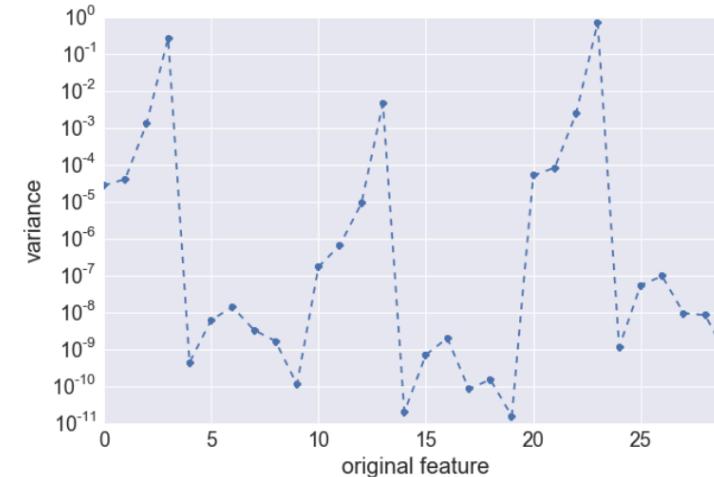
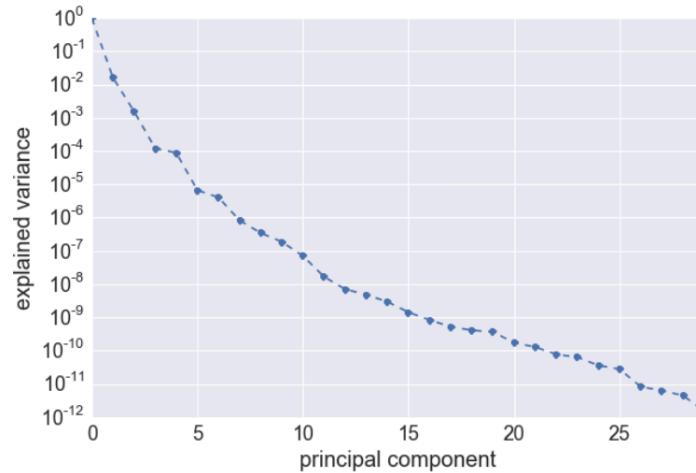
SVD Preprocessing

