

$$N \gg p$$

$$\text{Rank } r = p$$

Remember:

SVD and PCA

$$X_{N \times p} = Y_{N \times p}^T = \begin{pmatrix} \vec{y}_1^T \\ \vdots \\ \vec{y}_N^T \end{pmatrix}$$

$$\vec{y}_i \in \mathbb{R}^p$$

$$\tilde{Y}_{N \times p}^T = \left(I - \frac{1}{N} 11^T \right) Y^T$$

centered
data
matrix

$$\sum_{j=1}^p \sigma_j \vec{u}_j \vec{v}_j^T = \tilde{Y}_{N \times p}^T = U_{N \times p} \Sigma_{p \times p} V_{p \times p}^T$$

SVD

$$U = (\vec{u}_1, \dots, \vec{u}_p)_{N \times p}$$

$$V = (\vec{v}_1, \dots, \vec{v}_p)_{p \times p}$$

SVD provides
a new basis $\{\vec{v}_j\}_{j=1}^p$
orthonormal

our new
basis

\downarrow E-Y

closest
approximation
of rank (k)

For a
given k :

$$\hat{Y}_{N \times p}^T = \left[U_k \Sigma_k \right]_{N \times k} V_{k \times p}^T$$

$$\left\{ \begin{aligned} (Y_{PCA}) &= [U_k \Sigma_k]_{N \times k} \\ &\text{coordinates of the} \\ V_k &= (\vec{v}_1, \dots, \vec{v}_k) \quad \vec{v}_i \in \mathbb{R}^p \end{aligned} \right\}$$

$$\text{where } \Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k)$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k$$

ordered
singular
values.

Examples overlaid \longrightarrow

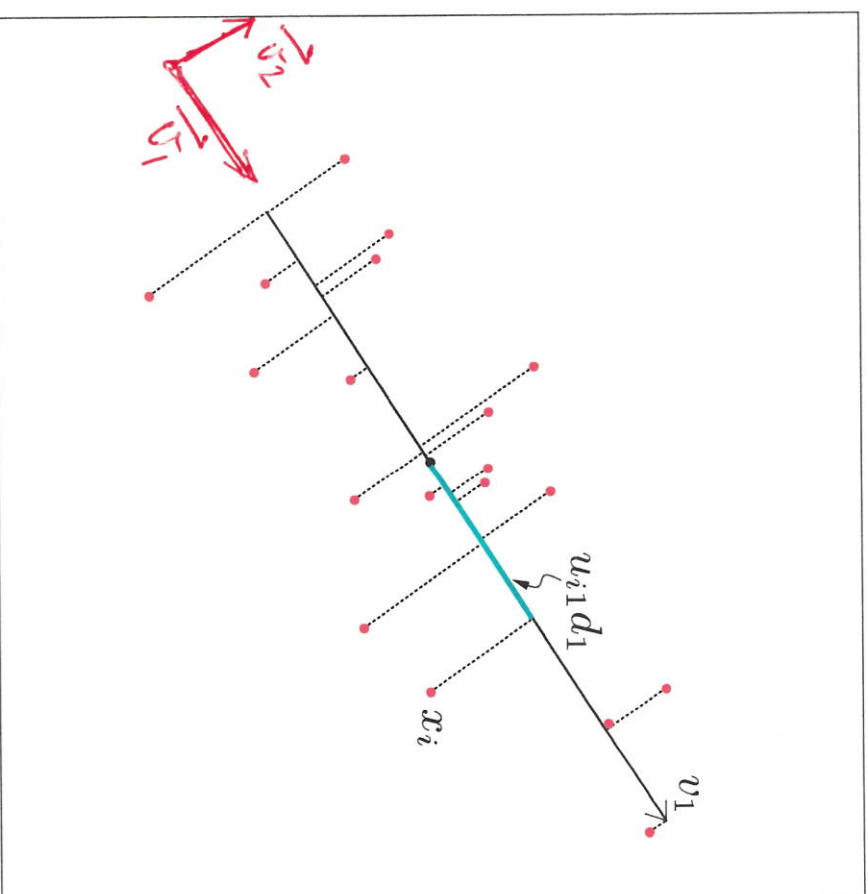
Illustration of PCA and SVD

\vec{v}_1 captures maximum
variability

From the SVD:

$$V = (\vec{v}_1 \vec{v}_2)$$

New basis



$$\vec{y}_i \in \mathbb{R}^2$$

$$\vec{y}_i = \sum_{i=1}^N \alpha_i^{(i)} \vec{v}_1 + \alpha_2^{(i)} \vec{v}_2$$

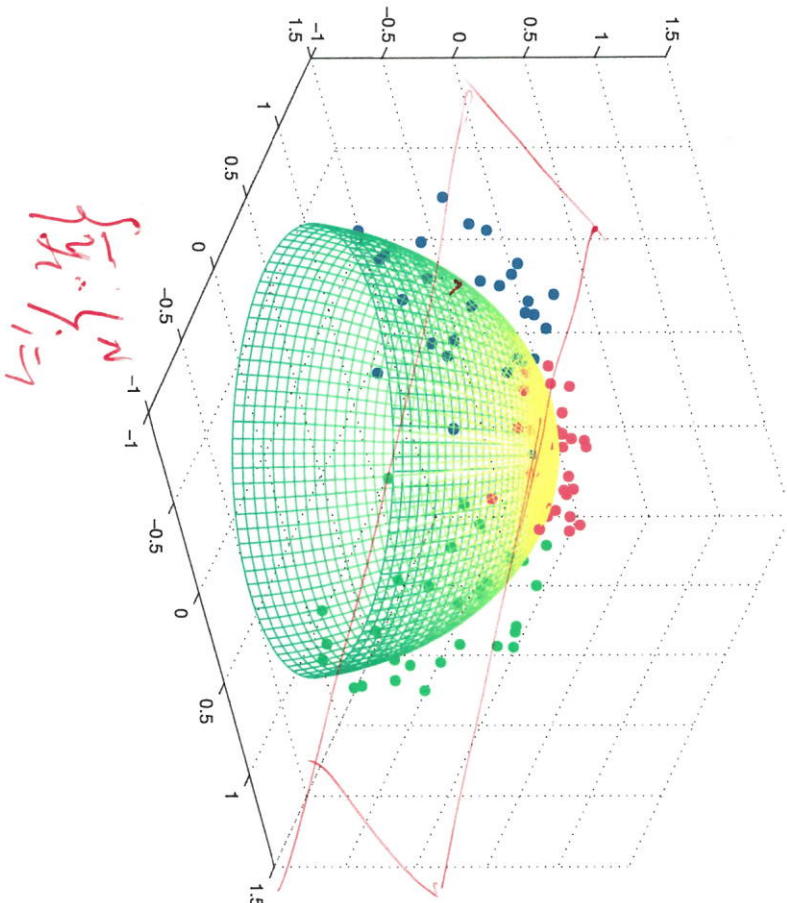
$$= \sigma_1 u_{i1} \vec{v}_1 + \sigma_2 u_{i2} \vec{v}_2$$

