# Reverse mode differentiation

Piece of computation:



Backpropagate derivatives wrt final scalar output

For any $Z$,
$$\bar{Z}_{ij} = \frac{\partial \text{output}}{\partial Z_{ij}}$$

## Use standard rules:
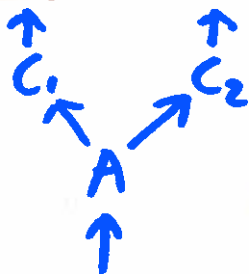
$C = \cos A \implies \bar{A} = \bar{C} \odot \sin A$ — Matlab .* Python *

$C = AB \implies \bar{A} = \bar{C} B^T, \quad \bar{B} = A^T \bar{C}$

$C = A + B \implies \bar{A} = \bar{C}, \quad \bar{B} = \bar{C}$

$C = A^T \implies \bar{A} = \bar{C}^T$

...

Stored in forward pass

Passed in by backprop

## Multiple children



$$\bar{C}_1 \quad \bar{C}_2$$

$$\bar{A} = \bar{A}_1 + \bar{A}_2$$

Apply rules separately for children and add

# Matrix multiplication

$$C = A B \qquad O(LMN)$$

$L \times N \qquad L \times M \quad M \times N$

$$C_{\ell n} = \sum_m A_{\ell m} B_{mn}$$

LN terms $\qquad$ cost $O(M)$ each

Square matrix-matrix multiply $\quad O(N^3)$

There are $O(N^2)$ numbers in the matrices

$\left[\text{See also} \quad \text{tutorial 5, } Q1 \right]$

# Autoencoder   (Unsupervised)



$\underline{f}$

$h$

$\underline{x}$
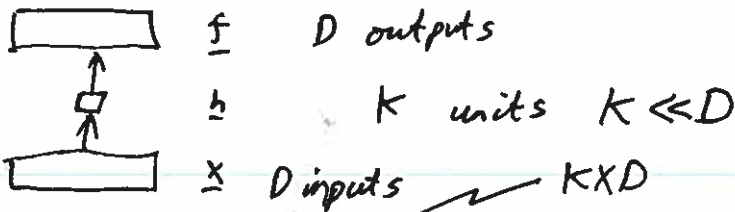
Learning task

$$\underline{f} \approx \underline{x}$$

not
useful

```
def autoencode (x):
    return x
    h = np.dot (I, x)
    f = np.dot (I, h)
    return f
```

# Dimensionality Reduction



$\underline{f}$     D outputs

$h$      K units   $K \ll D$

$\underline{x}$   D inputs       $K \times D$

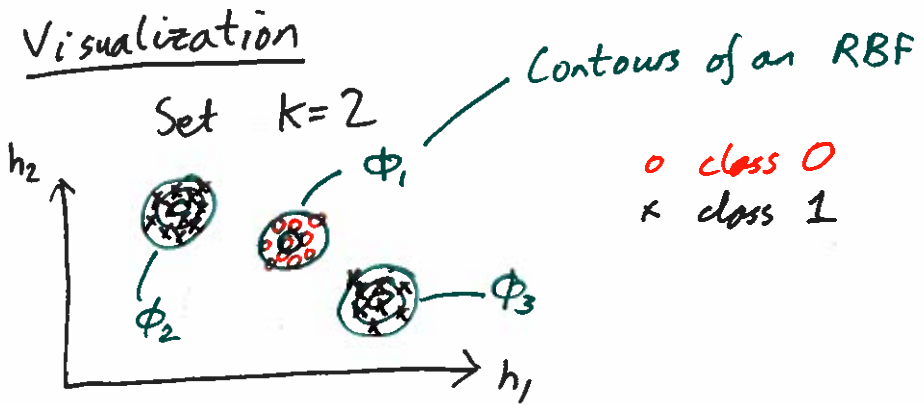$$\underline{h} = g^{(1)} \left( W^{(1)} \underline{x} + \underline{b}^{(1)} \right)$$

