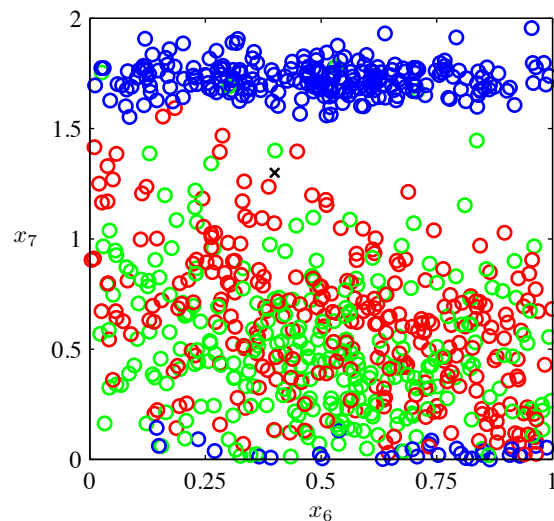
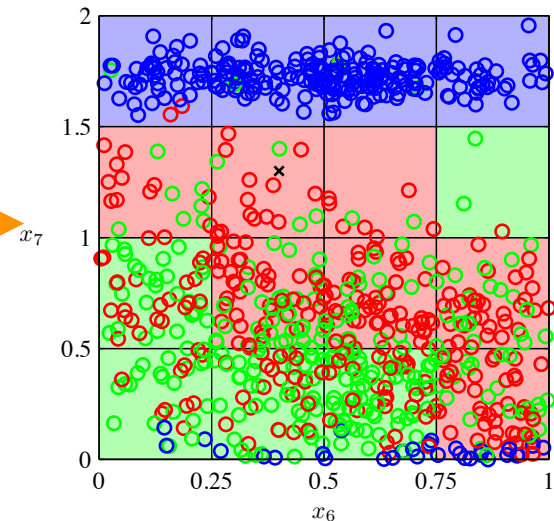


# What is curse of dimensionality?

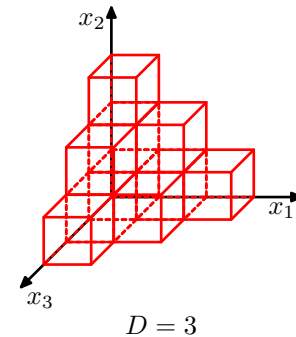
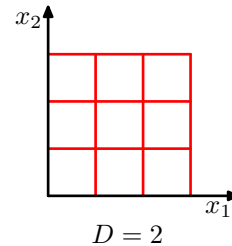
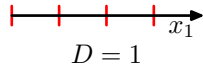
**Ex: a simple classification scheme (related to decision tree)**



divide up  
into small cells



# # of cells grows exponentially



# of cells  $r^D$

$r$ : number of divisions in each dimension

That is a lot even if  $D$  is just moderately large!

So to cover the whole space reasonably well,  
you need exponentially number of training data points!

# Another example

## Nearest neighbor classification

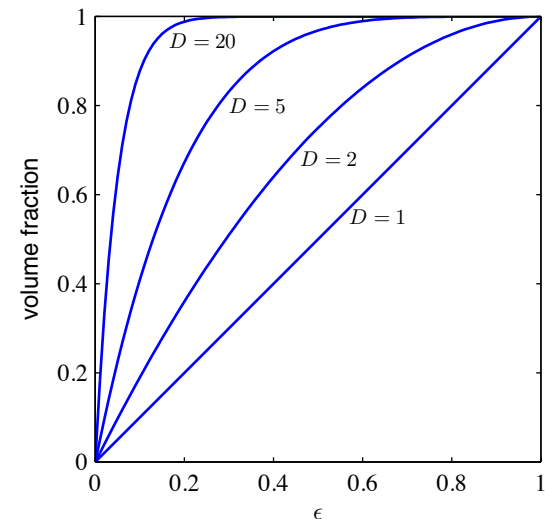
Say we want points of different classes are at least  $\epsilon$  away from each other

To put things in scale (otherwise, you can scale  $\epsilon$  arbitrarily), assume all points are at most 1 away from origin.

