

Dimensionality Reduction and Principle Component Analysis

Felix Brockherde 2015, TU Berlin and MPI Halle,
Stefan Haufe 2016, TU Berlin



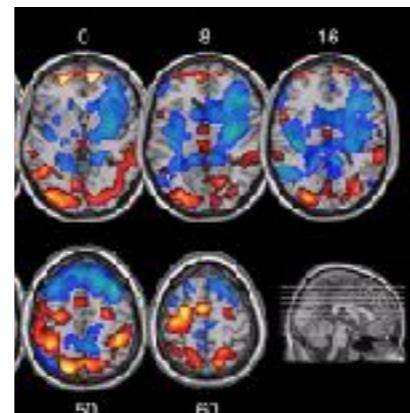
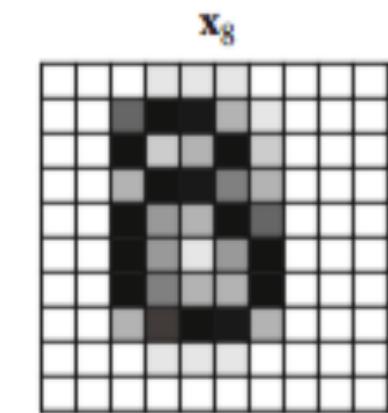
MAX-PLANCK-GESELLSCHAFT

This Lecture

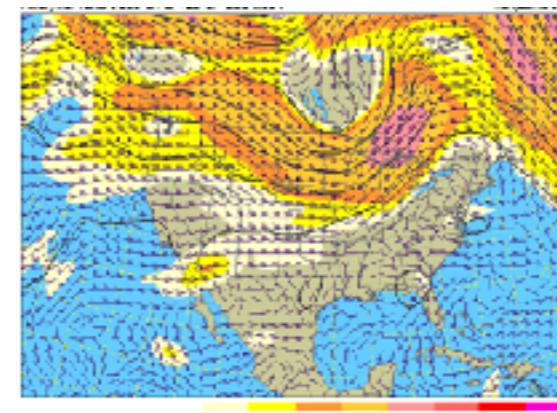
1. Dimensionality reduction
2. Principle Component Analysis
 1. What are Principle Components?
 2. How to find/calculate them
 1. Lagrange Multipliers
 3. What can we do with them? / Applications

Curse of dimensionality

In many machine learning problems, the data are high-dimensional



From Limb and Braun, 2008



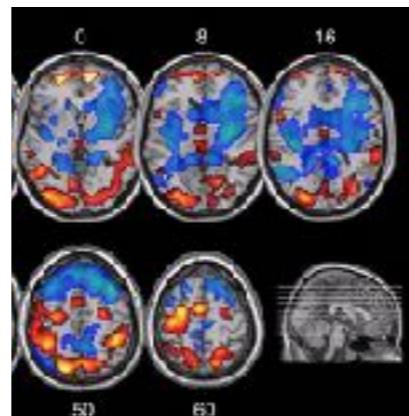
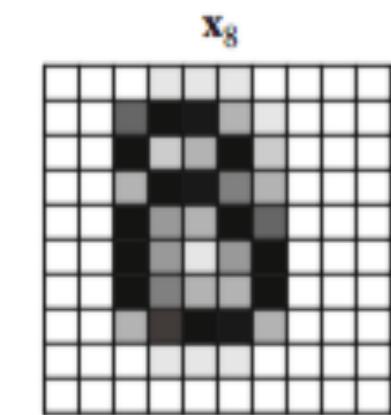
University of Colorado



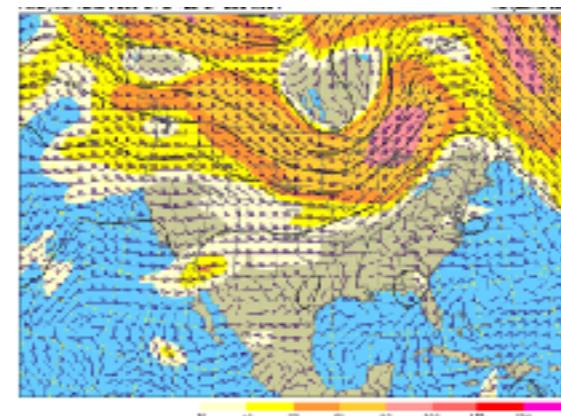
IEEE Spectrum

Curse of dimensionality

In many machine learning problems, the data are high-dimensional



From Limb and Braun, 2008



University of Colorado



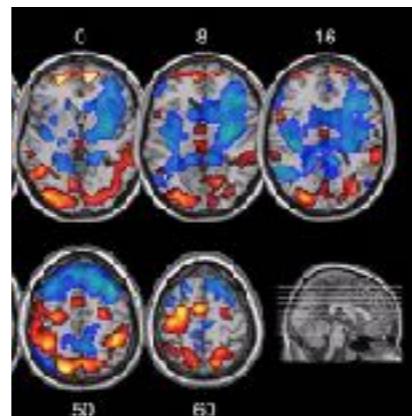
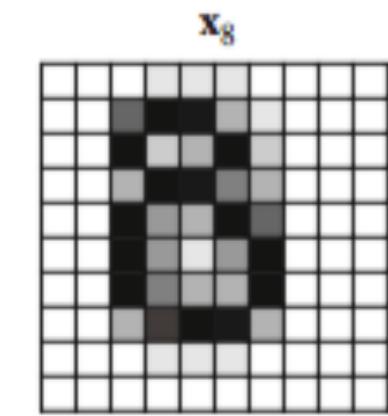
IEEE Spectrum

Standard regression/classification techniques can become

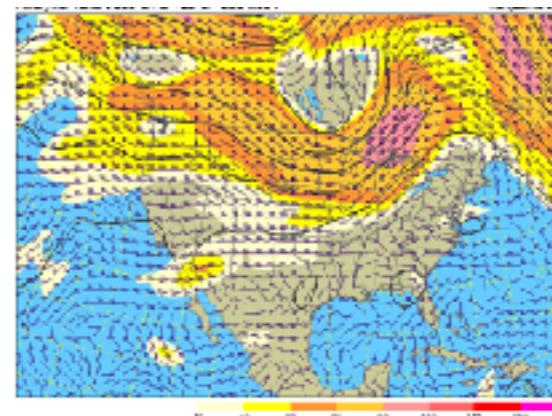
- ill-defined for $M \gg N$
 - ill-conditioned/numerically unstable even for $M < N$

Curse of dimensionality

In many machine learning problems, the data are high-dimensional



From Limb and Braun, 2008



University of Colorado



IEEE Spectrum

Standard regression/classification techniques can become

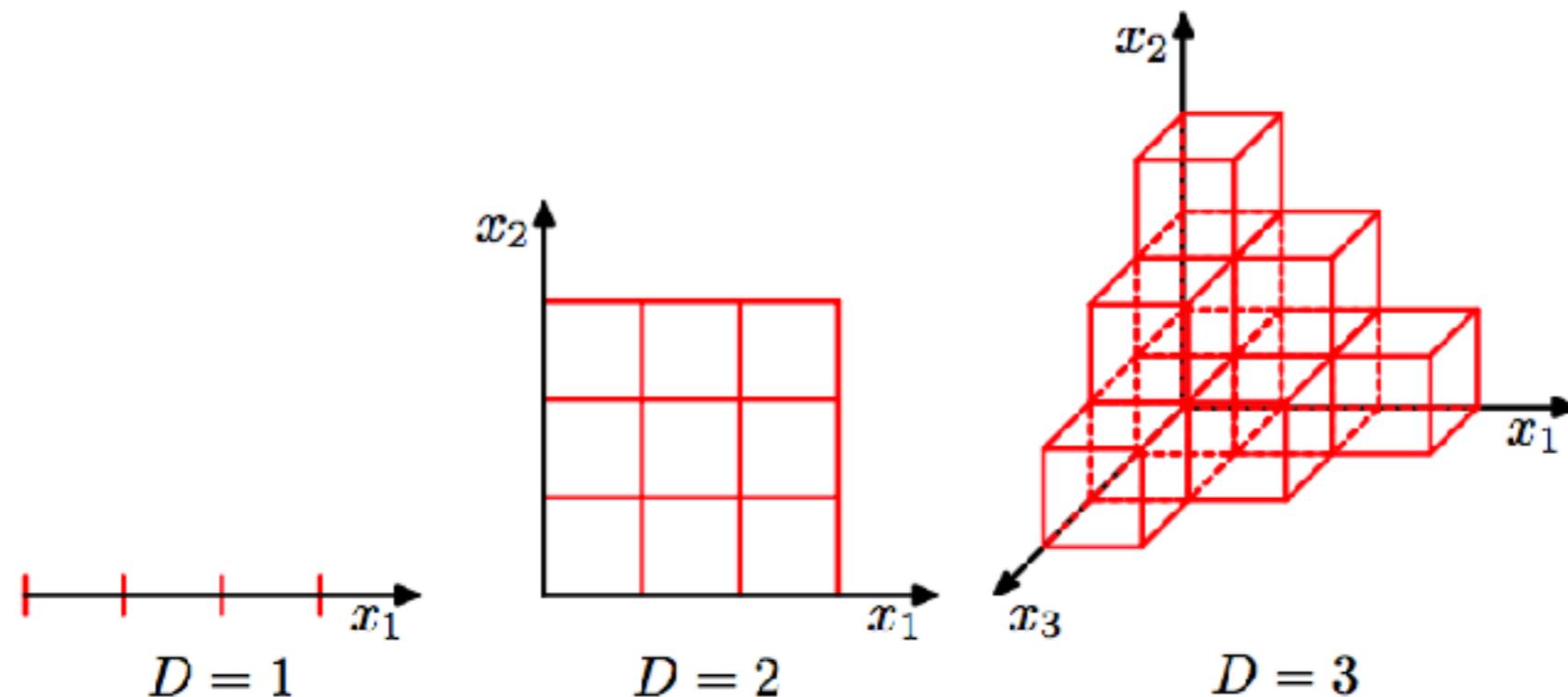
- ill-defined for $M \gg N$
- ill-conditioned/numerically unstable even for $M < N$