$$\underline{f} = g^{(2)} \left( W^{(2)} \underline{h} + \underline{b}^{(2)} \right) \leftarrow \text{decoder}$$
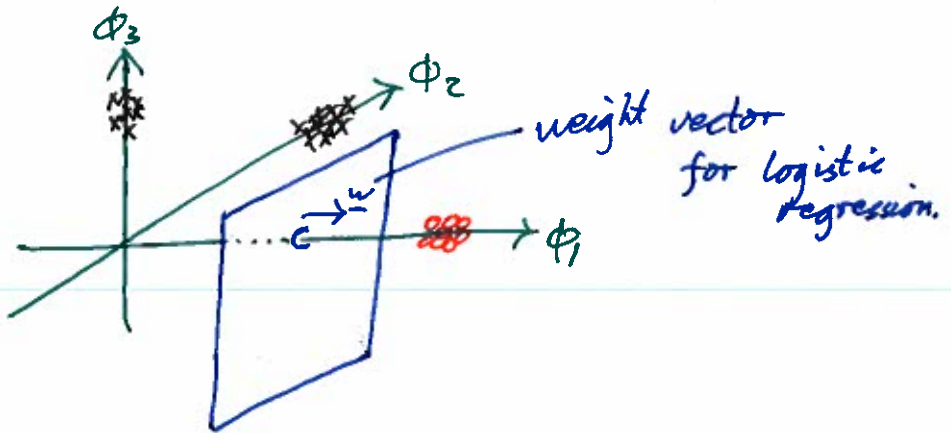
Encoder could be useful

Train this thing on a huge dataset

Use $\underline{h}(\underline{x}; W^{(1)}, \underline{b}^{(1)})$ to transform data

# Visualization

### Set  k=2

Contours of an RBF

o  class 0
x  class 1



We might want to increase dim. of data

weight vector
for logistic
regression.

# Sparse Autoencoder

Encourage most elements of $\underline{h}$
   to be close to zero — sparse

# Denoising Autoencoder

While training   we mask out (delete)
some of the inputs,   set $x_d$ to zero
   $\wedge$
   some

$\underline{m}$   mask vector, of random 0's and 1s

Cost on
an example
$$\| \underline{f}(\underline{x}^{(n)} \odot \underline{m}) - \underline{x}^{(n)} \|^2$$

Cost function   $\displaystyle\sum_{\underline{m}} p(\underline{m}) \frac{1}{N} \sum_{n=1}^{N} \| \underline{f}(\underline{x}^{(n)} \odot \underline{m}) - \underline{x}^{(n)} \|^2$