How many points between 1 and  $1-\epsilon$  from the origin?

$$1 - (1 - \epsilon)^D$$

Namely, when  $D$  is high, you cannot figure out data points that are  $\epsilon$  away even  $\epsilon$  is pretty big



# An overview of dimensionality reduction

## Parametric approaches

Linear: **PCA**, Fisher LDA, NMF, random projection, CCA, etc

Nonlinear: Neural networks, generative topological mapping, ICA

## Nonparametric approaches

kernelized version of parametric methods: **k-PCA**, kICA, kCCA

manifold learning

Gaussian processes

# Linear dimensionality reductions

## Many examples

Principal component analysis (PCA)

Fisher discriminant analysis (FDA)

Nonnegative matrix factorization (NMF)

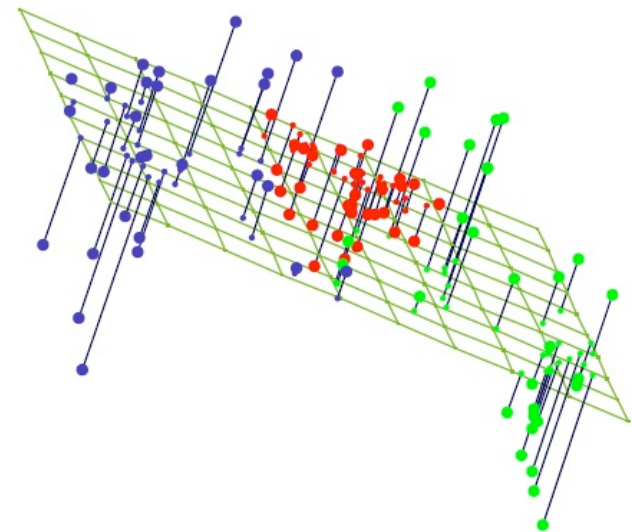
## Framework

$$\mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{y} \in \mathbb{R}^M$$

$$D \gg M$$

$$\mathbf{y} = \mathbf{U}^\top \mathbf{x}$$

linear transformation of original  
space

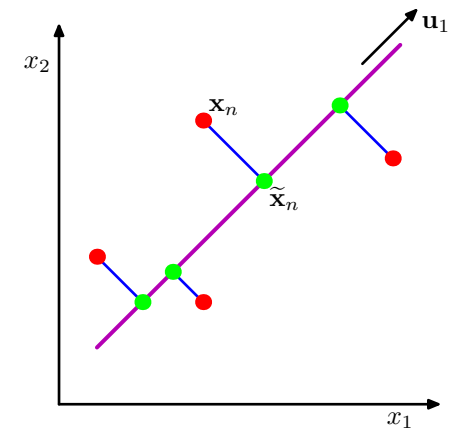
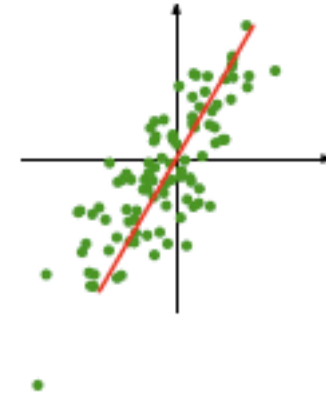


# PCA: first principal component

## Objective:

maximize variance of points projected on the line

## Derivation on whiteboard



# The framework of PCA

## Assumption:

Centered inputs (if not, subtract mean)

$$\mathbf{x} \in \mathbb{R}^D \quad \sum_i x_i = 0$$


Projection into subspace

$$\mathbf{U} \in \mathbb{R}^{D \times M} \quad \mathbf{y}_i = \mathbf{U}^T \mathbf{x}_i$$
$$\mathbf{U}^T \mathbf{U} = \mathbf{I}$$