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Singular,
ill-conditioned

Curse of dimensionality

When the dimensionality increases, the volume of the space increases so fast that the available data becomes sparse



The amount of data needed for a reliable result often grows exponentially with the dimensionality

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Linear regression setting: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Linear regression setting: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Prior: $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 I)$

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Linear regression setting: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Prior: $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 I)$

= favor a “simple” model

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Linear regression setting: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Prior: $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 I)$

= favor a “simple” model

Maximum a-posteriori (MAP) approach: $\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} p(\boldsymbol{\beta} | \mathcal{D})$

$$= \arg \max_{\boldsymbol{\beta}} p(\mathcal{D} | \boldsymbol{\beta}) p(\boldsymbol{\beta})$$

$$= \arg \max_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \frac{\sigma_{\boldsymbol{\epsilon}}^2}{\sigma_{\boldsymbol{\beta}}^2} \|\boldsymbol{\beta}\|^2$$

$$= (\mathbf{X}^\top \mathbf{X} + \frac{\sigma_{\boldsymbol{\epsilon}}^2}{\sigma_{\boldsymbol{\beta}}^2} I)^{-1} \mathbf{X}^\top \mathbf{y}$$

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Linear regression setting: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Prior: $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 I)$

= favor a “simple” model

Maximum a-posteriori (MAP) approach: $\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} p(\boldsymbol{\beta} | \mathcal{D})$

$$= \arg \max_{\boldsymbol{\beta}} p(\mathcal{D} | \boldsymbol{\beta}) p(\boldsymbol{\beta})$$

$$= \arg \max_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \frac{\sigma_{\boldsymbol{\epsilon}}^2}{\sigma_{\boldsymbol{\beta}}^2} \|\boldsymbol{\beta}\|^2$$

$$= (\mathbf{X}^\top \mathbf{X} + \frac{\sigma_{\boldsymbol{\epsilon}}^2}{\sigma_{\boldsymbol{\beta}}^2} I)^{-1} \mathbf{X}^\top \mathbf{y}$$

Regular,
better-conditioned

Regularization

Idea: impose constraints on the parameters to stabilize solution.

One way: introduce prior probability

Linear regression setting: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Prior: $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 I)$

= favor a “simple” model

Maximum a-posteriori (MAP) approach: $\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} p(\boldsymbol{\beta} | \mathcal{D})$

$$= \arg \max_{\boldsymbol{\beta}} p(\mathcal{D} | \boldsymbol{\beta}) p(\boldsymbol{\beta})$$

$$= \arg \max_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \frac{\sigma_{\boldsymbol{\epsilon}}^2}{\sigma_{\boldsymbol{\beta}}^2} \|\boldsymbol{\beta}\|^2$$

$$= (\mathbf{X}^\top \mathbf{X} + \frac{\sigma_{\boldsymbol{\epsilon}}^2}{\sigma_{\boldsymbol{\beta}}^2} I)^{-1} \mathbf{X}^\top \mathbf{y}$$

Regular,
better-conditioned

= Ridge Regression

Dimensionality reduction

Problem: we still have to invert an M-dimensional matrix → expensive

Goal: reduce data to features most relevant for the learning task

Dimensionality reduction

Problem: we still have to invert an M-dimensional matrix → expensive

Goal: reduce data to features most relevant for the learning task

One way: significance test for single features

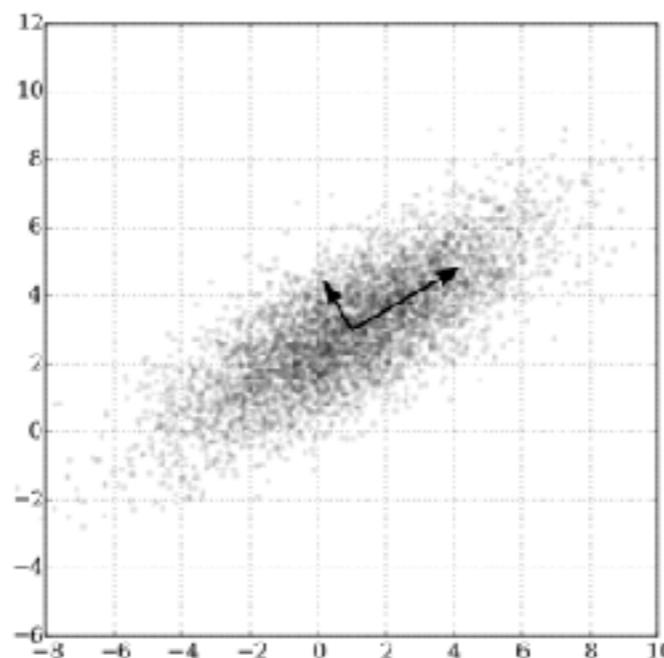
Dimensionality reduction

Problem: we still have to invert an M-dimensional matrix → expensive

Goal: reduce data to features most relevant for the learning task

One way: significance test for single features

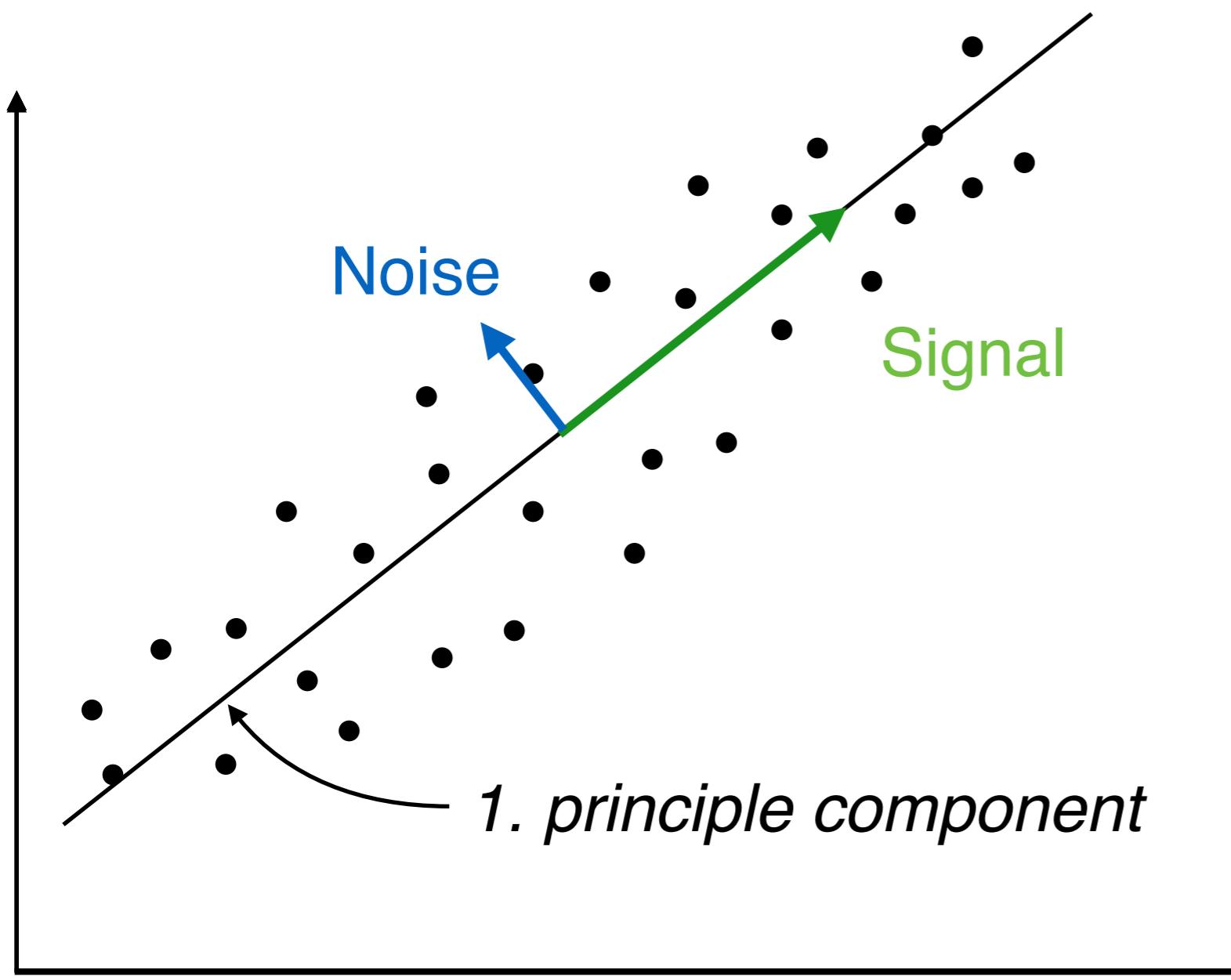
Another way: find ‘relevant’ directions/subspaces in correlated data



Why dimensionality reduction?

- **Visualization:**
Insights into high-dimensional structures in the data
- **Better Generalization:**
Fewer dimensions → less chances of overfitting
- **Speeding up** learning algorithms:
Most algorithms scale badly with increasing data dimensionality
- **Data compression:**
Less storage requirements

Principle Component Analysis (PCA)



Which line fits data best?