$$= (\vec{v}_1^T \vec{y}_i) \vec{v}_1 + (\vec{v}_2^T \vec{y}_i) \vec{v}_2$$

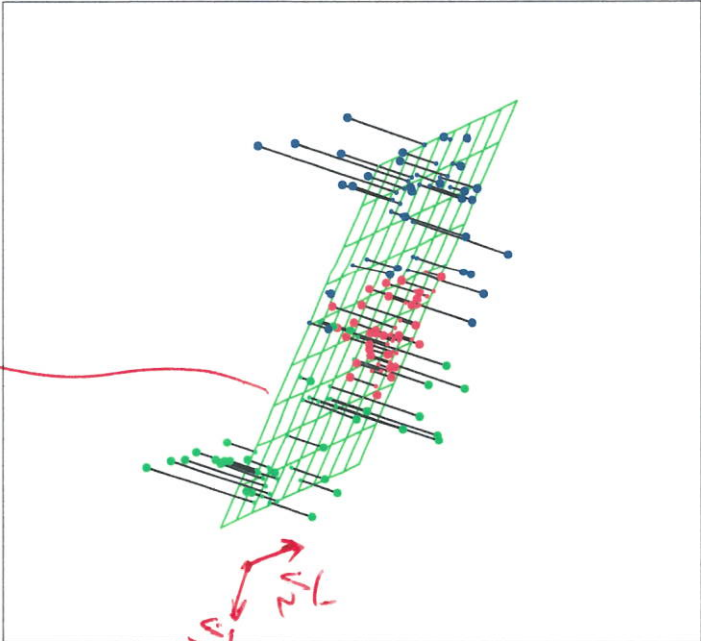
If we approximate by first component, we minimize the sum of perpendicular errors.

Illustration of PCA and SVD

$$\vec{x}_i \in \mathbb{R}^3$$



$$\sum_{i=1}^N \vec{x}_i \vec{x}_i^T$$

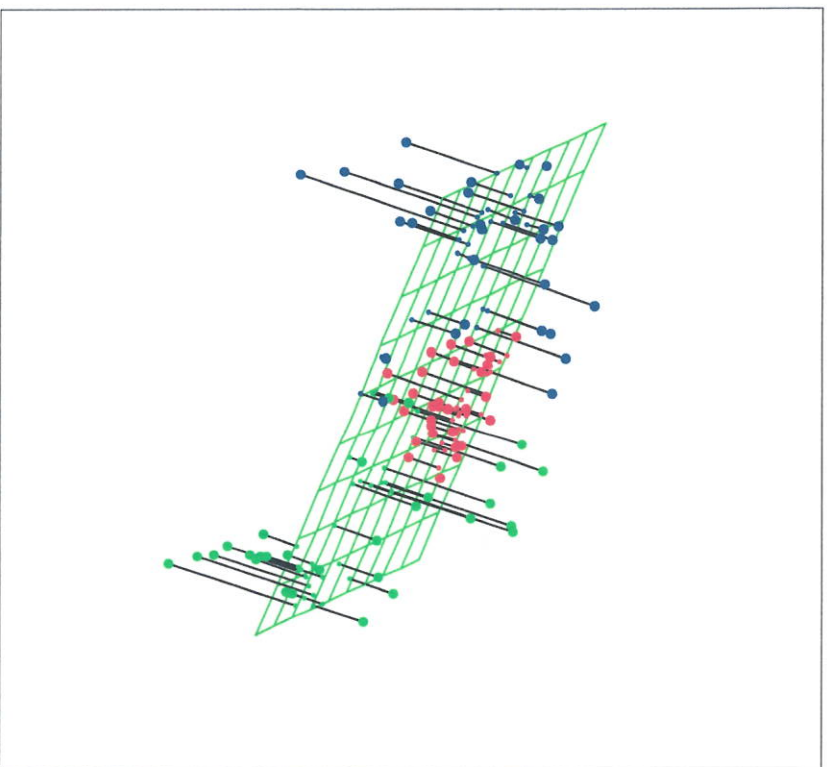


$$\vec{v}_1, \vec{v}_2$$

Best approximation given by hyperplane (\vec{v}_1, \vec{v}_2)