## Solution

Computer covariance matrix

each row is a data sample


$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

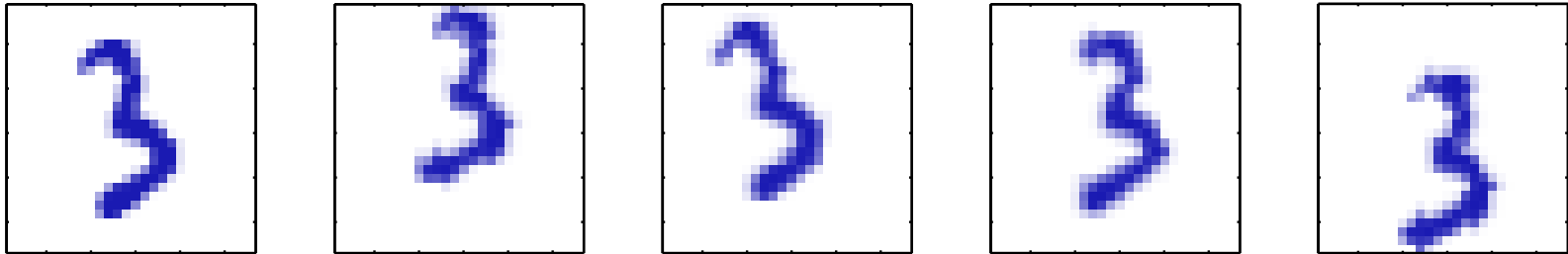
Diagonalize S:  $(\mathbf{u}_d, \lambda_d)$   $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_M$

Use top D eigenvectors to form U

$$\mathbf{U} = (\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_M)$$

# Examples of running PCA

## Original Images



## Eigenvectors

they look like blurred original images

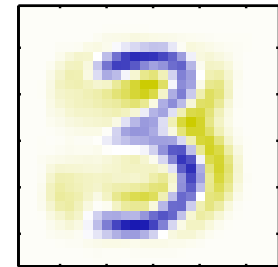
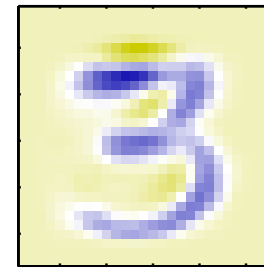
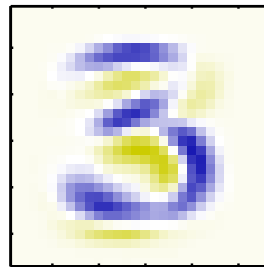
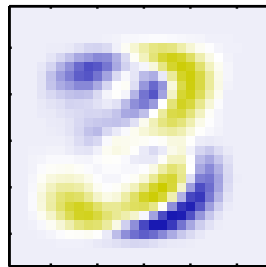
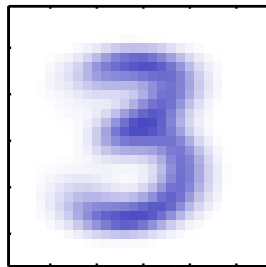
Mean