The line w that minimizes the **noise** and maximizes the **signal** [Pearson 1901].

Problem Definition

Given $x_1, \dots, x_n \in \mathbb{R}^d$ find a k -dimensional subspace, so that the data projected on that subspace

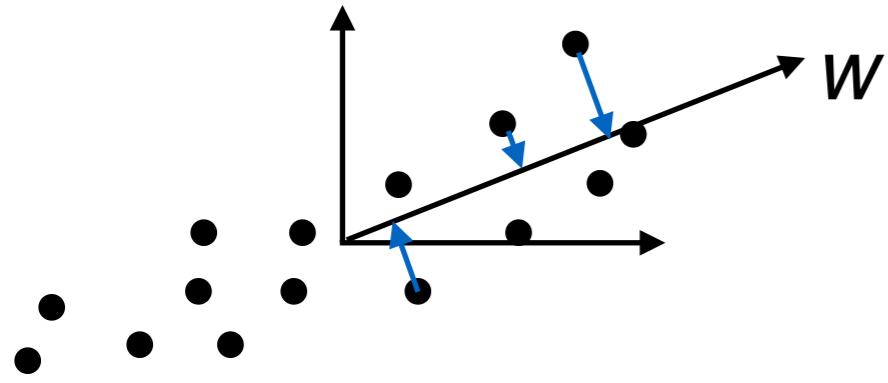
- 1) is as close to the original data as possible (minimum noise)
- 2) has maximum variance (maximum signal).



2 Motivations:
same result?

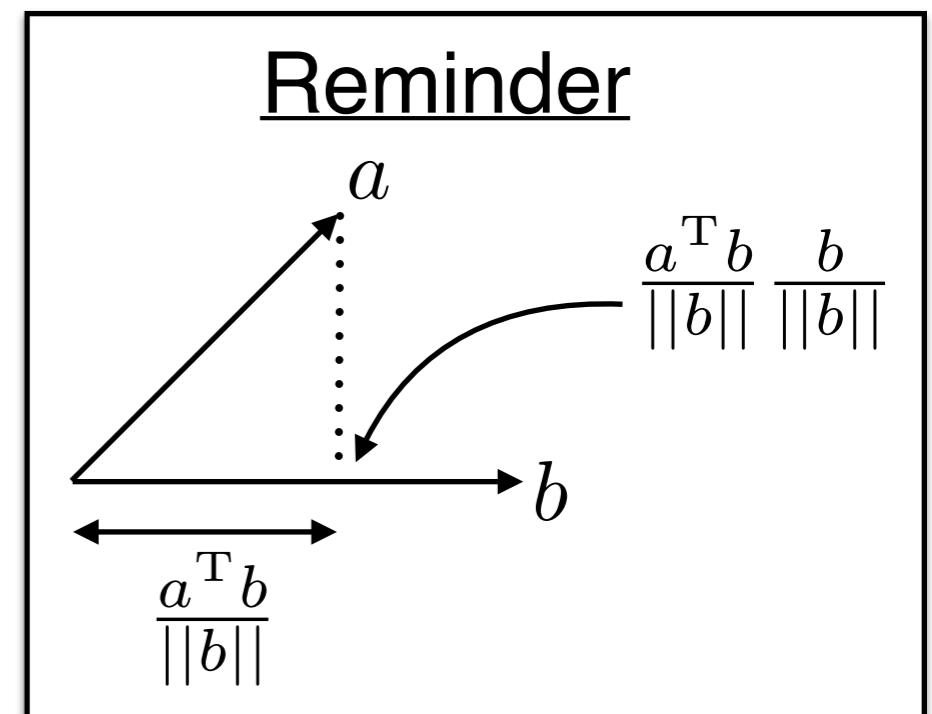
Assume **the data is centered**: $\mu = \frac{1}{n} \sum_{i=1}^n x_i = 0$

Minimum distance

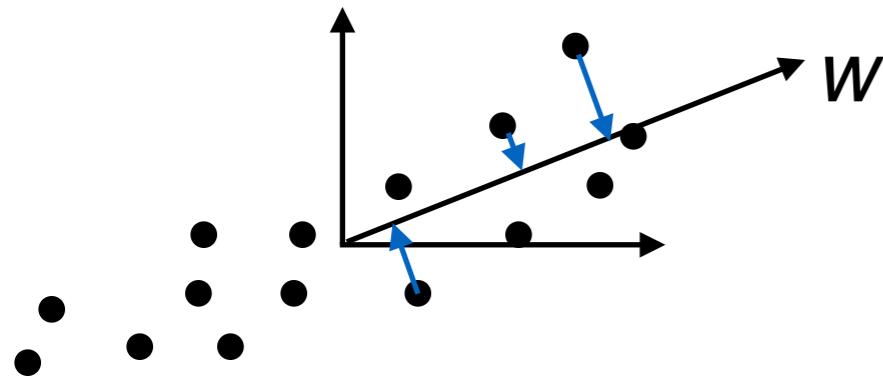


We are only interested in the direction of w , not its length, i.e. set $\|w\| = w^T w = 1$.

$$\begin{aligned} & \min_w \sum_{i=1}^n \|(w^T x_i)w - x_i\|^2 \\ &= \min_w \sum_{i=1}^n w^T x_i x_i^T w w^T w - 2w^T x_i x_i^T w + x_i^T x_i \\ &= \min_w - \sum_{i=1}^n w^T x_i x_i^T w \\ &= \max_w \sum_{i=1}^n w^T x_i x_i^T w \\ &= \max_w w^T X X^T w \end{aligned}$$



Maximum Variance



$$\max_w \text{Var}(w^T X) = E[(w^T X)^2]$$

$$= \max_w \frac{1}{n} \sum_{i=1}^n w^T x_i x_i^T w$$

$$= \max_w w^T X X^T w$$

✓ (same as
minimum
distance)

Introduce the *scatter matrix* $S = X X^T$. This leads to the constrained optimization problem:

$$\max_w w^T S w$$

$$\text{s.t. } \|w\| = 1$$

Lagrange Multipliers

- Imagine you have a constrained optimization problem:

Read: subject to

$$\begin{aligned} & \max_x f(x) && \text{Objective function} \\ & \text{s.t. } g(x) = 0 && \text{Constraint} \end{aligned}$$

- Example:

