

Principal Component Analysis

EE698V - Machine Learning for Signal Processing

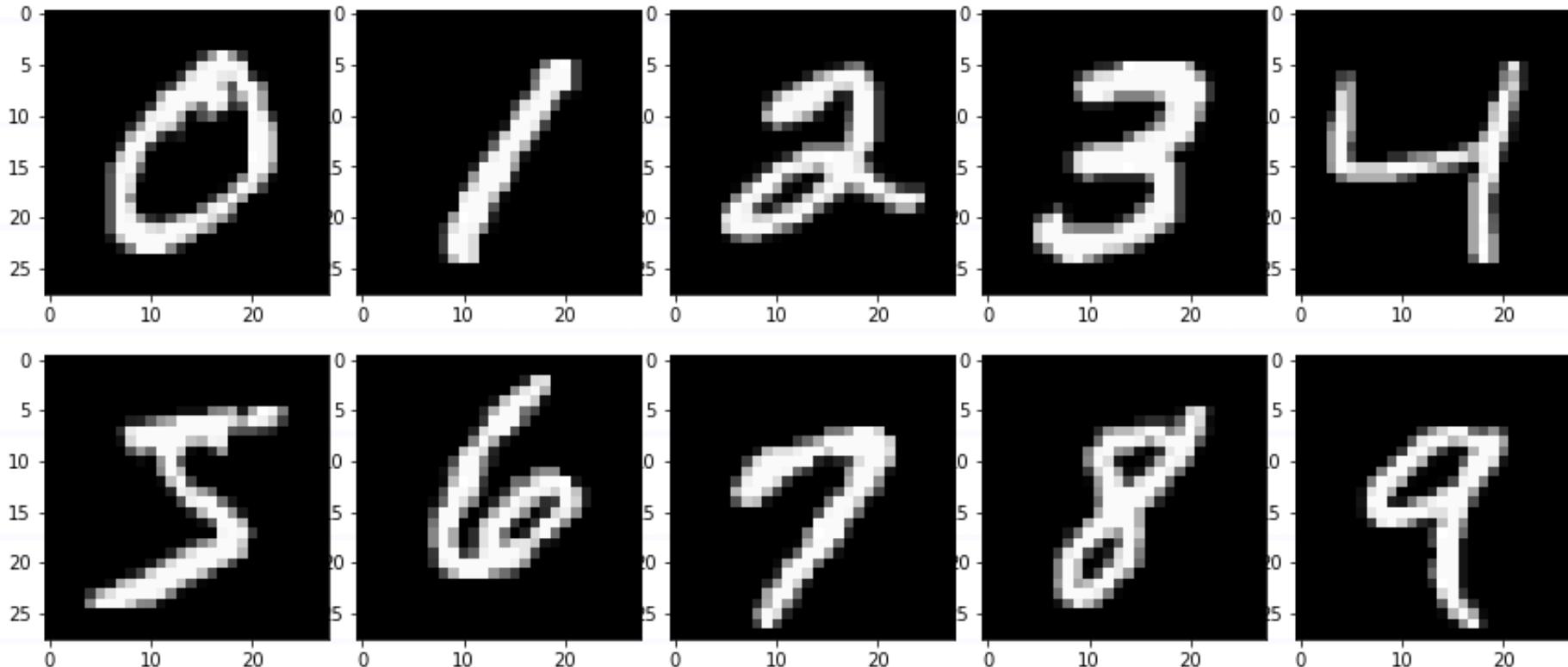
Vipul Arora



References

- PRML Section 12.1 **(highly recommended)**

Image Classification



Variations



Image Classification

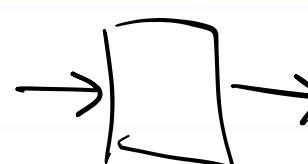
- What is one sample here?