Monte Carlo
   $\approx$   pick random $\underline{m}$
      random $n$

# Principal Components Analysis (PCA)

Linear auto-encoder $\quad g^{(1)}(a) = g^{(2)}(a) = a$

$$\underline{h} = V^T(\underline{x} - \bar{\underline{x}})$$

Training set mean

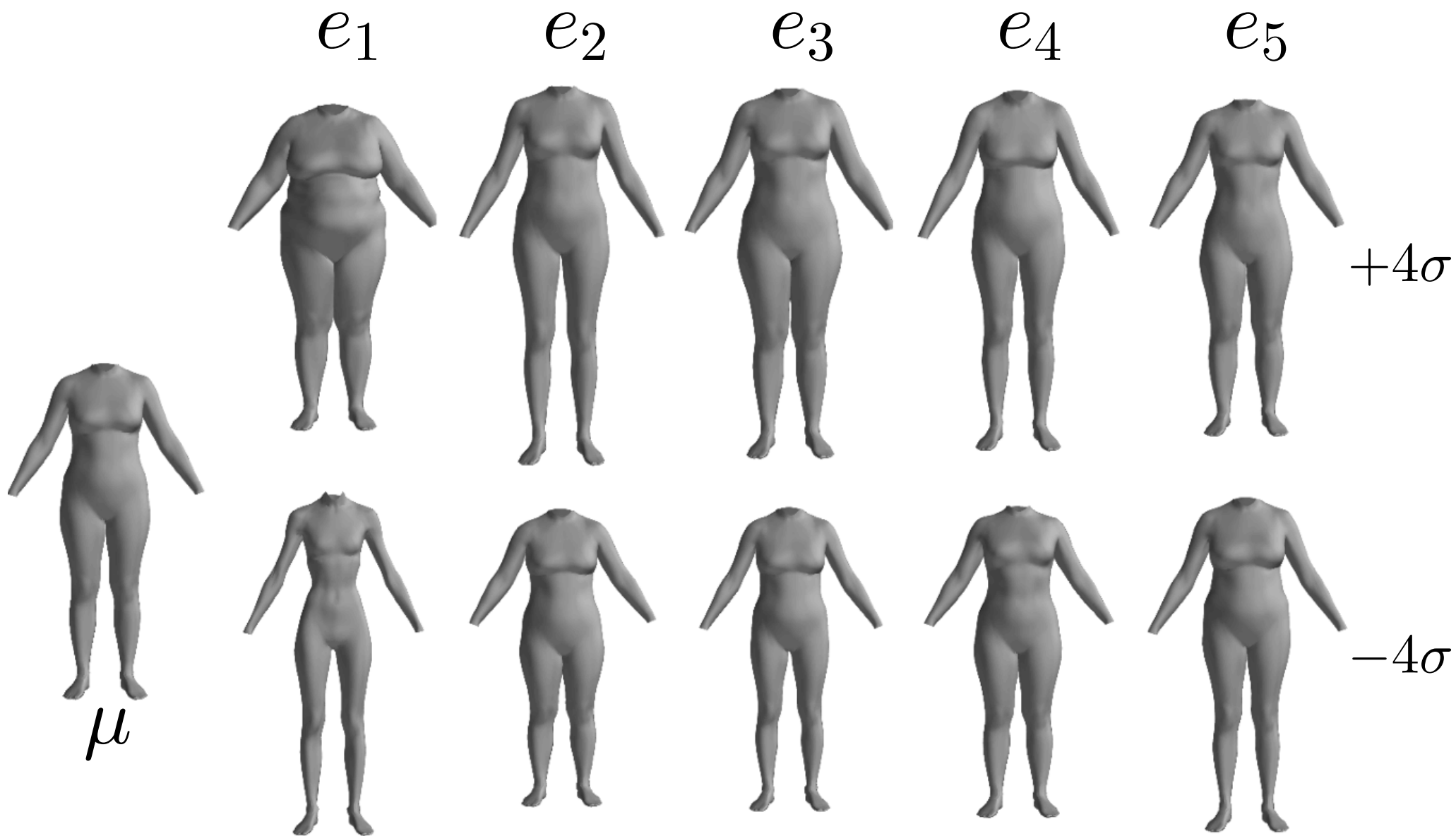$$\underline{f} = V\underline{h} + \bar{\underline{x}}$$

shared $D \times K$ matrix

## PCA advantages

- Fit columns of $V$ to be eigenvectors of the covariance of data
  (No SGD!)

- Same answer every time.