$$\max_{\mathbf{x}} 1 - x_1^2 - x_2^2$$

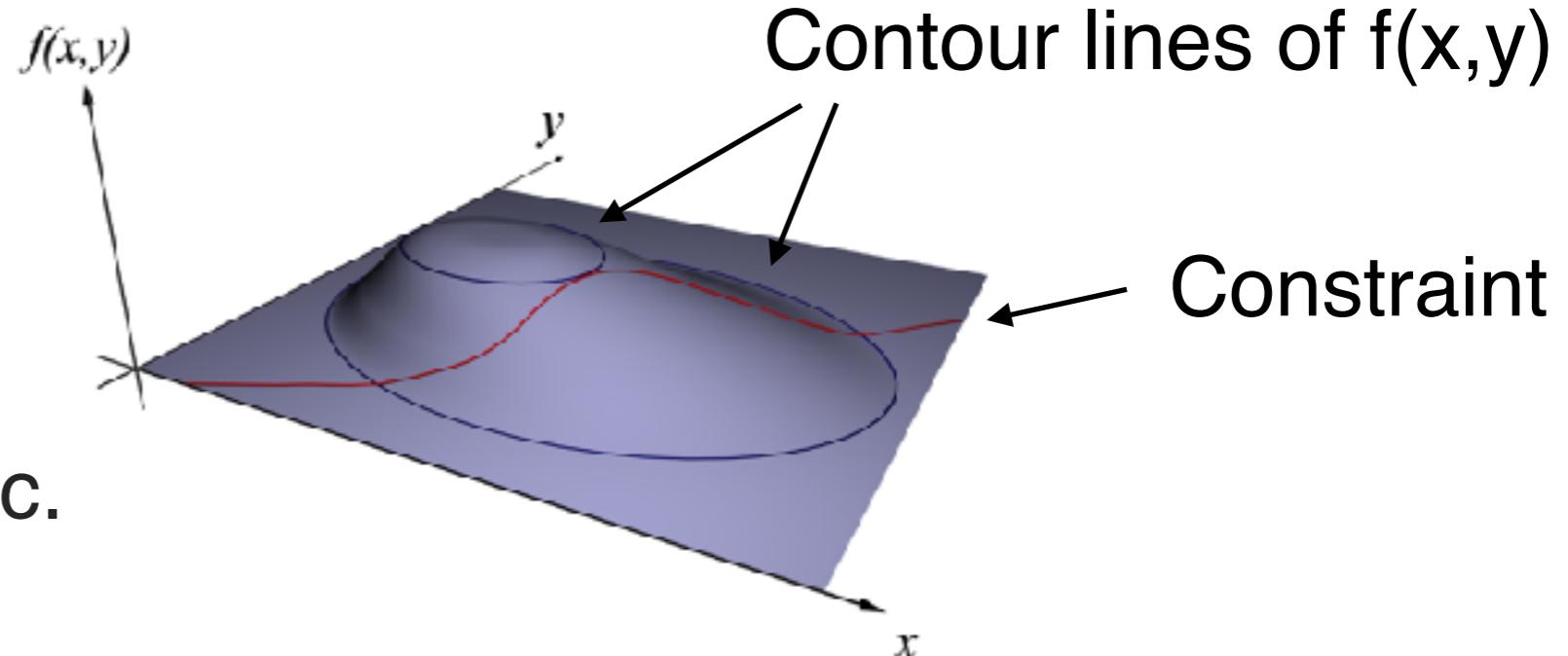
$$\text{s.t. } x_1 + x_2 - 1 = 0$$

Geometric intuition

$$\max_x f(x)$$

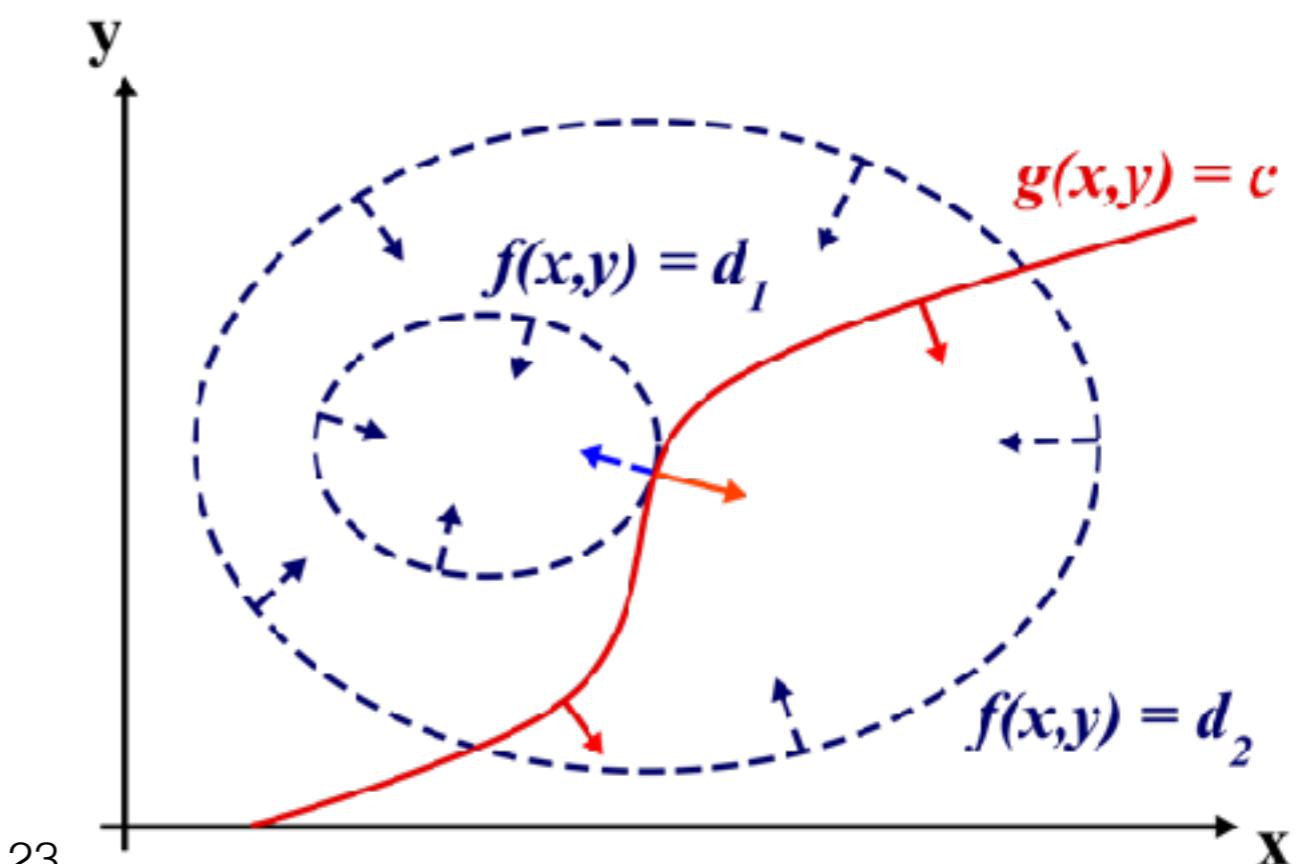
$$\text{s.t. } g(x) = 0$$

Intuition 1: ∇g must be normal to the line $g(x, y) = c$.



Intuition 2: At an optimum $f(x^*,y^*)$ can not increase in the direction of a neighboring point where $g(x, y) = c$. Thus, at (x^*, y^*) ∇f must be perpendicular to the line $g(x, y) = c$.

It follows: In (x^*, y^*) the gradients ∇g and ∇f are parallel!

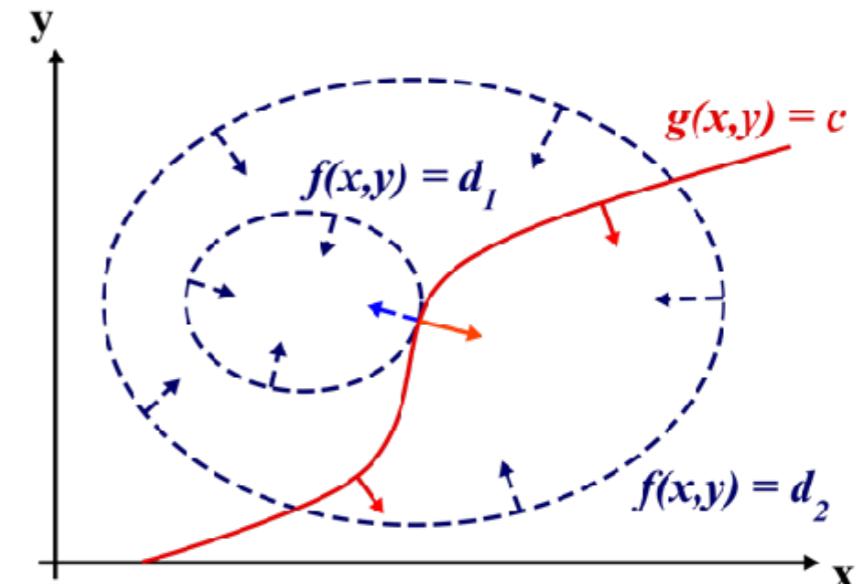


The Lagrange Multiplier

In (x^*, y^*) the gradients ∇g and ∇f are parallel. Thus for a maximum:

$$\nabla f = \lambda \nabla g$$

Lagrange Multiplier



Lagrangian function:

$$\mathcal{L}(x, \lambda) = f(x) - \lambda g(x)$$

$\nabla L = 0$ is a necessary (but not sufficient) condition for the optimization solution.

Thus, to solve a constrained optimization problem, we can define the Lagrangian and look for the points where its gradient vanishes.

Example

$$\max_{\mathbf{x}} 1 - x_1^2 - x_2^2$$

$$\text{s.t. } x_1 + x_2 - 1 = 0$$

1. Define Lagrangian:

$$\mathcal{L}(x, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1)$$

2. Set gradient of Lagrangian to zero:

$$-2x_1 + \lambda = 0$$

$$-2x_2 + \lambda = 0$$

$$x_1 + x_2 - 1 = 0$$

3. Solve equations: $(x_1^*, x_2^*) = \left(\frac{1}{2}, \frac{1}{2}\right)$ $\lambda = 1$

PCA optimization problem

$$\begin{aligned} \max_w \quad & w^T S w \\ \text{s.t.} \quad & \|w\| = 1 \end{aligned}$$

1. Define Lagrangian: $\mathcal{L}(w, \lambda) = w^T S w + \lambda(1 - w^T w)$

2. Compute gradient: $\frac{\partial \mathcal{L}(w, \lambda)}{\partial w} = 2S w - 2\lambda w$

3. Set to zero: $S w = \lambda w$

PCA optimization problem

$$\max_w w^T S w$$

$$\text{s.t. } \|w\| = 1$$

1. Define Lagrangian: $\mathcal{L}(w, \lambda) = w^T S w + \lambda(1 - w^T w)$

2. Compute gradient: $\frac{\partial \mathcal{L}(w, \lambda)}{\partial w} = 2S w - 2\lambda w$

3. Set to zero: $S w = \lambda w$



This is an
Eigenvalue
problem!

Eigenvalue problems

$$Sw = \lambda w$$

$$(S - \lambda I)w = 0$$

Eigenvalue problems

$$Sw = \lambda w$$

$$(S - \lambda I)w = 0$$

Solutions are found at roots of characteristic polynomial

$$p(\lambda) := \det(S - \lambda I) = 0$$

Eigenvalue problems

$$Sw = \lambda w$$

$$(S - \lambda I)w = 0$$

Solutions are found at roots of characteristic polynomial

$$p(\lambda) := \det(S - \lambda I) = 0$$

In general: d roots (eigenvalues) $\lambda_1, \dots, \lambda_d$

Corresponding eigenvectors: w_1, \dots, w_d

Eigenvalue problems

$$Sw = \lambda w$$

$$(S - \lambda I)w = 0$$

Solutions are found at roots of characteristic polynomial

$$p(\lambda) := \det(S - \lambda I) = 0$$

In general: d roots (eigenvalues) $\lambda_1, \dots, \lambda_d$

Corresponding eigenvectors: w_1, \dots, w_d

Leads to the decomposition $W^\top SW = \Lambda$ $S = W\Lambda W^\top$

where $W = [w_1, \dots, w_d]$, $W^\top W = I$ is orthogonal and
 $\Lambda, \Lambda_{ii} = \lambda_i$ is diagonal.

PCA algorithm

The i-th eigenvalue is the variance in the direction of the i-th eigenvector:

$$\text{Var}(w_i^T x) = \frac{1}{n} w_i^T S w_i = \frac{1}{n} \lambda_i w_i^T w_i = \frac{1}{n} \lambda_i$$

PCA algorithm

The i-th eigenvalue is the variance in the direction of the i-th eigenvector:

$$\text{Var}(w_i^T x) = \frac{1}{n} w_i^T S w_i = \frac{1}{n} \lambda_i w_i^T w_i = \frac{1}{n} \lambda_i$$

The direction of largest variance corresponds to the largest eigenvector (= the eigenvector with largest eigenvalue).