3
3

X : 28×28 matrix

\wedge

$x[i, j]$

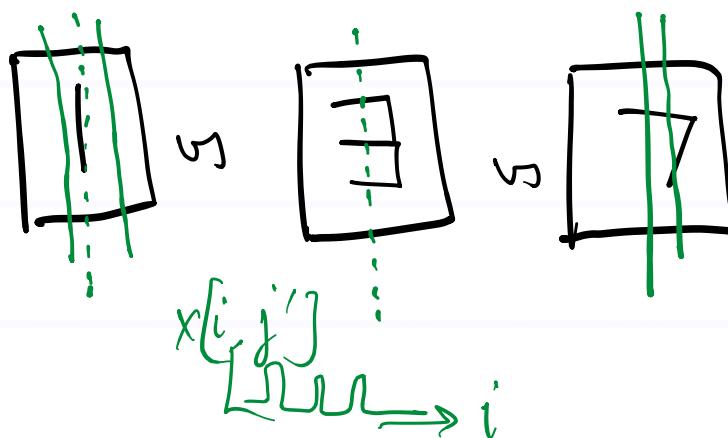


class $y \in \{0, \dots, 9\}$

$x[i, :]$

$x[:, i]$

✓ $x[:, :, :]$



$y = [\dots]$

$x[:, j]$

28 cols.



$y_t = 10 \times 1$

$y = [10 \times 28]$

$y_t = \max(y, \text{axis}=1)$

sum
avg

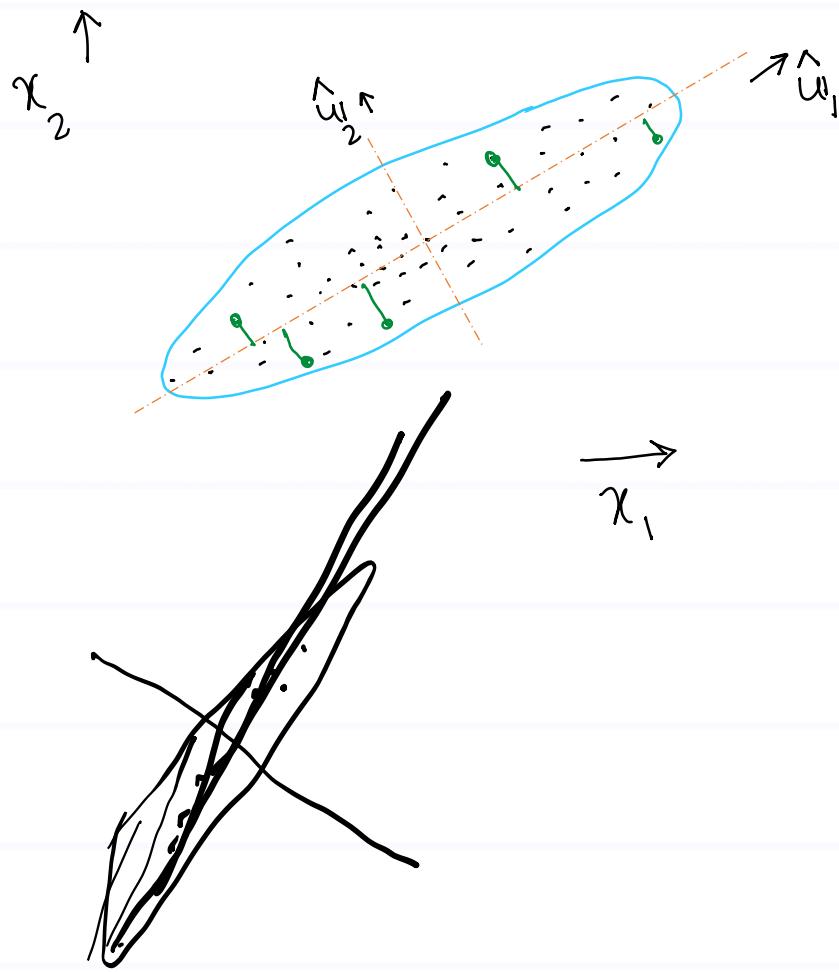
$y_t = 10 \times 1 \rightarrow$ to decide class

$c = \operatorname{arg\max}_c y[c]$

Classification

- Despite all variations in x , the information lies in a space of much lower dimensionality
- Instead of processing large sized samples, can we transform them to a lower dimensional space?
- E.g., we extracted height and width of flowers!