$$\lambda_1 = 3.4 \cdot 10^5$$

$$\lambda_2 = 2.8 \cdot 10^5$$

$$\lambda_3 = 2.4 \cdot 10^5$$

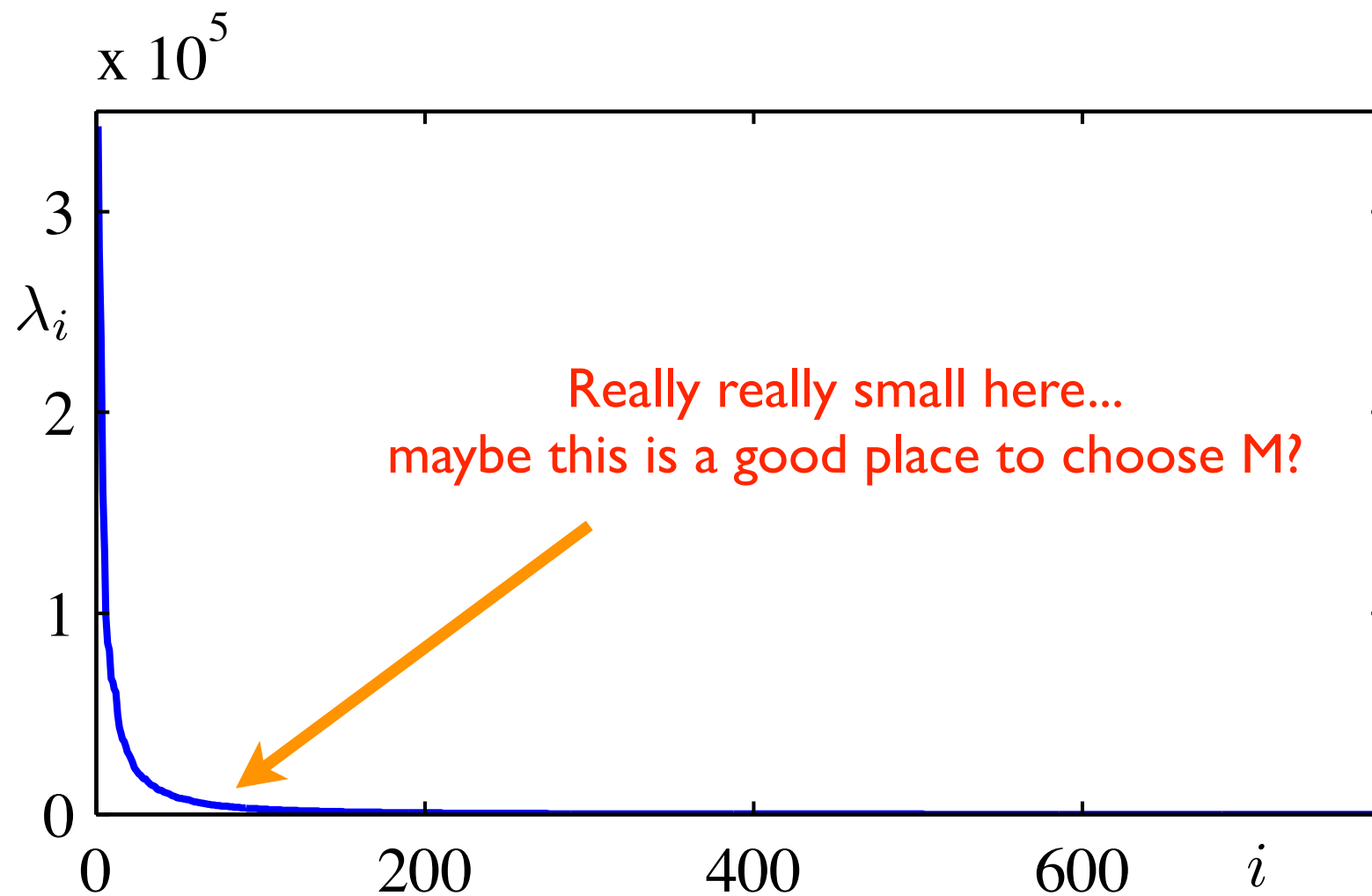
$$\lambda_4 = 1.6 \cdot 10^5$$



Used to centralize inputs



# Eigenspectrum of the covariance matrix

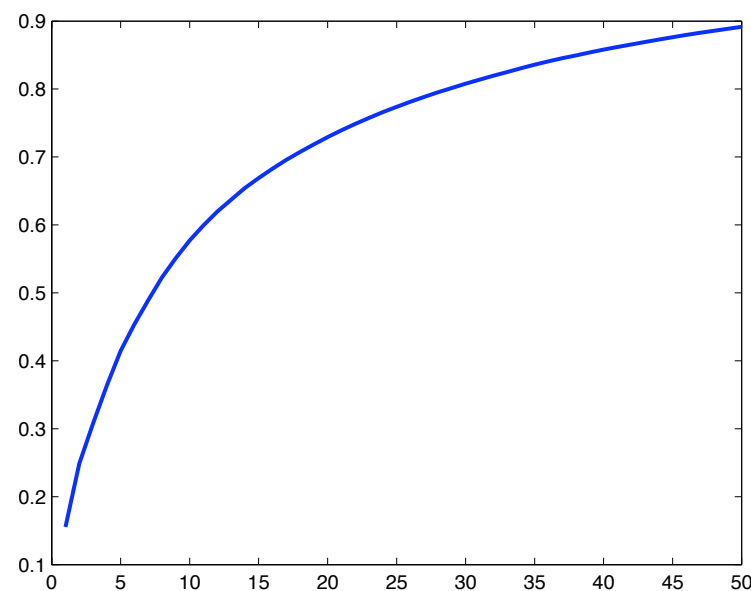
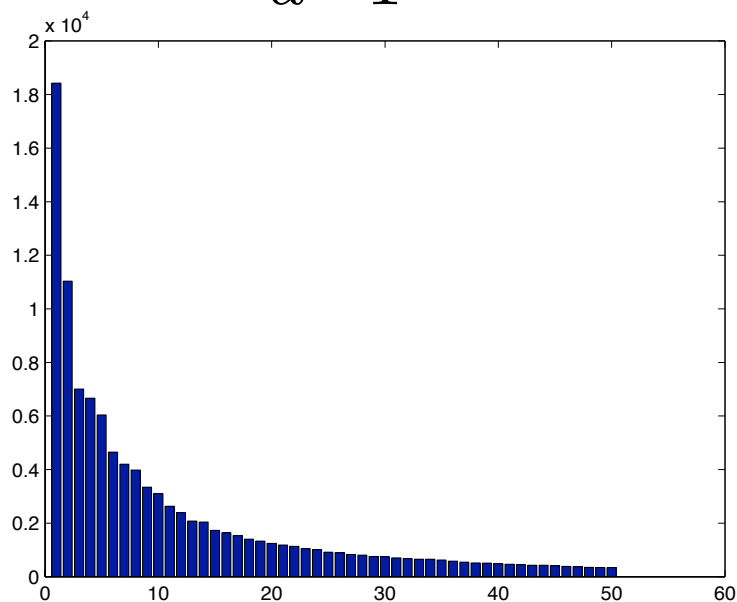