PCA algorithm

The i -th eigenvalue is the variance in the direction of the i -th eigenvector:

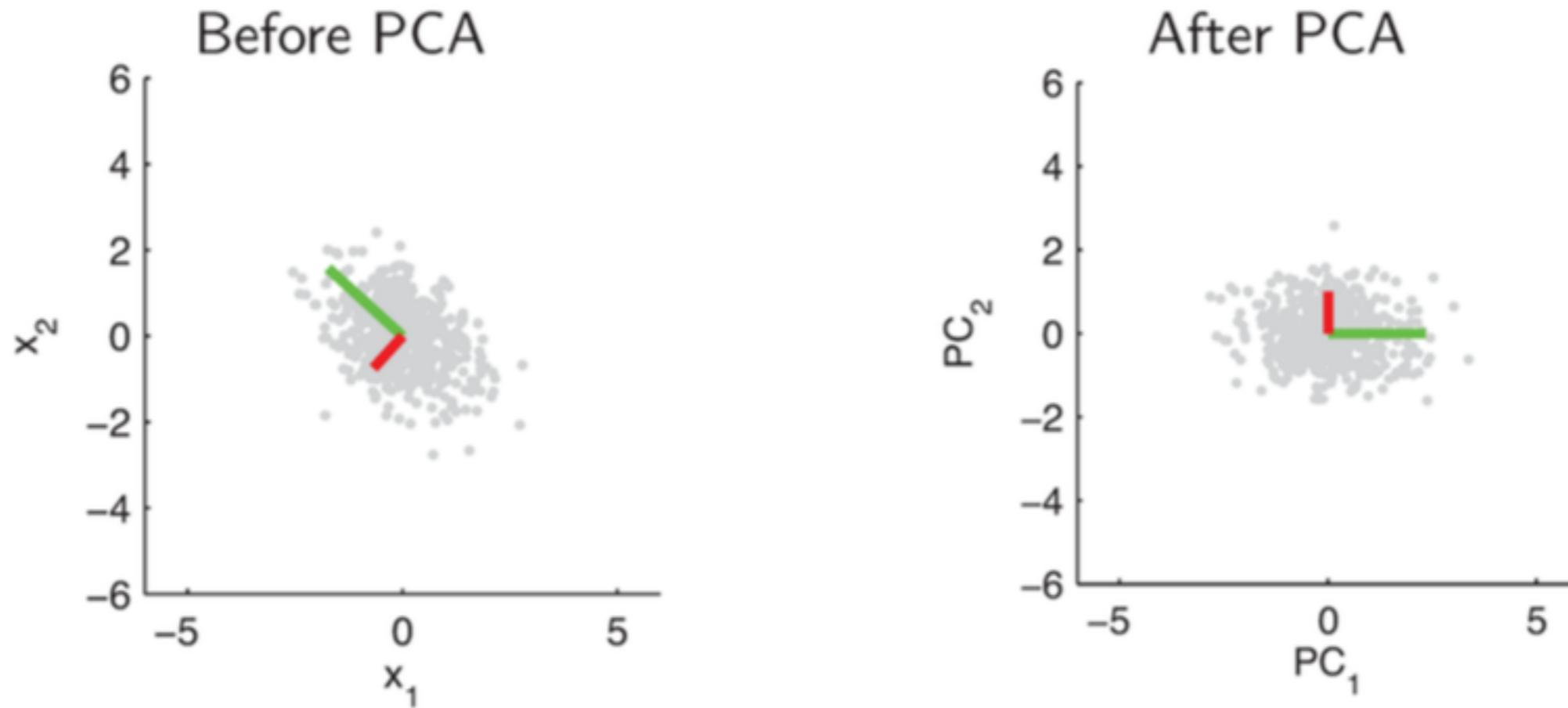
$$\text{Var}(w_i^T x) = \frac{1}{n} w_i^T S w_i = \frac{1}{n} \lambda_i w_i^T w_i = \frac{1}{n} \lambda_i$$

The direction of largest variance corresponds to the largest eigenvector (= the eigenvector with largest eigenvalue).

Algorithm 1: Principal Component Analysis

Require: data $x_1, \dots, x_N \in \mathbb{R}^d$, number of principal components k

- 1: # Center Data
 - 2: $X = X - 1/N \sum_i x_i$
 - 3: # Compute Covariance Matrix
 - 4: $C = 1/N X X^T$
 - 5: # Compute largest k eigenvectors
 - 6: $W = \text{eig}(C)$
 - 7: **return** W
-



PCA rotates data into new coordinate system with the directions of largest variances being the new coordinate axes.

Questions

1. What happens when the data are not centered?
2. What happens if data are already standardized (features have zero mean and unit variance)?

Solving PCA with SVD

A singular value decomposition factorizes a matrix as:

$$U\Lambda V = M$$

where

U are the Eigenvectors of MM^* .

V are the Eigenvectors of M^*M .

The square roots of the Eigenvalues of MM^* are on the diagonal of Λ .

Solving PCA with SVD

A singular value decomposition factorizes a matrix as:

$$U\Lambda V = M$$

where

U are the Eigenvectors of MM^* .

V are the Eigenvectors of M^*M .

The square roots of the Eigenvalues of MM^* are on the diagonal of Λ .

Instead of calculating the EV-decomposition of S , we can compute the SVD-decomposition of X . This is computationally much more stable.

E.g. Läuchli Matrix

$$\begin{pmatrix} 1 & 1 & 1 \\ \epsilon & 0 & 0 \\ 0 & \epsilon & 0 \\ 0 & 0 & \epsilon \end{pmatrix}^\top$$

SVD in numpy

numpy.linalg.svd



`numpy.linalg.svd(a, full_matrices=1, compute_uv=1)`

[\[source\]](#)

Singular Value Decomposition.

Factors the matrix `a` as `u * np.diag(s) * v`, where `u` and `v` are unitary and `s` is a 1-d array of `a`'s singular values.

Parameters : `a : (..., M, N) array_like`

A real or complex matrix of shape `(M, N)`.

`full_matrices : bool, optional`

If True (default), `u` and `v` have the shapes `(M, M)` and `(N, N)`, respectively. Otherwise, the shapes are `(M, K)` and `(K, N)`, respectively, where `K = min(M, N)`.

`compute_uv : bool, optional`

Whether or not to compute `u` and `v` in addition to `s`. True by default.

Returns :

`u : {(..., M, M), (... , M, K)} array`

Unitary matrices. The actual shape depends on the value of `full_matrices`. Only returned when `compute_uv` is True.

`s : (... , K) array`

The singular values for every matrix, sorted in descending order.

`v : {(..., N, N), (... , K, N)} array`

Unitary matrices. The actual shape depends on the value of `full_matrices`. Only returned when `compute_uv` is True.

Raises : `LinAlgError`

If SVD computation does not converge.

Power iteration

The SVD has computational complexity $O(\min(n^2d, d^2n))$.

But: We often only need a few largest principle components and not all of them.

Algorithm 6.1 Simple vector iteration or power iteration

- 1: Choose a starting vector $\mathbf{x}^{(0)} \in \mathbb{F}^n$ with $\|\mathbf{x}^{(0)}\| = 1$.
 - 2: $k = 0$.
 - 3: **repeat**
 - 4: $k := k + 1$;
 - 5: $\mathbf{y}^{(k)} := A\mathbf{x}^{(k-1)}$;
 - 6: $\mu_k := \|\mathbf{y}^{(k)}\|$;
 - 7: $\mathbf{x}^{(k)} := \mathbf{y}^{(k)} / \mu_k$;
 - 8: **until** a convergence criterion is satisfied
-

The solution is the first eigenvector of A .

Power iteration intuition

Assume a diagonalizable matrix $A = U\Lambda U^T$.

The power iteration computes after step k: $A^k x$
where x is a random start vector.

We have $A^k = U\Lambda^k U^T$ because $U^T U = I$.