- The solutions for different $k$ they're nested

  $h_1(\underline{x})$ it's the same for all $k$
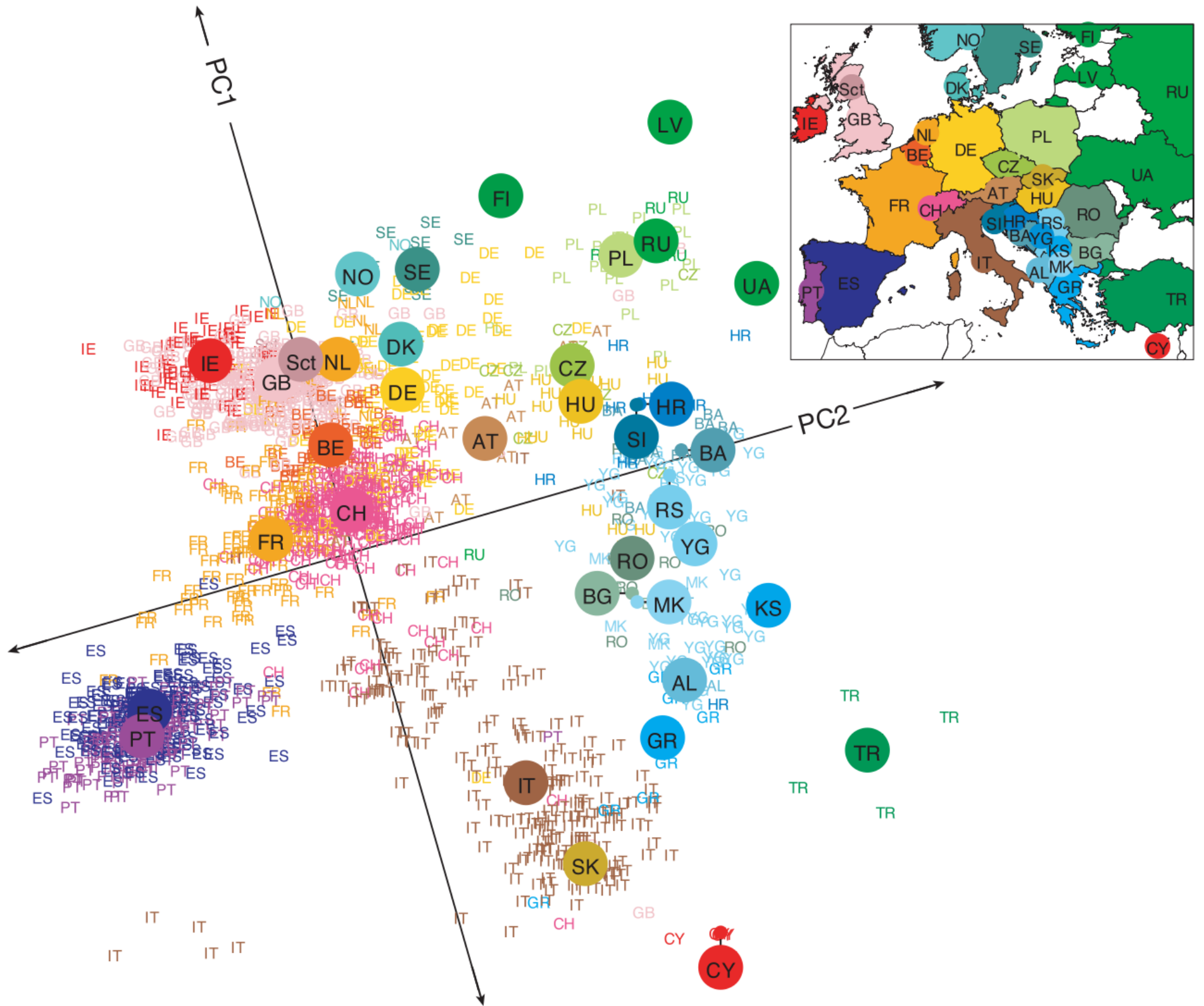  $h_2(\underline{x})$ " " " " " $k \geqslant 2$

# PCA applied to bodies



$e_1$    $e_2$    $e_3$    $e_4$    $e_5$

$+4\sigma$

$-4\sigma$

$\mu$

Freifeld and Black, ECCV 2012

# PCA applied to DNA

Carefully selected both individuals and features

1,387 individuals

197,146 single nucleotide polymorphisms (SNPs)