(a)

# A slightly sensible approach

common choice is 95% or 90%

$$\frac{\sum_{d=1}^M \lambda_d}{\sum_{d=1}^D \lambda_d} \geq \text{Threshold}$$



# Application of PCA

## Preprocessing

Diagonalize data

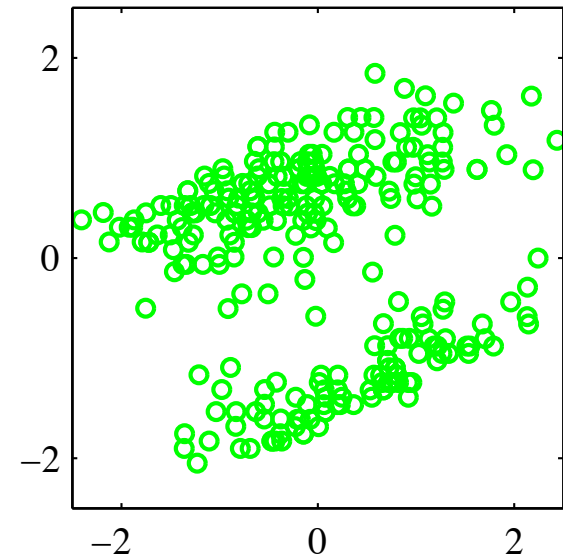
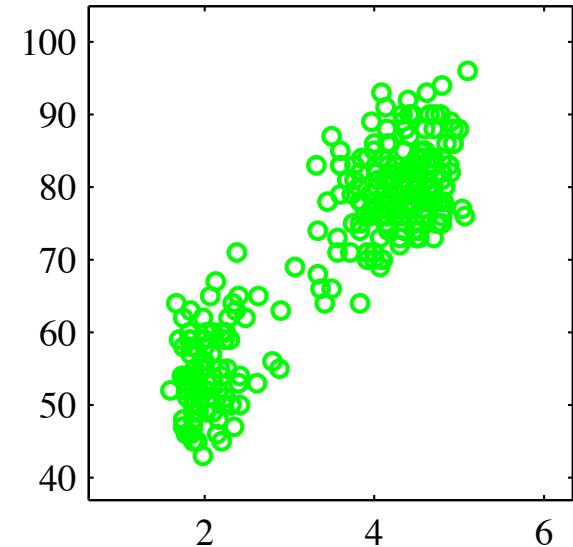
$$y_i = U^T x_i$$

Normalize data (whitening)