Transform x into a new coordinate system: $\tilde{x} = U^T x$

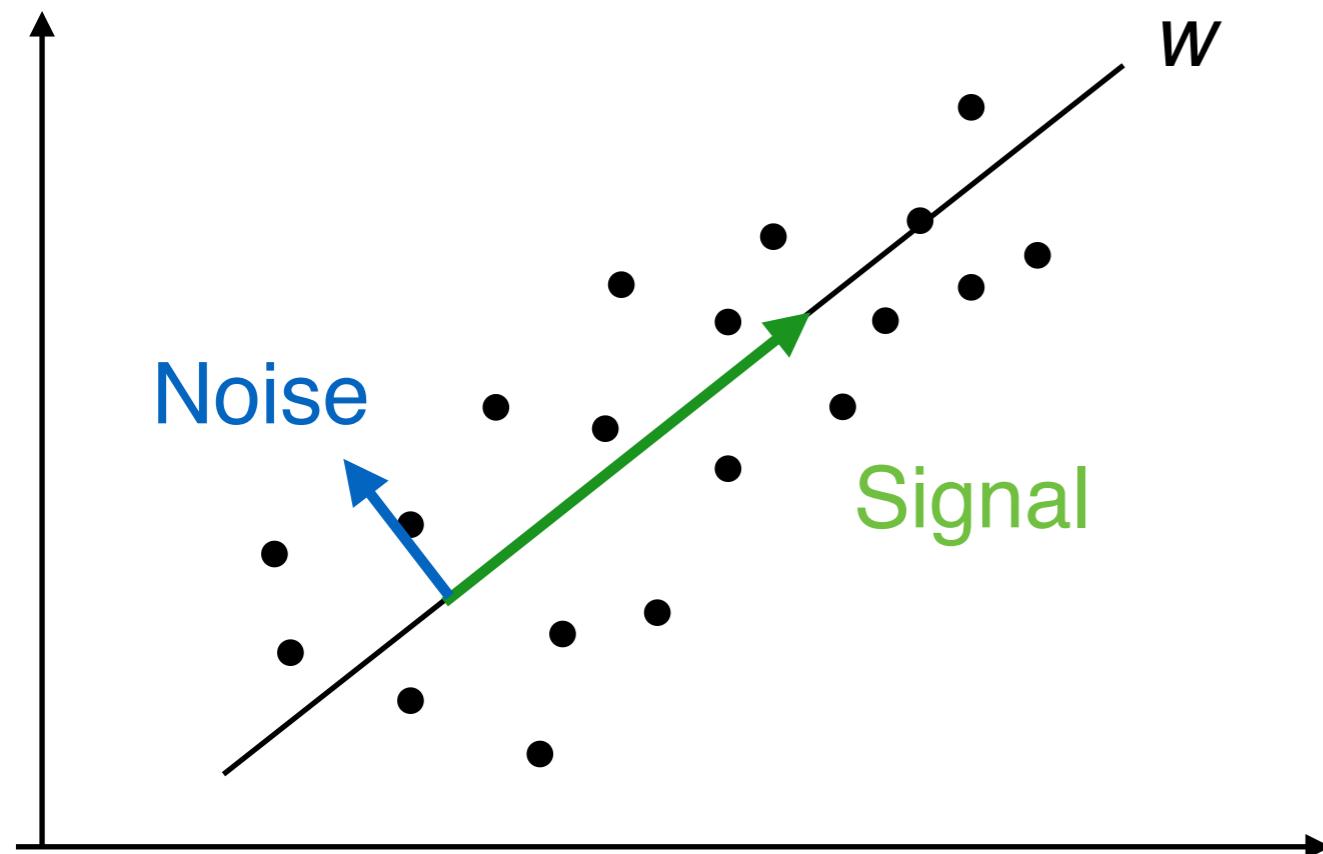
$$\begin{aligned} A^k x &= U\Lambda^k U^T x = U\Lambda^k \tilde{x} = \sum_{i=1}^d (\lambda_i^k \tilde{x}_i) u_i = \lambda_1^k \sum_{i=1}^d \left(\frac{\lambda_i^k}{\lambda_1^k} \tilde{x}_i\right) u_i \\ \Rightarrow \frac{A^k x}{\|A^k x\|} &\xrightarrow{k \rightarrow \infty} u_1 \end{aligned}$$

because $\frac{\lambda_i}{\lambda_1} \leq 1$ for $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$

Deflation

To find the second EV, remember $S = W\Lambda W^T = \sum_{i=1}^d \lambda_i w_i w_i^T$.
Thus, do power iteration on $\hat{S} = S - \lambda_1 w_1 w_1^T$.

Application: Dimensionality Reduction



Projection on k principle components:

$$\begin{pmatrix} w_1^T x \\ \vdots \\ w_k^T x \end{pmatrix} = \begin{bmatrix} | & & | \\ w_1 & \cdots & w_k \\ | & & | \end{bmatrix}^T x$$

We can reduce the dimensionality of our dataset by projecting on the first k principle components.

How much signal are we going to loose?

$$\sum_{i=k+1}^d \text{Var}(w^T x) = \frac{1}{n} \sum_{i=k+1}^d \lambda_i$$

number of components

Application: Eigenfaces



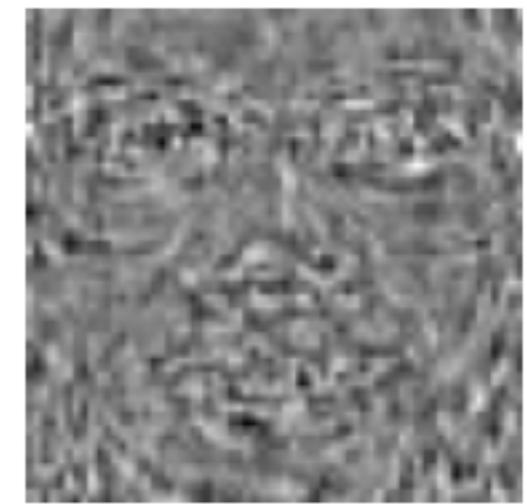
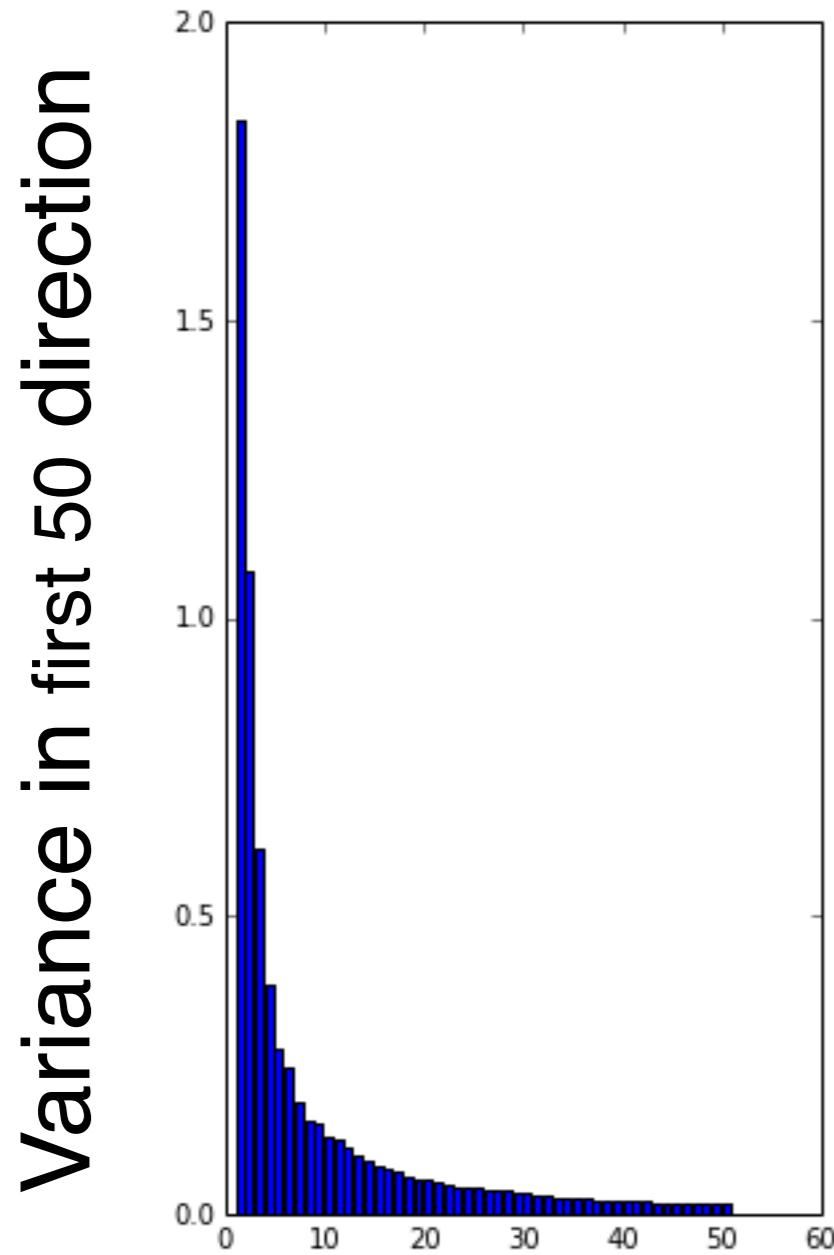
Idea: Faces look very similar in comparison to other random images. How many principle components would we need to accurately describe all faces?

An 64x64 pixel image of a face can be represented as a 4096 dimensional vector where each entry has the pixel's grayscale value.

Reference: Turk, Matthew A and Pentland, Alex P. Face recognition using eigenfaces. Computer Vision and Pattern Recognition, 1991.

Eigenfaces

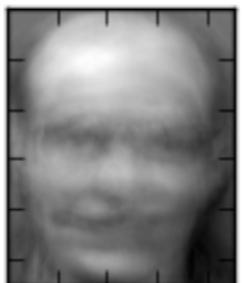
The principle components are directions in our faces space.
Thus, each principle component is a face representation, too.