Projection to subspace

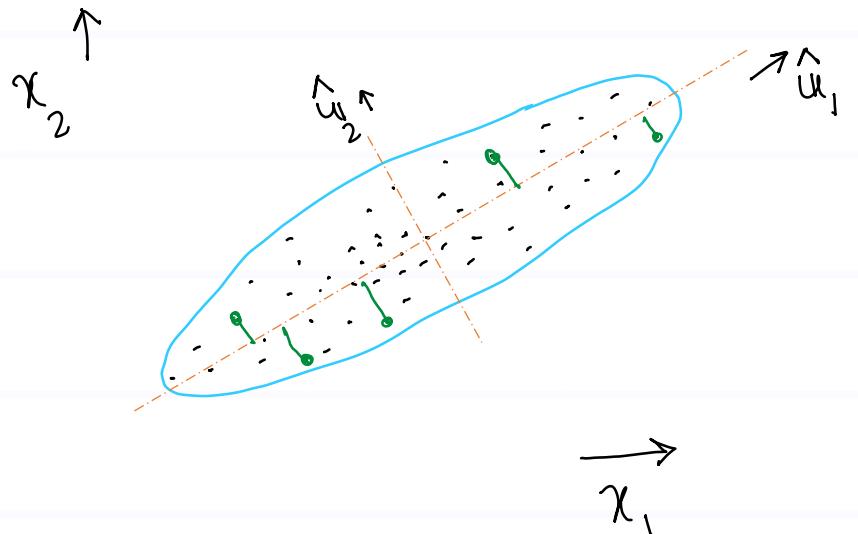


- The projection on u_1 axis contains more information than that on u_2 axis

- How?

$$x_i = (x_{1,i} \quad x_{2,i} \quad \dots \quad x_n)$$

Projection to subspace



- What is the projection on u_1 axis?

$$x^\top u_1$$

- What is the projection on u_2 axis?

$$x^\top u_2$$

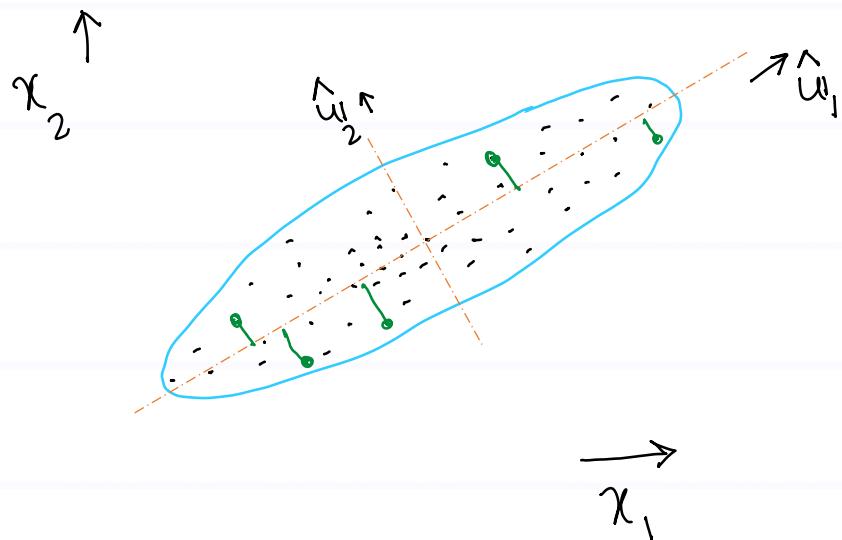
- So, we can write x as:

$$x = (x^\top u_1) u_1 + (x^\top u_2) u_2$$

more info

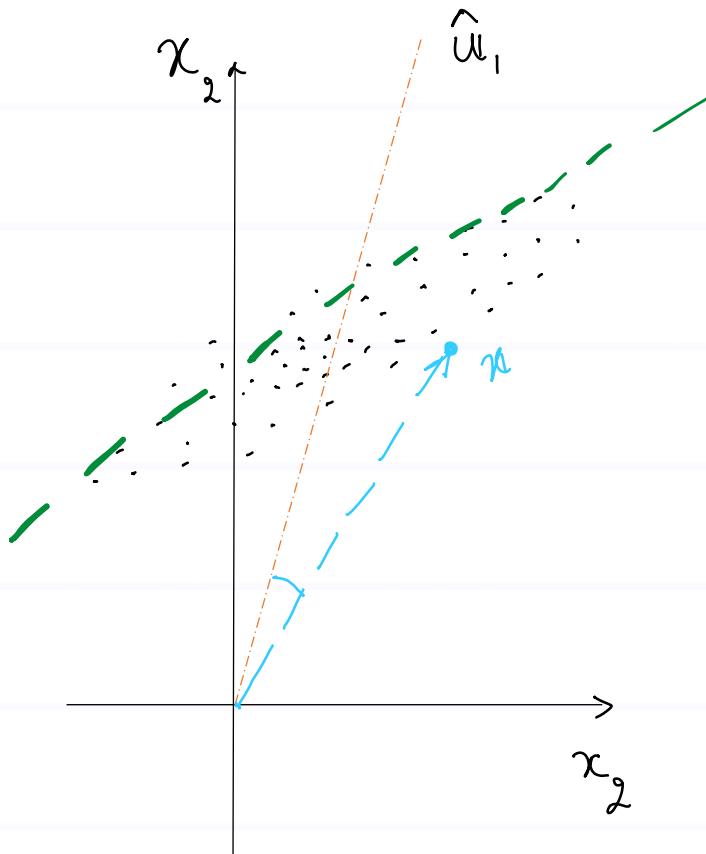
less info

Finding u_1



- u_1 is the axis along which there is maximum variation
- So, we can discard u_2

Finding u_1



- But hold on ...
- What is the direction of maximum variation here?

$$\underbrace{(x^T u_1)}_{\max} u_1 + \underbrace{(x^T u_2)}_{\min} u_2$$

- What is the error?

origin should be mean

Finding u_1

- Mean shifting: make the data zero-mean first

How to find optimal u axes?

- Minimize the information lost

$$x = \sum_{i=1}^D (x^T u_i) u_i ; \quad x \in \mathbb{R}^D$$

$$\tilde{x} = \sum_{i=1}^M (x^T u_i) u_i ; \quad M < D$$

How to find optimal u axes?

- Mean squared error ... again!

$$x - \tilde{x} = \sum_{i=M+1}^D (x^T u_i) u_i$$

$$E = \frac{1}{N} \sum_{n=1}^N \left(x_n - \tilde{x}_n \right)^2 = \frac{1}{N} \sum_{n=1}^N \left(\sum_{i=m+1}^D \left(x_n^T u_i \right) u_i \right)^2$$

$$= \frac{1}{2} \sum_{n=1}^N \left[\sum_{i=m+1}^D \underbrace{\left(\mathbf{x}_n^\top \mathbf{u}_i \right)^2 \mathbf{u}_i^\top \mathbf{u}_i}_{= 1 \text{ if } \mathbf{u}_i \text{'s}} + \sum_{i=m+1}^D \sum_{j=m+1, j \neq i}^D \underbrace{\left(\mathbf{x}_n^\top \mathbf{u}_i \right) \left(\mathbf{x}_n^\top \mathbf{u}_j \right) \mathbf{u}_i^\top \mathbf{u}_j}_{= 0 \text{ if } \mathbf{u}_i \text{'s}} \right]$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D \left(x_n^T w_i \right)^2$$

How to find optimal u axes?