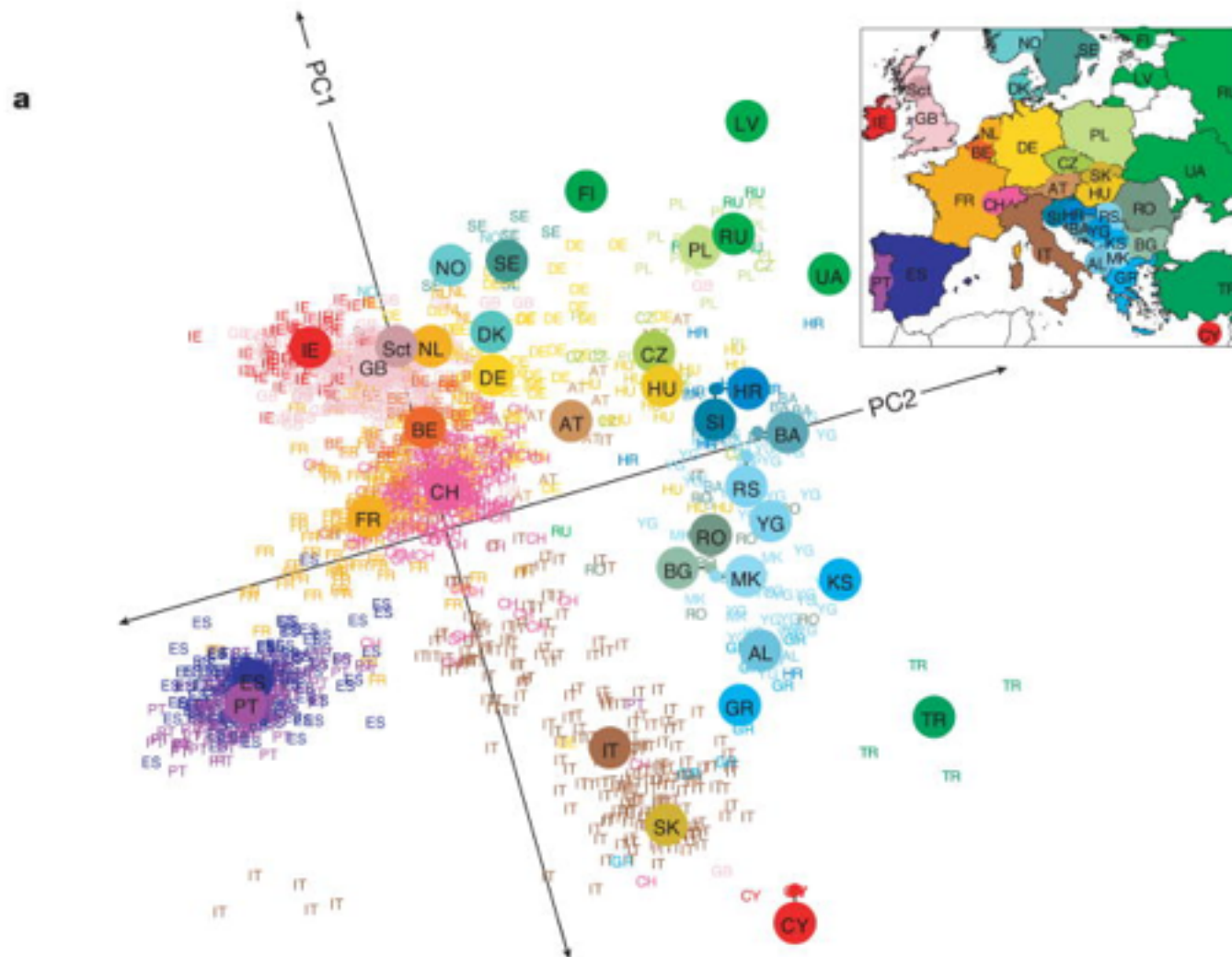
$$y_i = \lambda^{-1/2} U^T x_i$$

Benefits:

- 1) depress noisy features
- 2) couple with other models



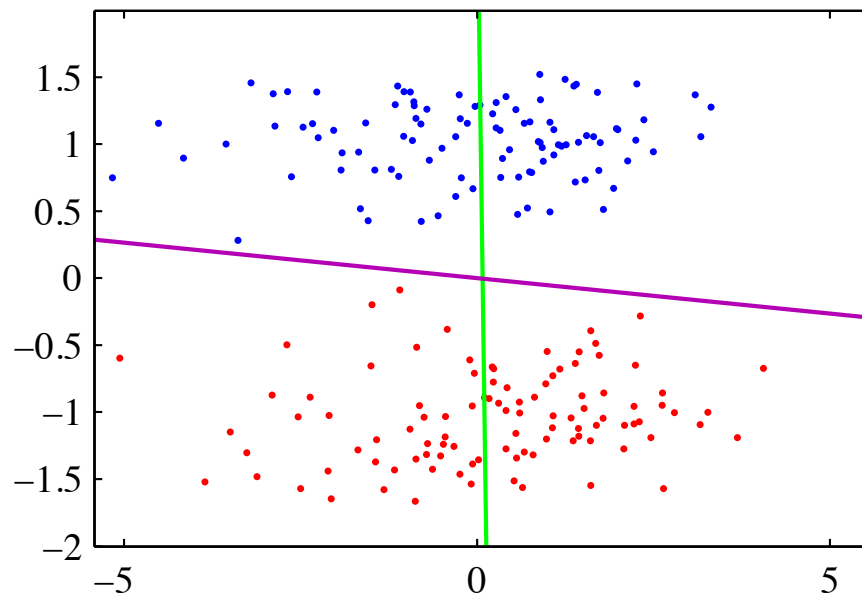
# Visualizing data with PCA



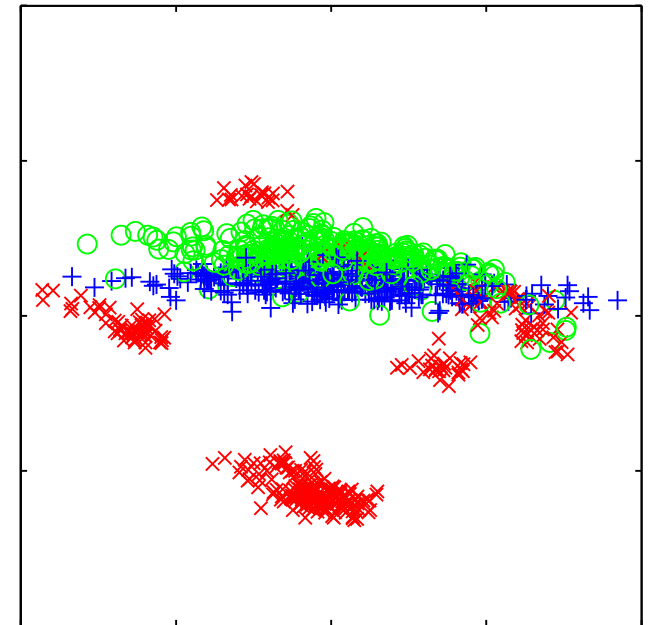
# Does PCA help for classification?

**Not always!**

no help



help



Visualization

# The second interpretation of PCA

## minimum reconstruction error

Consider a basis of a subspace with M-dimension where  $M < D$

$$\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M \in \mathbb{R}^D$$

Consider the approximation error by linearly combining the basis

$$\sum \left\| \mathbf{x}_n - \sum_m \alpha_{nm} \mathbf{u}_m \right\|_2^2$$

If we minimize the error by optimizing both the basis and the coefficients, we get

$$\alpha_{nm} = \mathbf{x}_n^T \mathbf{u}_m$$

$\mathbf{u}_m$  is the m-th largest eigenvector of S(ample) covariance matrix

# Unexplained variance and residual error

## Unexplained variance

Each projection direction  $u_i$  contributes  $\lambda_i$  variance

With  $D$ -dimensional data,  $M$  projection directions, the “leftover” is

$$\sum_{d=1}^D \lambda_d - \sum_{m=1}^M \lambda_m = \sum_{m=M+1}^D \lambda_m$$

## Reconstruction error

$$\sum \left\| \mathbf{x}_n - \sum_m \alpha_{nm} \mathbf{u}_m \right\|_2^2$$

The two are exactly the same!!!

# Issues with PCA

## Choose dimensionality

ad hoc approach: we have seen that

more principled approach: bayesian PCA, etc

## Missing values

what if data is missing?