Eigenfaces

Reconstruction AT&T Facedatabase

Eigenvectors #10



Eigenvectors #30



Eigenvectors #50



Eigenvectors #70



Eigenvectors #90



Eigenvectors #110



Eigenvectors #130



Eigenvectors #150



Eigenvectors #170



Eigenvectors #190



Eigenvectors #210



Eigenvectors #230



Eigenvectors #250



Eigenvectors #270



Eigenvectors #290



Eigenvectors #310



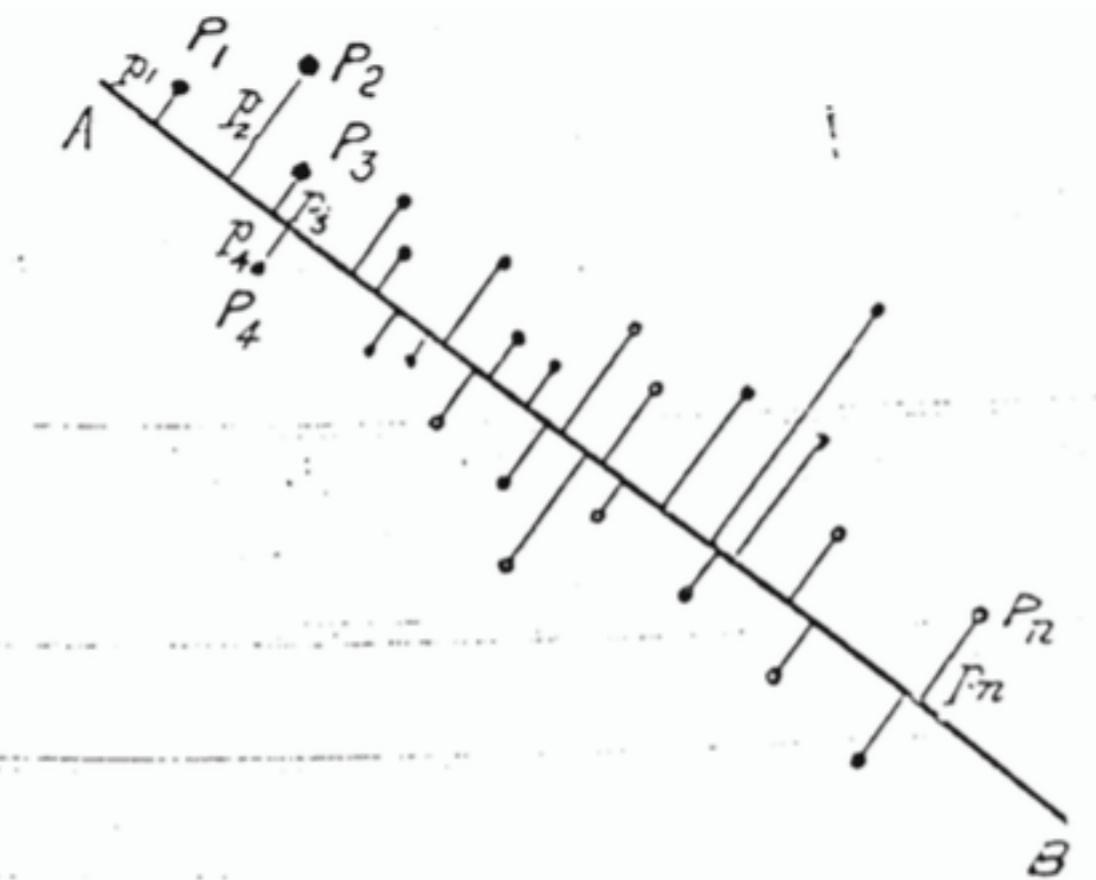
Application: Denoising

We can use PCA to denoise data:

Step 1: Reduce dimensionality to filter out
“noise” directions: $(w_1^T x, \dots, w_k^T x)^T$

Step 2: Project back into original space:

$$\sum_{i=1}^k (w_i^T x) w_i$$



Step 3: Undo centering:

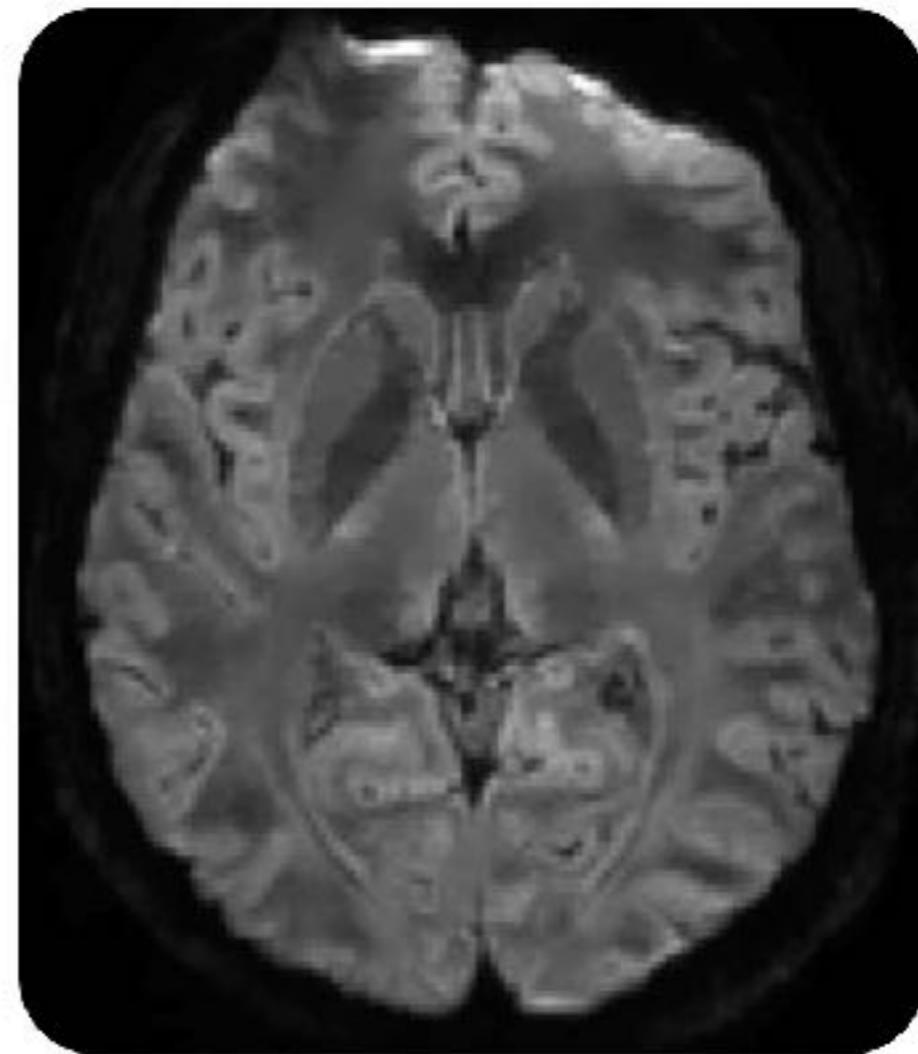
$$\sum_{i=1}^k (w_i^T x) w_i + \sum_{i=1}^n x_i$$

Application: Denoising

Original DWI



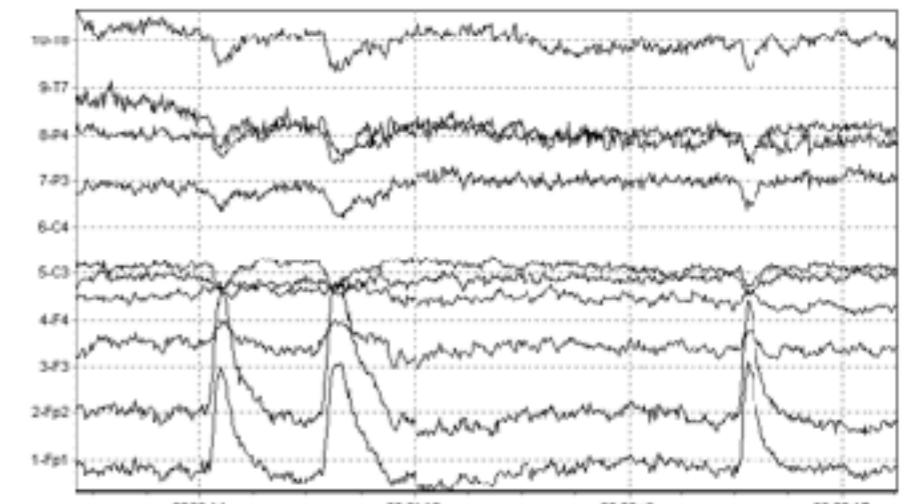
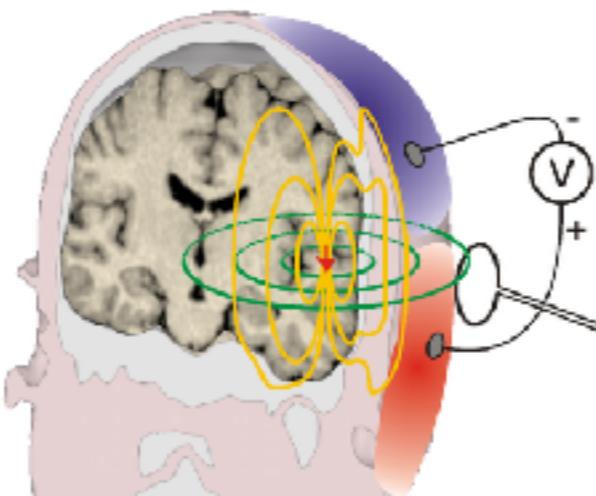
Denoised DWI using the LPCA filter