Illustration of PCA and SVD

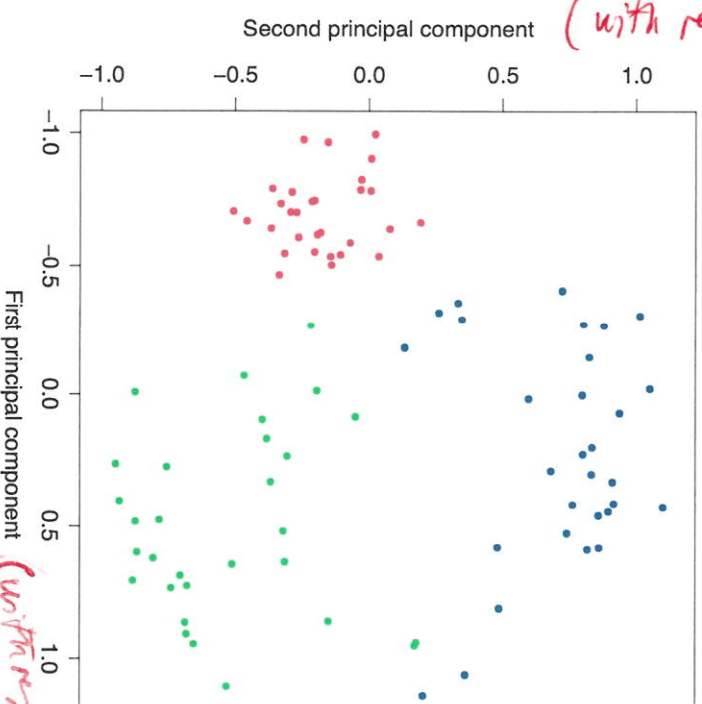
$p=3$



$\tilde{Y}_{N \times p}^T$

mapped to

(with respect to \tilde{V}_2)



$(Y_{PCA})_{N \times 2} = \text{coordinates in } (\tilde{V}_1, \tilde{V}_2)$

Illustration of PCA and SVD



Sample:

$$\{\vec{y}_i\}_{i=1}^{130}$$

$N = 130$ images

vectors:
 $\vec{y}_i \in \mathbb{R}^{16 \times 16}$

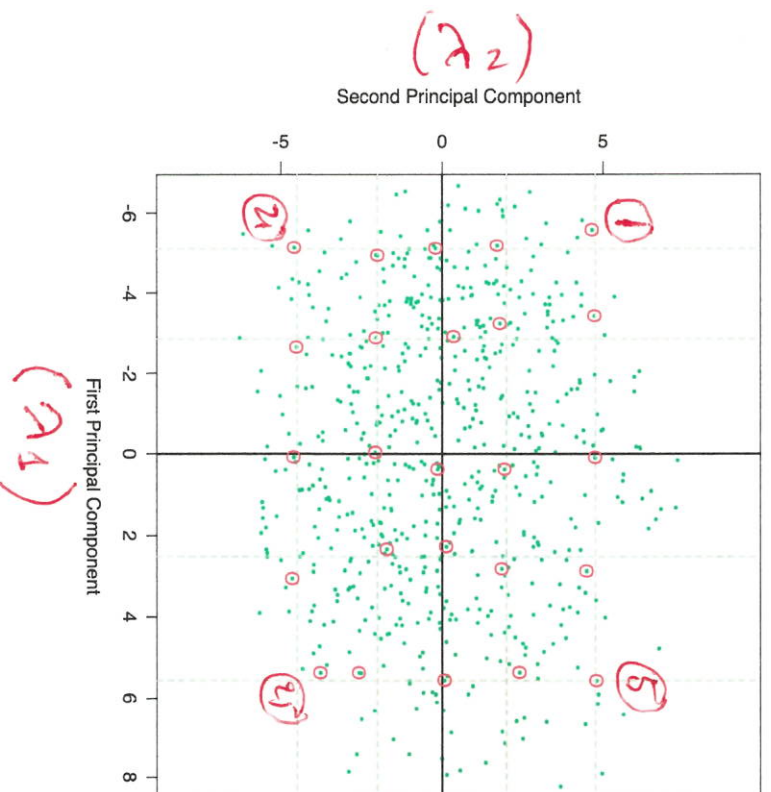
Each \vec{y}_i contains
the grayscale of
the images with
256 pixels.

Illustration of PCA and SVD

Vectors expressed in terms of the first two U_i

$$\begin{aligned} \hat{f}(\lambda) &= \bar{x} + \lambda_1 \vec{v}_1 + \lambda_2 \vec{v}_2 \\ &= \text{[image of digit 3]} + \lambda_1 \cdot \text{[image of digit 3]} + \lambda_2 \cdot \text{[image of digit 3]} \end{aligned}$$

Illustration of PCA and SVD



PC1
 (The bottom part of the 3 becomes more curved)

PC2
 The 3 becomes less thick

First extension:

PCA is about linear projections that minimise a quadratic.
How do we extend?