Centering: taking the mean out



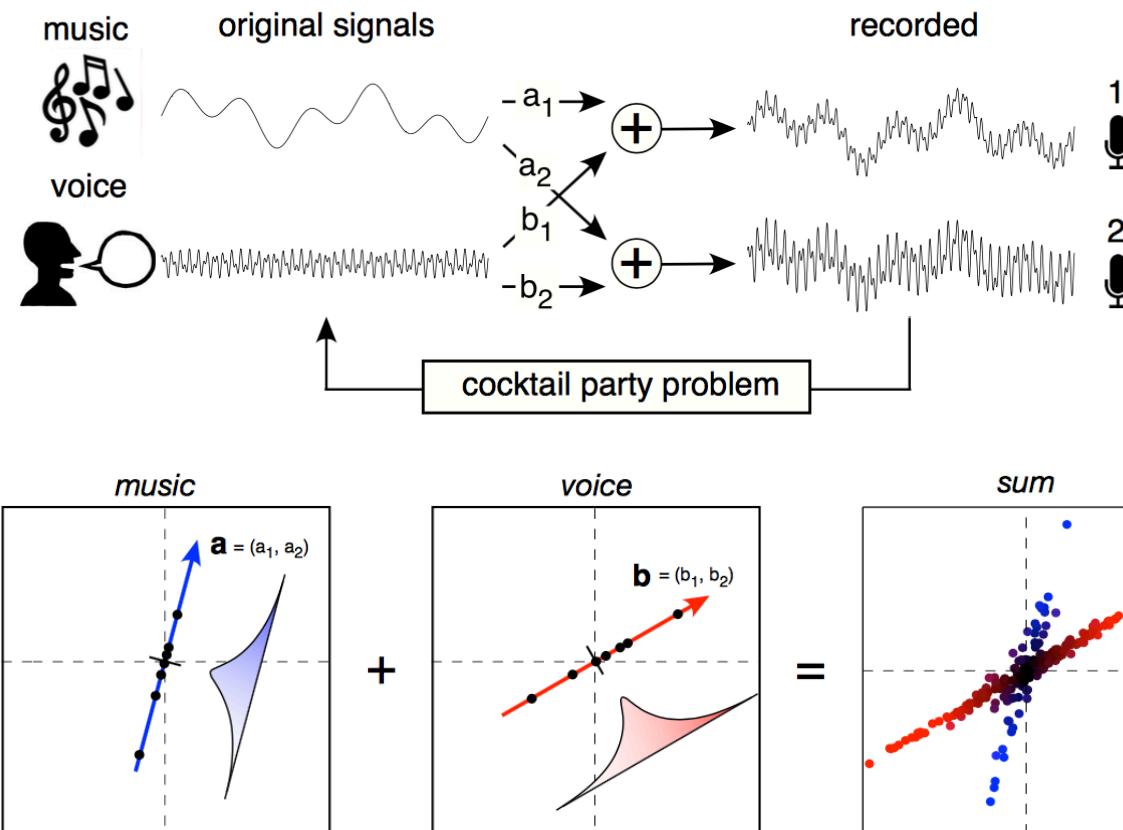
SVD Preprocessing

- Scaling or normalizing: use the units where the data are all similar in numbers, otherwise rescale