Mann et al., 2013

Application: EEG artifacts

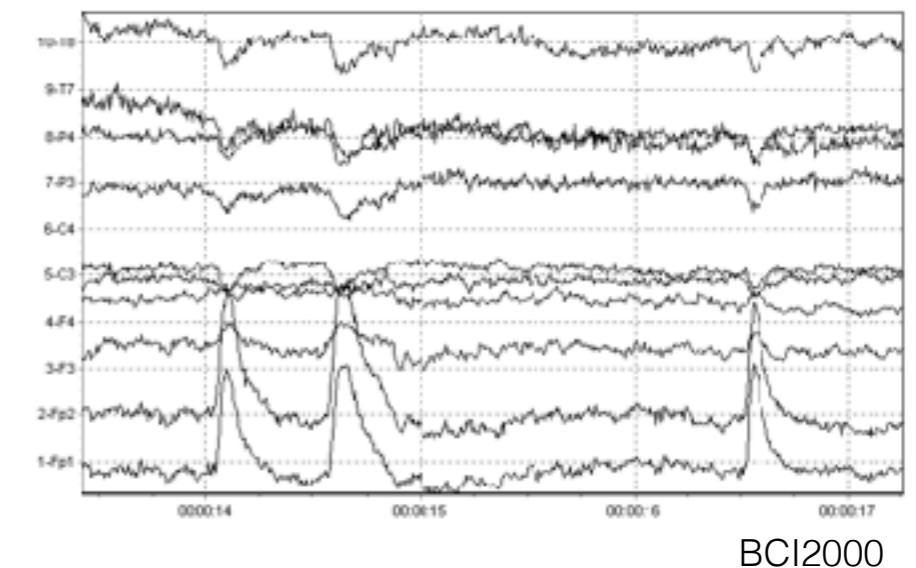
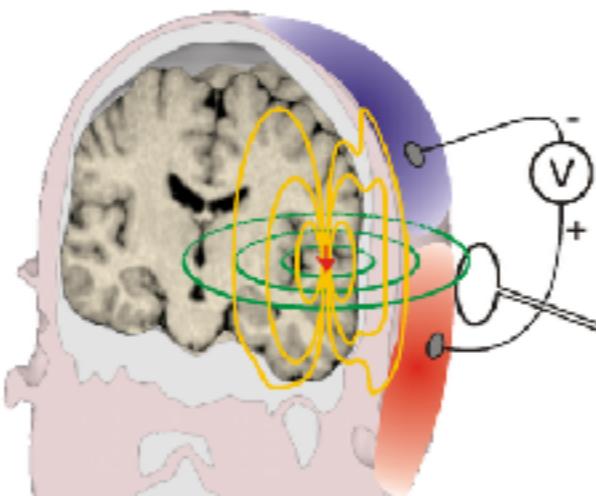
In electroencephalographic (EEG) recordings, eye blink artifacts can be stronger than the neuronal activity.



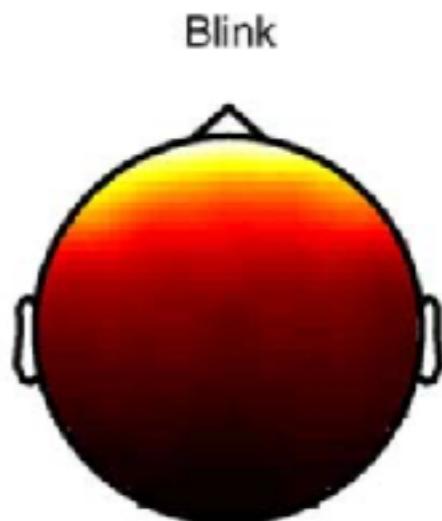
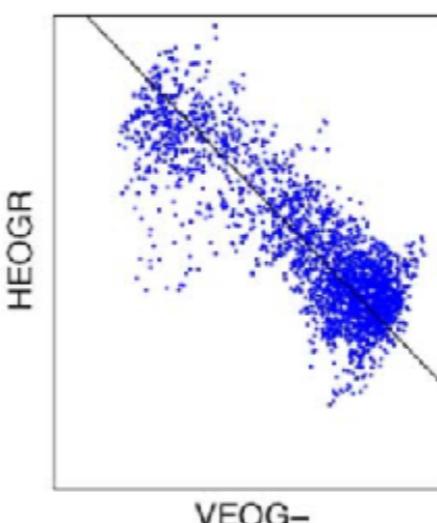
BCI2000

Application: EEG artifacts

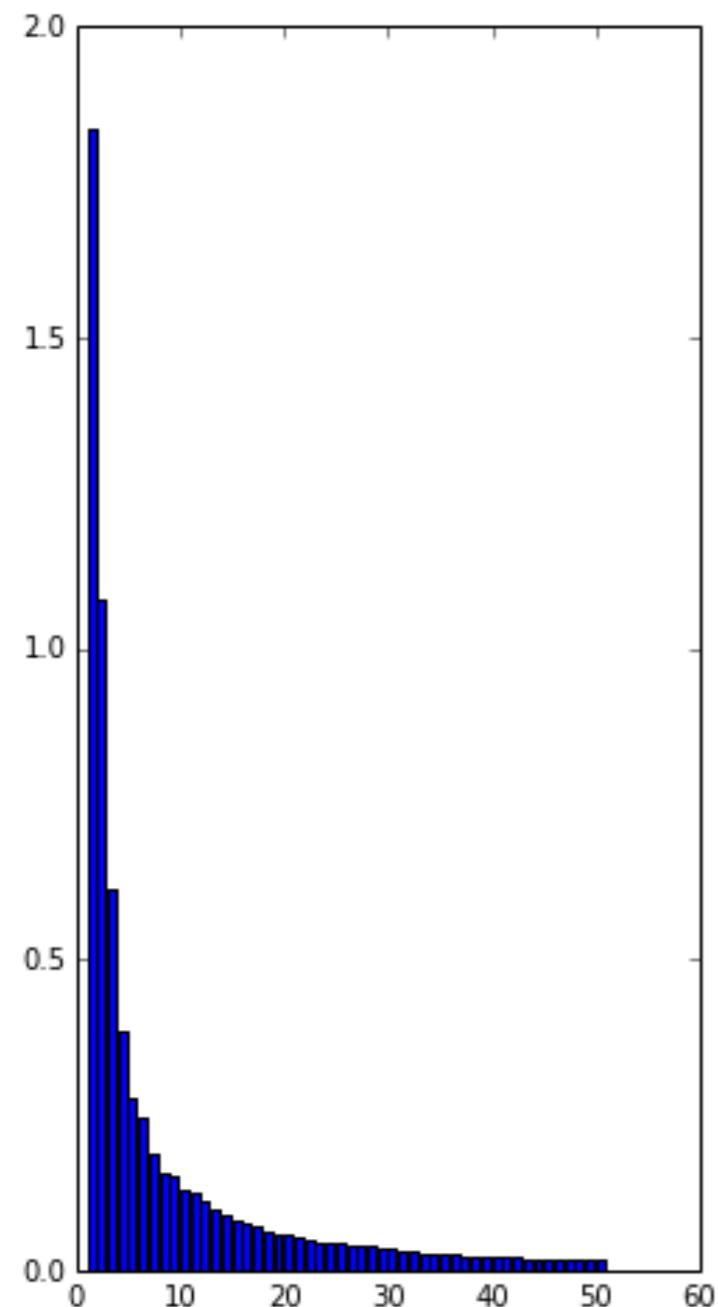
In electroencephalographic (EEG) recordings, eye blink artifacts can be stronger than the neuronal activity.



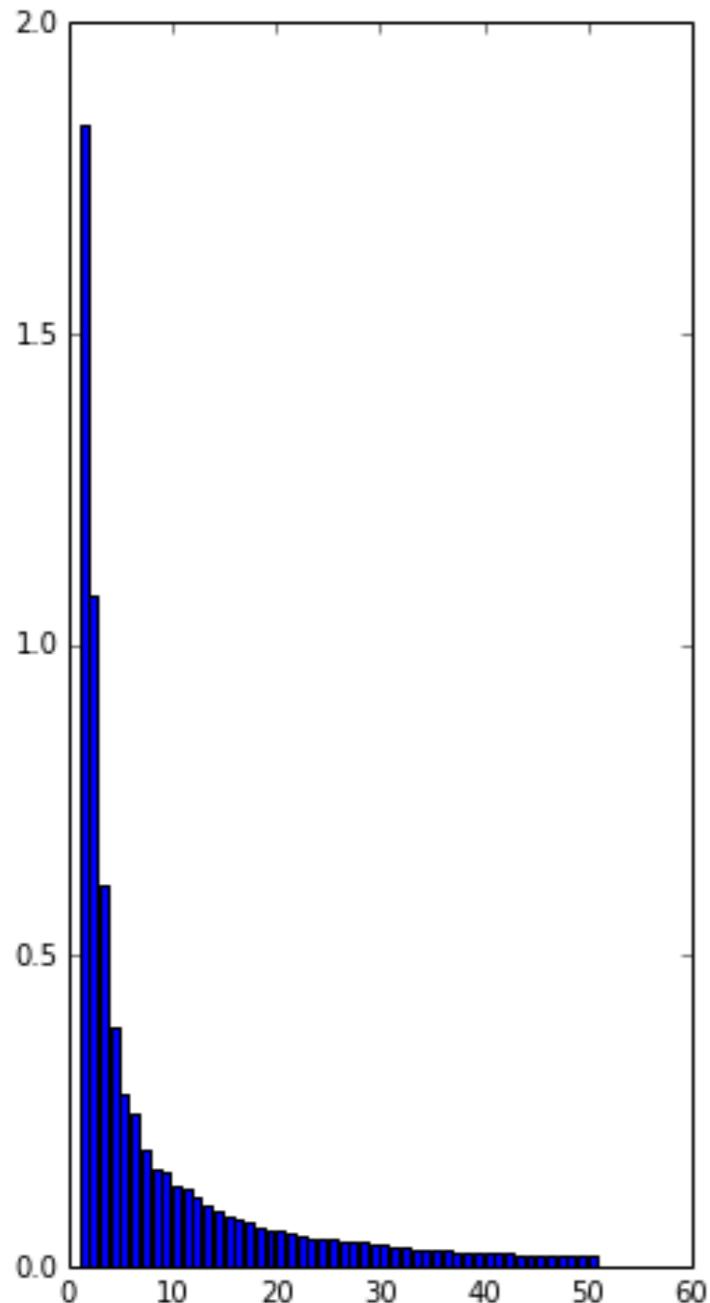
→ reasonable to remove first principal components



How many principal components?



How many principal components?



Heuristics

How robust is PCA to outliers?