Probabilistic PCA: a little bit later



# Computational complexity

## For large input dimensionality $D$

the covariance matrix is  $D \times D$

finding eigenvalue and eigenvector for the covariance matrix can be very expensive

# The trick

section 12.1.4

## Same set of eigenvalues

This is called Multidimensional scaling (MDS)

$$\frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{u} = \lambda \mathbf{u} \rightarrow \frac{1}{N} \mathbf{X} \mathbf{X}^T \boxed{\mathbf{X} \mathbf{u}} = \lambda \boxed{\mathbf{X} \mathbf{u}}$$

## Same low dimensional representation

$$\mathbf{u} \propto \mathbf{X}^T \mathbf{v}$$


## Different computational cost

PCA's matrix scales quadratically in D

MDS's matrix scales quadratically in N

Big win for MDS when D is much greater than N !

## Use nonlinear kernel?

$$\frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{u} = \lambda \mathbf{u} \rightarrow \frac{1}{N} \mathbf{X} \mathbf{X}^T \mathbf{X} \mathbf{u} = \lambda \mathbf{X} \mathbf{u}$$



Gram matrix  
(matrix of inner products!)

Can we use kernel here?  
Just like we did for kernel regression?

## Why? ---- Linear is not enough

Issue:  
projection on any 1D line  
will cause mixing of data  
from different classes.



# We need nonlinear projection

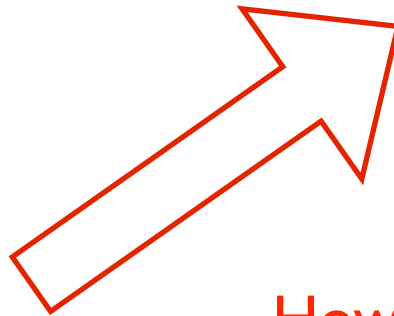
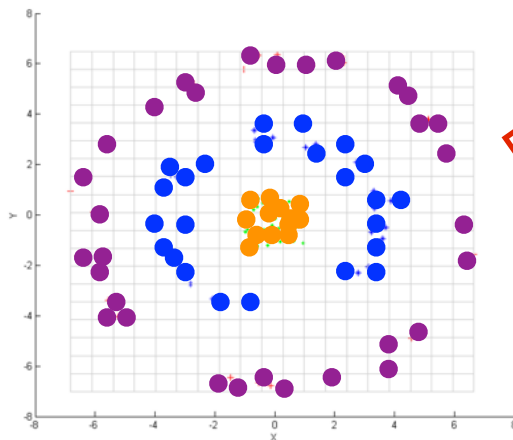
## Intuition:

Find a good space through

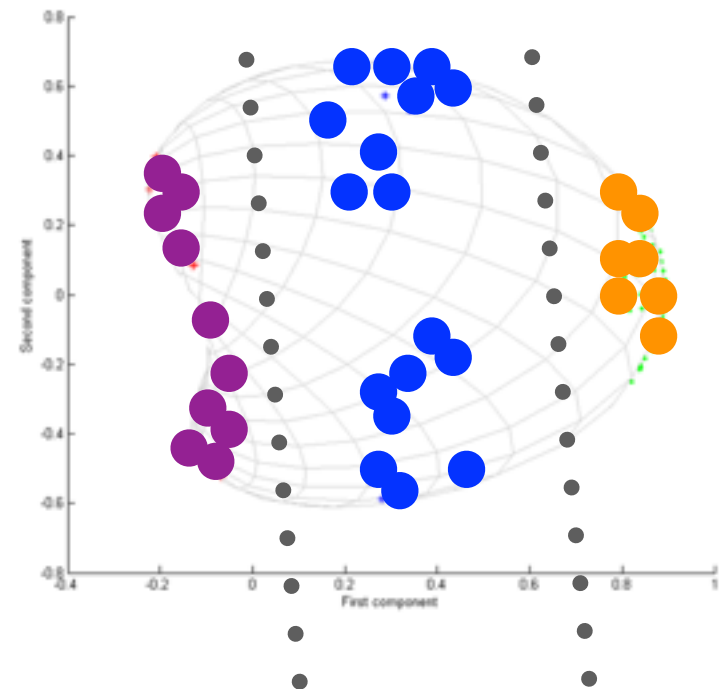
**nonlinear** mapping;

Then do **linear** projection

(as in PCA)



✓ Data separated!



How to find this  
nonlinear mapping?

# Details of kernel PCA

## Derivation

see note

## Demo