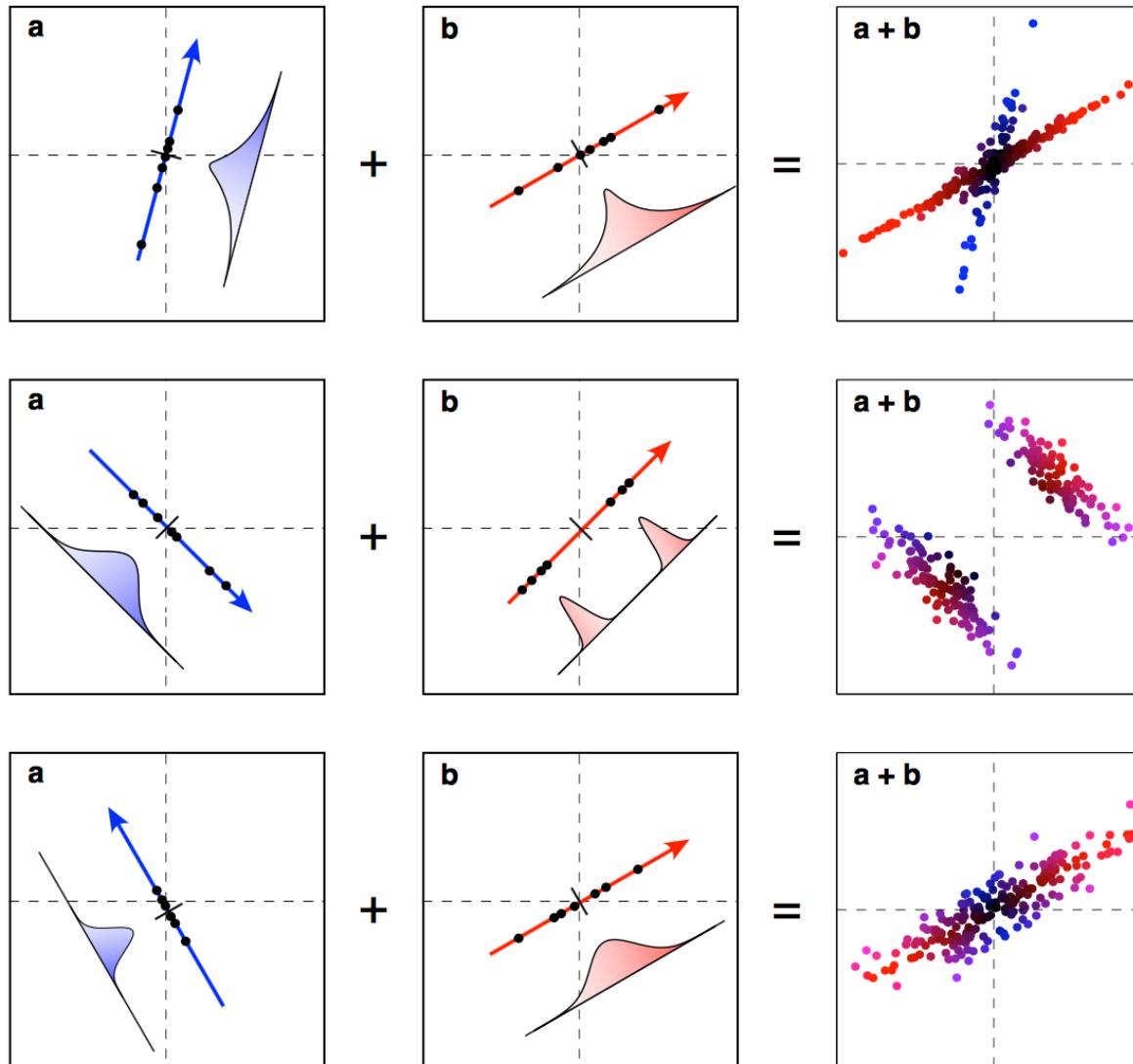


Independent Component Analysis (ICA): cocktail party problem

- We have 2 sources of sound and 2 microphones and we would like to separate the two