$$E = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D \underbrace{\left(\mathbf{x}_n^\top \mathbf{u}_i \right)^\top \left(\mathbf{x}_n^\top \mathbf{u}_i \right)}_{\mathbf{u}_i^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{u}_i}$$

$$\begin{aligned} &= (\mathbf{u}_i^\top \mathbf{x}_n)^\top (\mathbf{u}_i^\top \mathbf{x}_n) \\ &= \mathbf{x}_n^\top \mathbf{u}_i \mathbf{u}_i^\top \mathbf{x}_n \end{aligned}$$

$$E = \sum_{i=M+1}^D \mathbf{u}_i^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)}_S \mathbf{u}_i$$

$S = \text{variance}$ [remember \mathbf{x}_n is zero mean]

Also, $\mathbf{u}_i^\top \mathbf{u}_i = 1$, so add a Lagrange multiplier term

How to find optimal u axes?

$$E = \sum_{i=M+1}^D \left[u_i^\top S u_i + \lambda_i (1 - u_i^\top u_i) \right]$$

Now minimize E w.r.t. u_i :

$$\frac{\partial E}{\partial u_i} = 2 S u_i - \lambda_i 2 u_i = 0$$

$$\Rightarrow S u_i = \lambda_i u_i \quad (\lambda_i \text{ is eigen value of } S)$$