Each person reduced to two(!) numbers with PCA
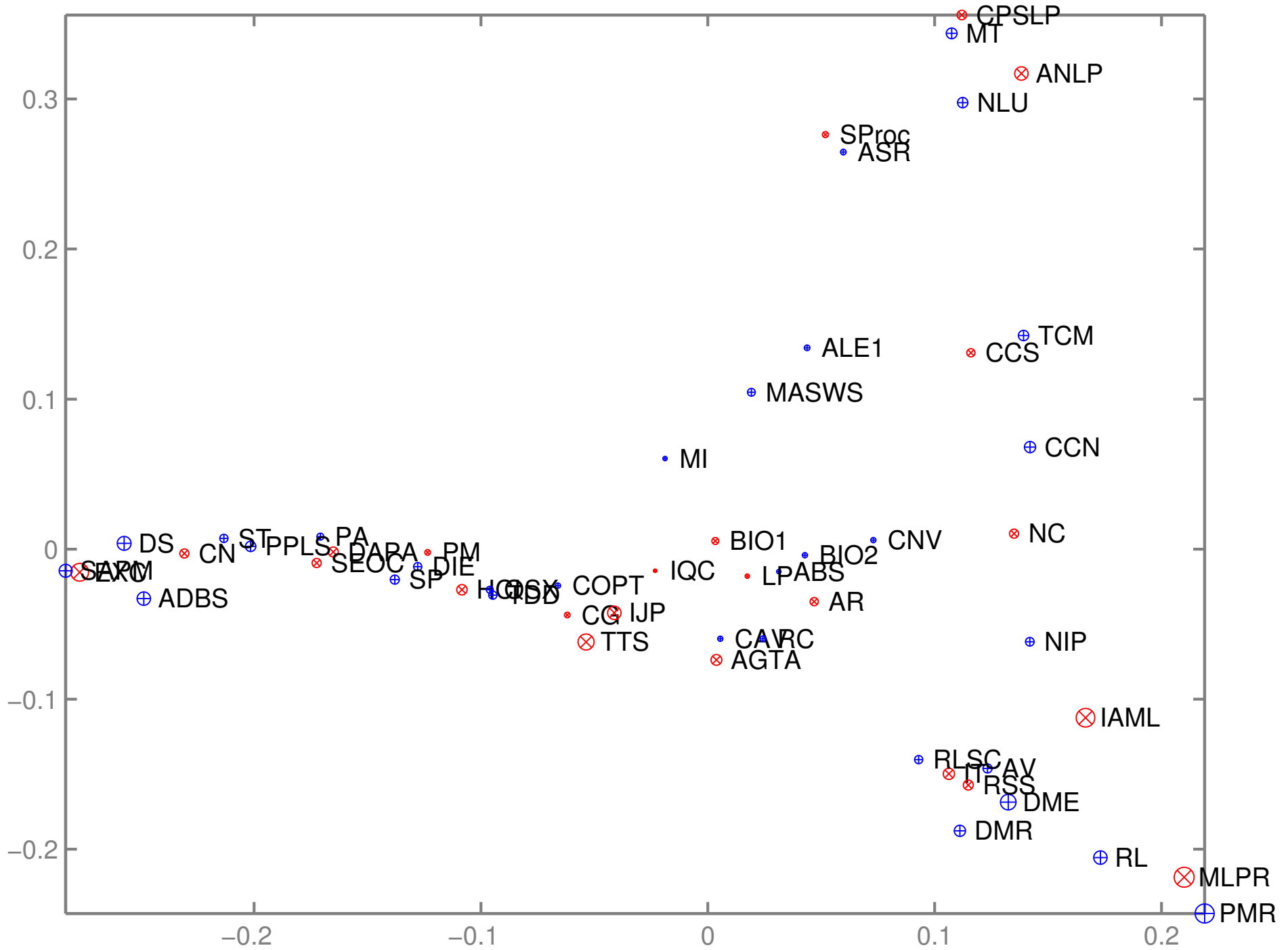
# MSc course enrollment data

Binary $S \times C$ matrix $M$

$M_{sc} = 1$,  if student $s$ taking course $c$

Each course is a length $S$ vector

. . . OR each student is a length $C$ vector

# PCA applied to MSc courses

# PCA applied to MSc students