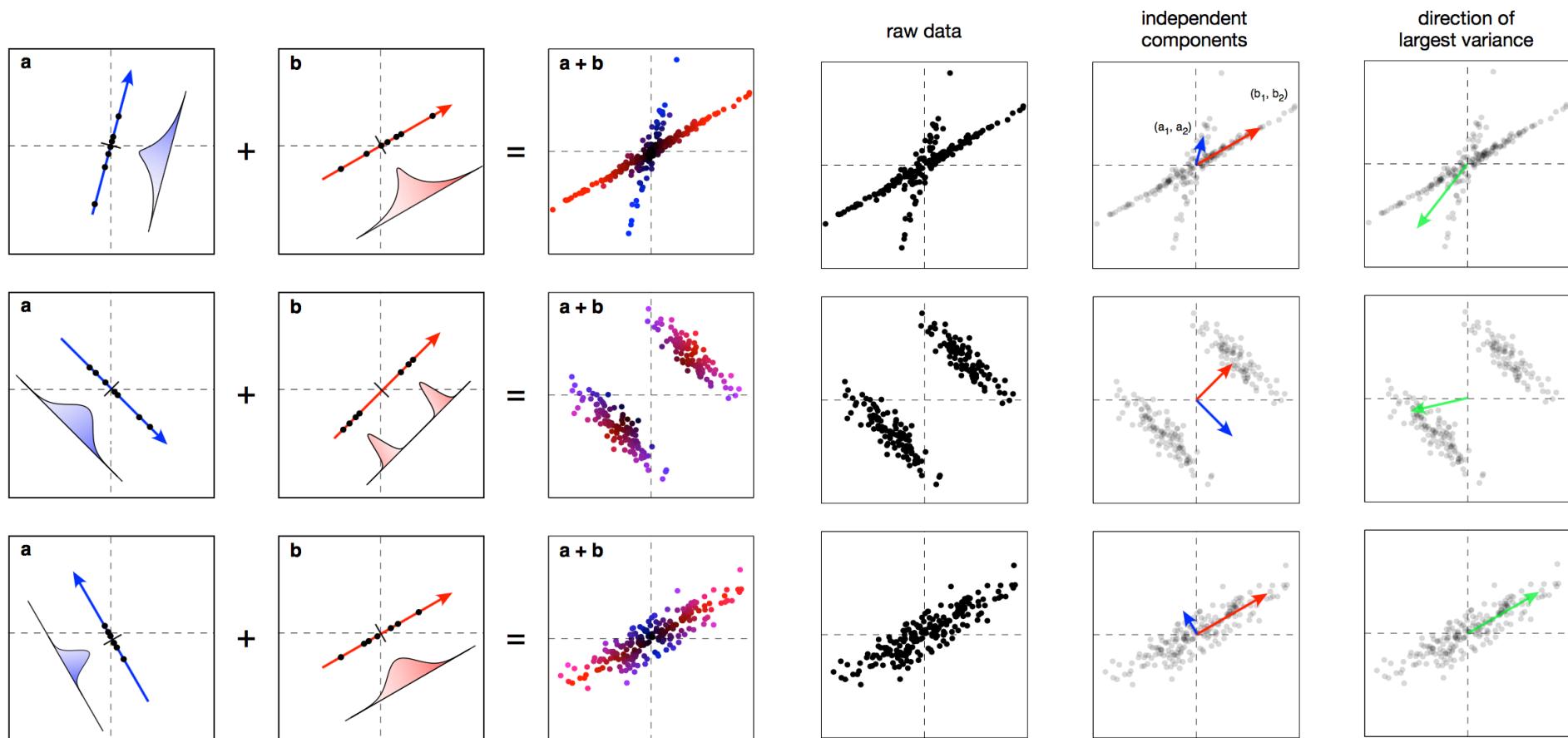


Many possible forms of independent sources



PCA vs. ICA

- They are different in general
- Mean is always subtracted
- PCA is not very meaningful for non-Gaussian distributions

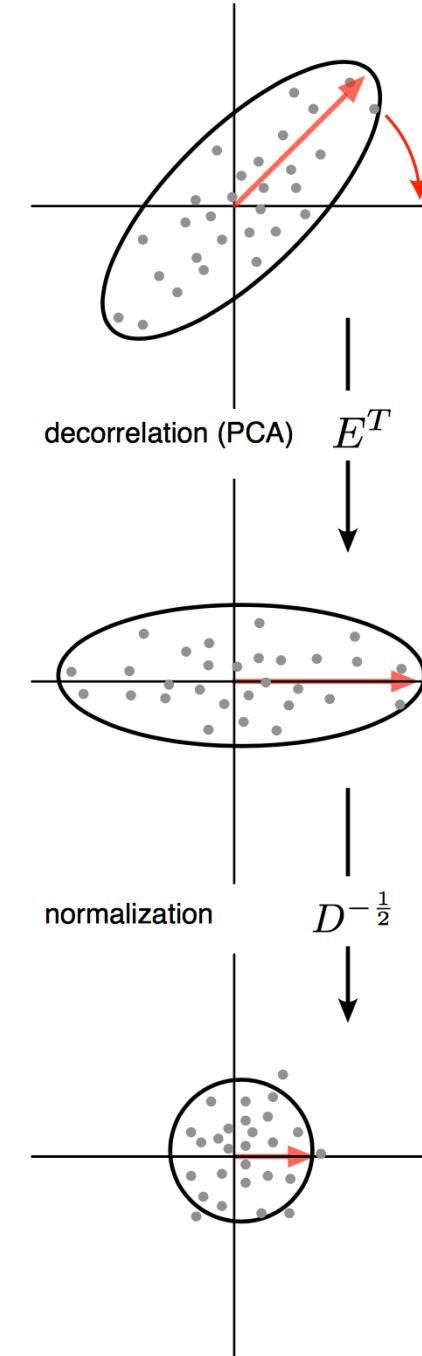


ICA Setup

- We have data \mathbf{x} and would like to reconstruct individual sources \mathbf{s} assuming they are statistically independent. We subtract the mean and assume linear relation $\mathbf{x} = \mathbf{As}$ but we do not know \mathbf{A} or \mathbf{s} . We assume linear reconstruction $\mathbf{s}' = \mathbf{Wx}$. Our task is to determine \mathbf{W} .
- We can assume SVD of $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$
- $\mathbf{W} = \mathbf{A}^{-1} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T$
- Let us whiten \mathbf{s} , $\mathbf{s}\mathbf{s}^T = \mathbf{I}$
- Then $\mathbf{x}\mathbf{x}^T = \mathbf{A}\mathbf{s}\mathbf{s}^T\mathbf{A}^T = \mathbf{U}\Sigma\mathbf{V}^T\mathbf{s}\mathbf{s}^T\mathbf{V}\Sigma^T\mathbf{U}^T = \mathbf{U}\Sigma\Sigma^T\mathbf{U}^T = \mathbf{U}\Sigma^2\mathbf{U}^T$
- We can diagonalize $\mathbf{x}\mathbf{x}^T = \mathbf{E}\mathbf{D}\mathbf{E}^T$, so we find SVD solution $\mathbf{W} = \mathbf{V}\mathbf{D}^{-1/2}\mathbf{E}^T$. We know \mathbf{E} and \mathbf{D} .