① "Nonlinearity": Kernel PCA

$$(\tilde{Y}^T \tilde{Y})_{ij} = (\tilde{\vec{y}}_i^T \tilde{\vec{y}}_j)$$

↓
PCA was about eigenvalues of $\tilde{Y}^T \tilde{Y}$ and $\tilde{Y} \tilde{Y}^T$

↓
which are "kernel" matrices.

If we assume a kernel function

$$k_{ij} = K(\tilde{\vec{y}}_i, \tilde{\vec{y}}_j)$$

↑
positive definite, etc. (see SVMs)

Then

Eigendecomposition of K (kernel matrix)

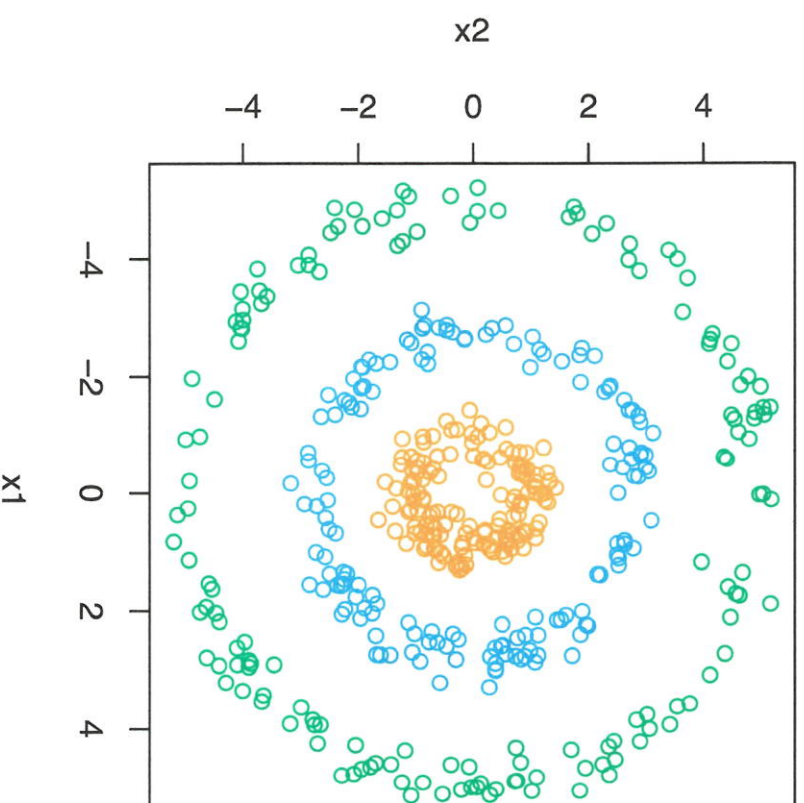
$$\text{e.g. } K(\vec{x}_i, \vec{x}_j) = e^{-\frac{\|\vec{x}_i - \vec{x}_j\|^2}{c}}$$

Radial Kernel

→
Example
overleaf

we can do spectral analysis of the kernel (which is nonlinear in the original coordinates)