left multiply with u_i^\top

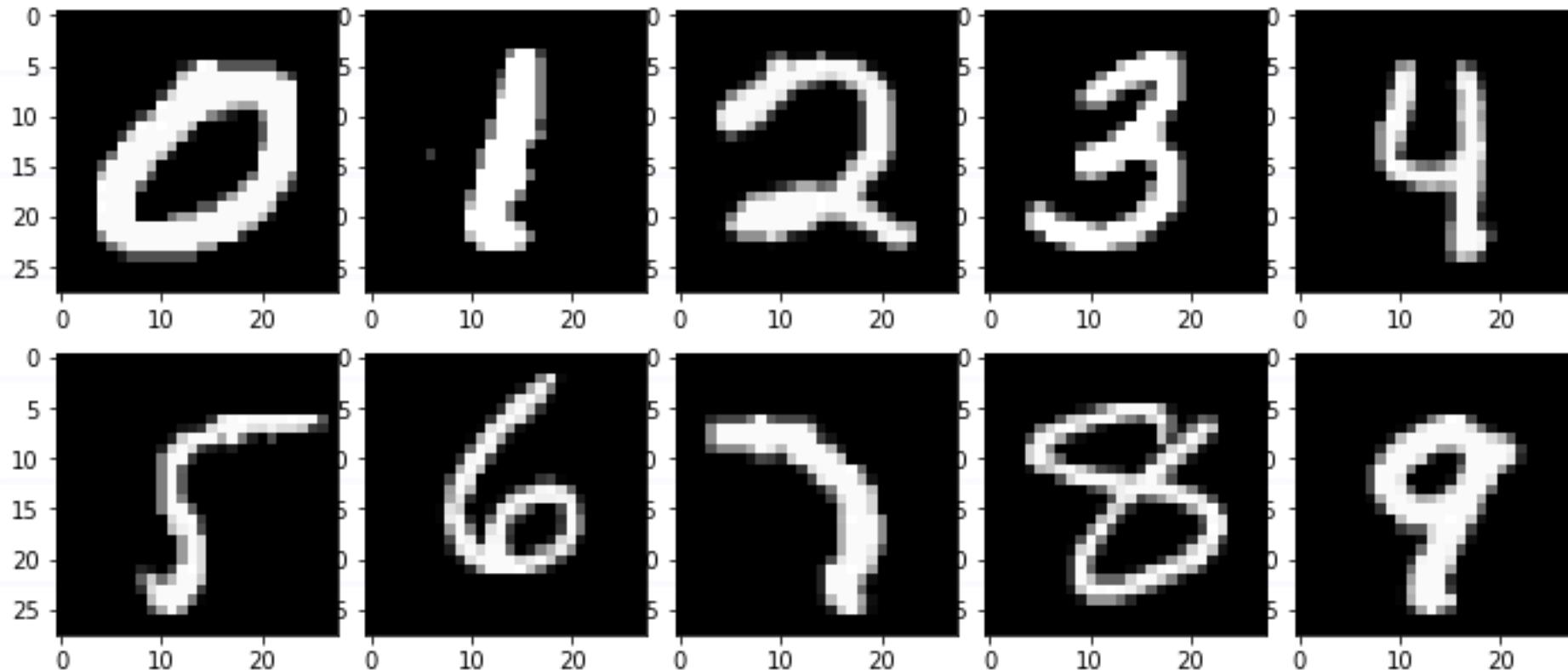
$$u_i^\top S u_i = \lambda_i$$

$$\Rightarrow E = \sum_{i=M+1}^D \lambda_i$$

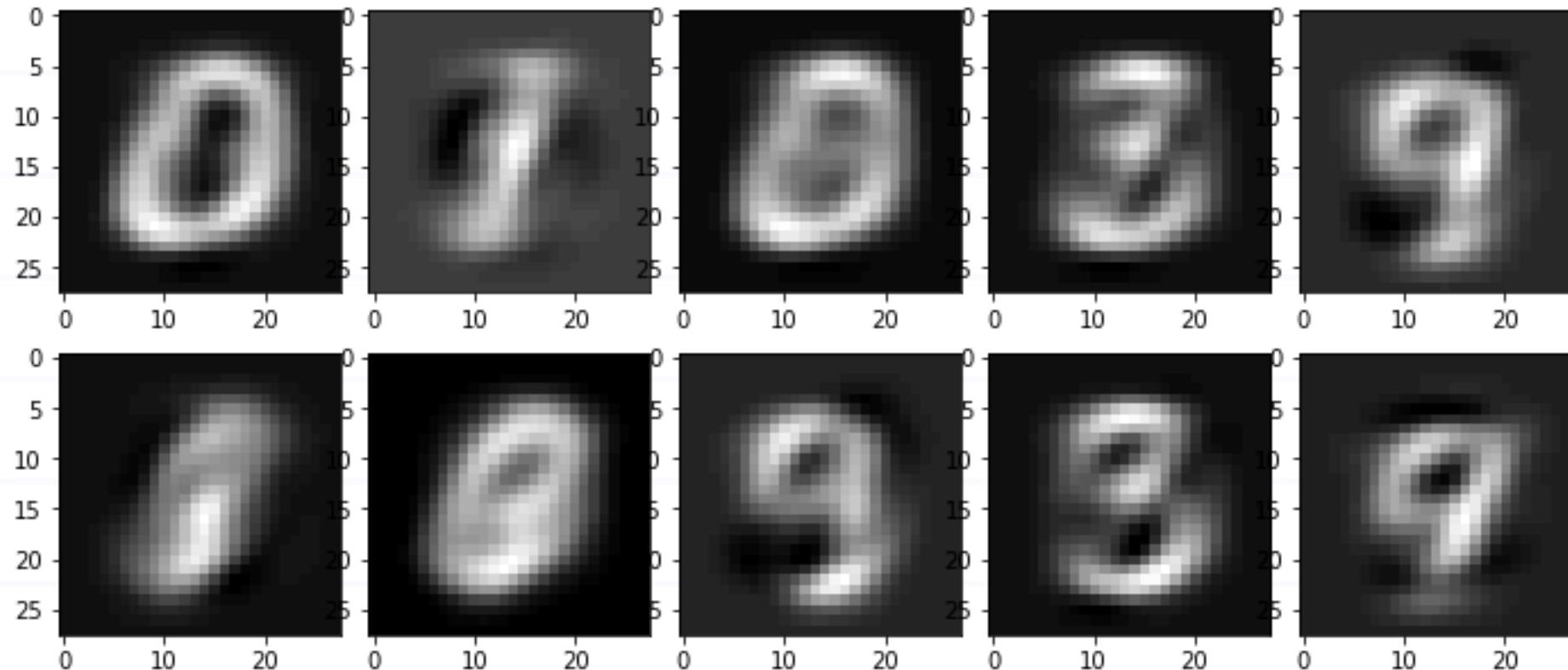
Principal Component Analysis

- Given data samples $s_n \in \mathbb{R}^D$
- Normalize: $x_n = s_n - \mathbb{E}[s]; \quad \mathbb{E}[s] = \frac{1}{N} \sum_{n=1}^N s_n$ $x_n \in \mathbb{R}^D$
- Obtain variance matrix $S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T$
- Eigen value decomposition of S to get $\lambda_i, u_i; i = 1, \dots, D$ with λ_i in decreasing order $u_i \in \mathbb{R}^D$
- Choose first M eigen vectors as the principal axes
- $\tilde{x}_n = \sum_{i=1}^M (x_n^T u_i) u_i$ $\tilde{x}_n \in \mathbb{R}^D$