ICA Continued

- We have identified \mathbf{U} and Σ of \mathbf{A} , without knowing \mathbf{A} . We first decorrelated \mathbf{A} via PCA (\mathbf{E}) and then whitened \mathbf{A} via $\mathbf{D}^{-1/2}$:
 $\mathbf{x}_w = \mathbf{D}^{-1/2} \mathbf{E}^T$ and $\mathbf{x}_w \mathbf{x}_w^T = \mathbf{I}$
- We need to find $\mathbf{s}' = \mathbf{V} \mathbf{x}_w$. Since \mathbf{V} is orthonormal this can be viewed as another rotation that does not change $\mathbf{s}' \mathbf{s}'^T = \mathbf{I}$.
- If the problem was Gaussian uncorrelated and statistically independent are the same: any \mathbf{V} would be as good as any other: rotation of a circle is still a circle.
- If it is non-Gaussian we can impose some other statistic to rotate \mathbf{x}_w into statistically independent components

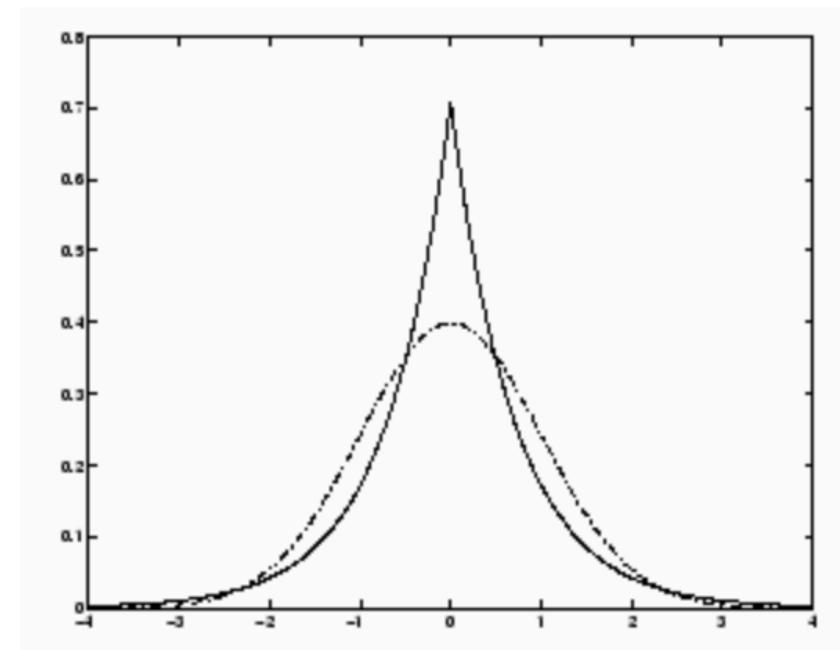


Measures of non-Gaussianity

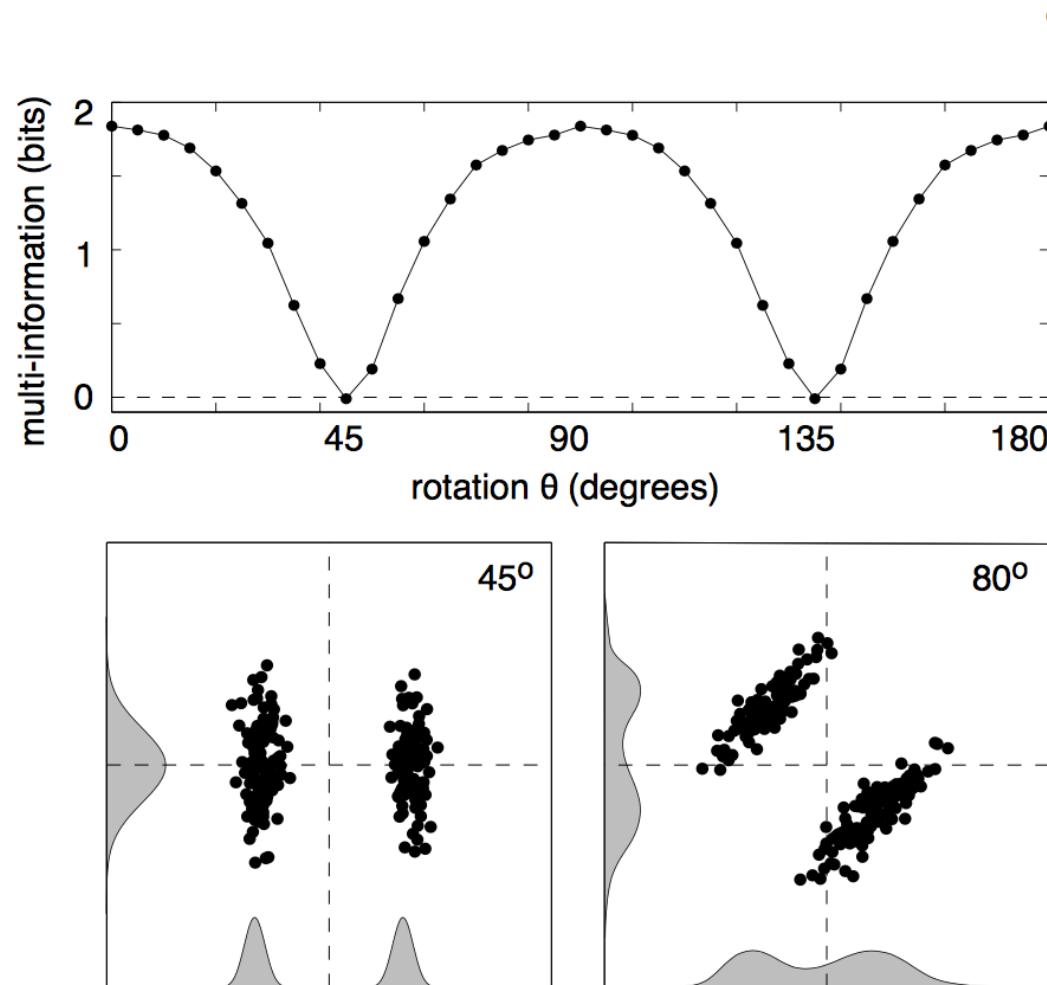
- Moments: skewness, kurtosis, etc.
- Since we want independence we want to set to 0 cross-terms, e.g.
- Kurtosis $\text{Kurt}_{12} = \langle \mathbf{x}_{w1}^2 \mathbf{x}_{w2}^2 \rangle - 1$
- In 2d \mathbf{V} can be represented as a rotation angle

$$\mathbf{v} = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

- Vary θ until Kurt_{12} is 0.



Doubly Peaked Example



This example uses multi-information: **information theory** (next lecture)