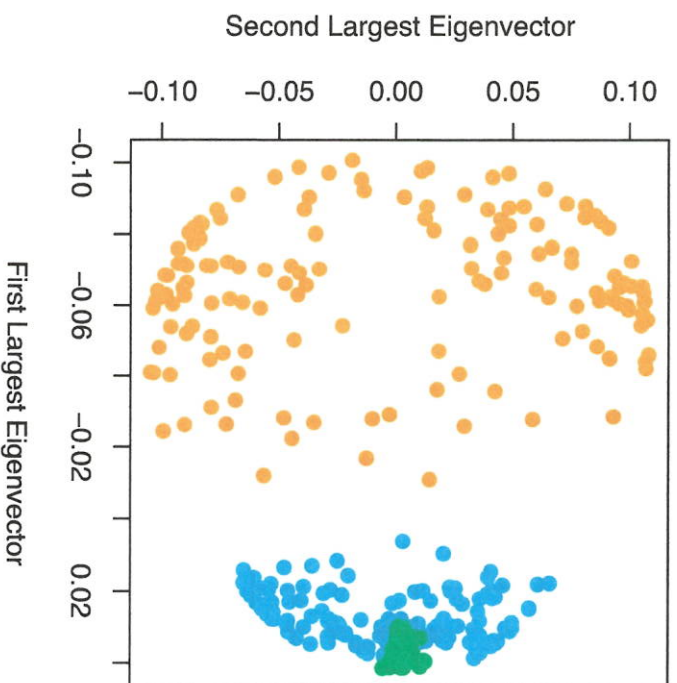
Representing these data in terms of standard
PCA does not
lead to a good mapping.



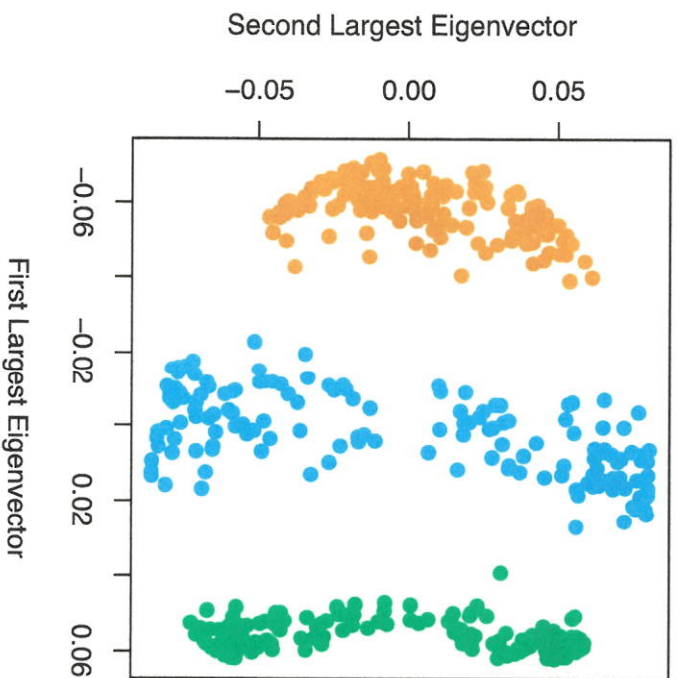
Kernel PCA

Increasing the hyperparameter c

Radial Kernel ($c=2$)



Radial Kernel ($c=10$)



Second extension: Interpretability through sparsity.

Sparse PCA

Our description from standard PCA was of the form:

$$(Y_{PCA})_{N \times K} = (U_K \Sigma_K)_{N \times K}$$

↓
coordinates in terms of

$$V_K = (\vec{v}_1 \dots \vec{v}_K)_{p \times K}$$

The basis V is non sparse in terms of our original descriptors

i.e. V is not a sparse matrix.

Original PCA:

This leads to diminished interpretability

This is an alternative description of PCA

$$\left\{ \begin{array}{l} \max_{\vec{v}} \vec{v}^T (\tilde{Y} \tilde{Y}^T) \vec{v} \\ \text{such that } \vec{v}^T \vec{v} = 1 \end{array} \right\} \quad \vec{v} \in \mathbb{R}^p$$

Sparser

by adding extra constraints.

$$\sum_{j=1}^p |v_j| \leq t$$

Similar to LASSO

$\{\vec{v}_{j,L}\}_{j=1}^p$
SVD LASSO basis

Find a basis that tries to minimize the quadratic error under sparsity constraints.