Main Decomposition of Linear Algebra

- 1. Singular Value Decomposition $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$
2. Gram-Schmidt Orthogonalization $\mathbf{A} = \mathbf{Q} \mathbf{R}$ (can be defined for $N \times M$)
3. Gaussian Elimination $\mathbf{A} = \mathbf{L} \mathbf{U}$
- We might prefer to separate out singular values σ_i and heights h_i and pivots d_i :
 1. $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$ with unit vectors in \mathbf{U} and \mathbf{V} .
The singular values σ_i are in Σ .
 2. $\mathbf{A} = \mathbf{Q} \mathbf{H} \mathbf{R}$ with unit vectors in \mathbf{Q} and diagonal 1's in \mathbf{R} .
The heights h_i are in \mathbf{H} .
 3. $\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{U}$ with diagonal 1's in \mathbf{L} and \mathbf{U} .
The pivots d_i are in \mathbf{D} .
- Each h_i tells the height of column i above the plane of columns 0 to $i - 1$.
The volume of the full N -dimensional box ($r=M=N$) comes from
 $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T = \mathbf{L} \mathbf{D} \mathbf{U} = \mathbf{Q} \mathbf{H} \mathbf{R}$:
- $|\det \mathbf{A}| = |\text{product of } \sigma\text{'s}| = |\text{product of } d\text{'s}| = |\text{product of } h\text{'s}|$.

Literature

- *Numerical Recipes*, Press et al., Chapter 2, 11, 15
- *Computational Physics*, M. Newman, Chapter 6.2