If we find a description in terms of a few vectors that are sparse, then the description is closer to identifying the important combinations of descriptors.

$$\hat{\vec{y}} = \hat{f}(\lambda) = \sum_{j=1}^k \lambda_j \underbrace{\vec{v}_j}_L$$

$k \ll$ from LASSO

Then the matrix V will be sparser and the \vec{v}_j will contain more zeros

$$\vec{v}_j = \begin{pmatrix} \boxed{1} \\ 0 \\ 0 \\ 0 \\ \boxed{1} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

coordinate for x_1

coordinate for x_j

Third extension: Interpretability through non-negativity
 Non-negative matrix factorisation (NMF)

Standard.
PCA/SVD

$$\hat{Y}^T = (U_K \Sigma_K) V_K^T$$

$N \times p \quad N \times K \quad K \times p$

↓ the same matrix

$$X = (U_K \Sigma_K) V_K^T$$

$N \times K \quad K \times p$

NMF =
 Generative to a factorisation: $X = W H^T$

$N \times K \quad K \times p$

where we require positivity

$$\begin{aligned} w_{ij} &\geq 0 & \forall i, j \\ h_{ij} &\geq 0 & \forall i, j \end{aligned}$$

It can be shown to be equivalent to

$$\max_{W, H} L(W, H) = \sum_{i=1}^n \sum_{j=1}^p \left[X_{ij} \log(WH)_{ij} - (WH)_{ij} \right]$$

$$X_{ij} \sim \text{Poi} \quad \text{with mean } (WH)_{ij}$$

What does ^{NMF} ~~it~~ give?

See over

for an example with images

Non negative matrix factorisation

Setup:

$$\{ \vec{y}_i \}_{i=1}^N \quad \vec{y}_i \in \mathbb{R}^p$$

$$N=2429$$

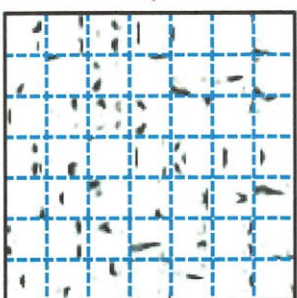
$$p=381=19 \times 19$$

$$k=49$$

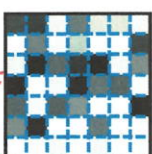
pixels of 49×49 we learn $H_{49 \times 381}$
in feature maps

Each of the squares is the 381-dimensional vector

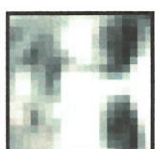
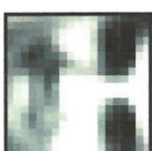
NMF



x



=

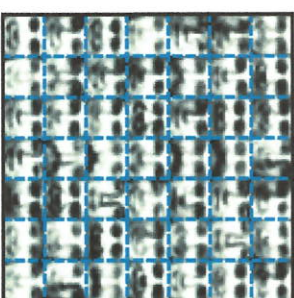


Original

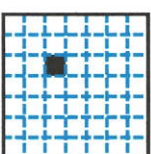
$$\vec{w}_i \text{ (column)} \\ W_{19 \times 49}$$

all positive

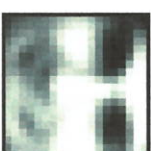
VQ = K-means



x



=

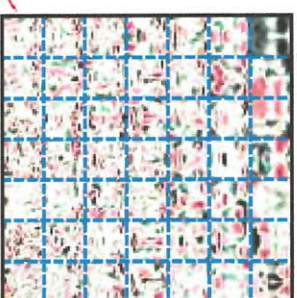


The 49 faces closest to 'original'

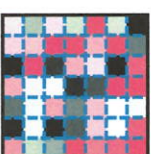
This whole set of squares is the $(V_k)_{p \times k}$ matrix

Each of the squares is a $\vec{v}_i \in \mathbb{R}^{381}$

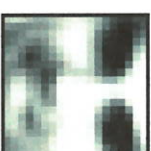
PCA



x



=



$$\vec{u}_i \text{ (column)} \\ U_{19 \times 49}$$

- Each \vec{y}_i is an image 19×19 pixels, which is vectorized into a 361-dimensional vector.
- Number of images in sample $N=2429$