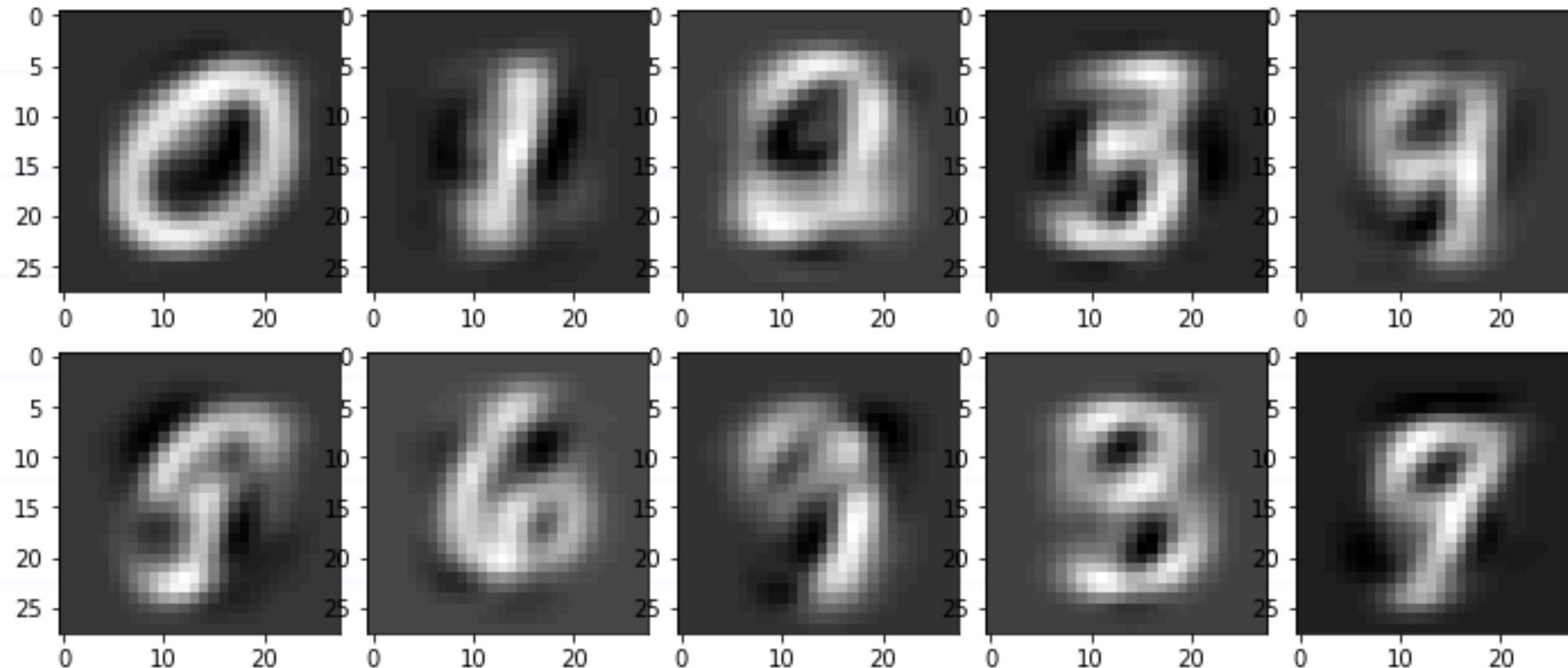
Original



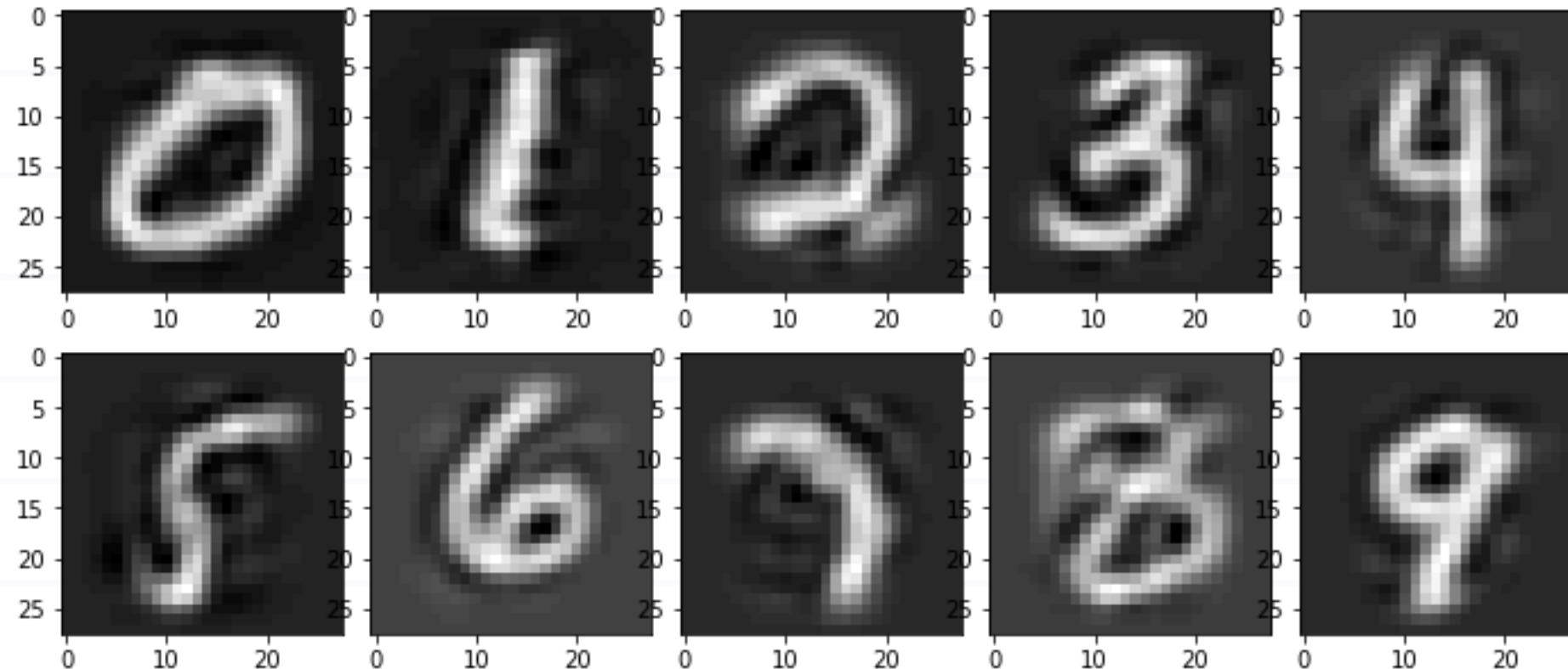
PCA with M=3



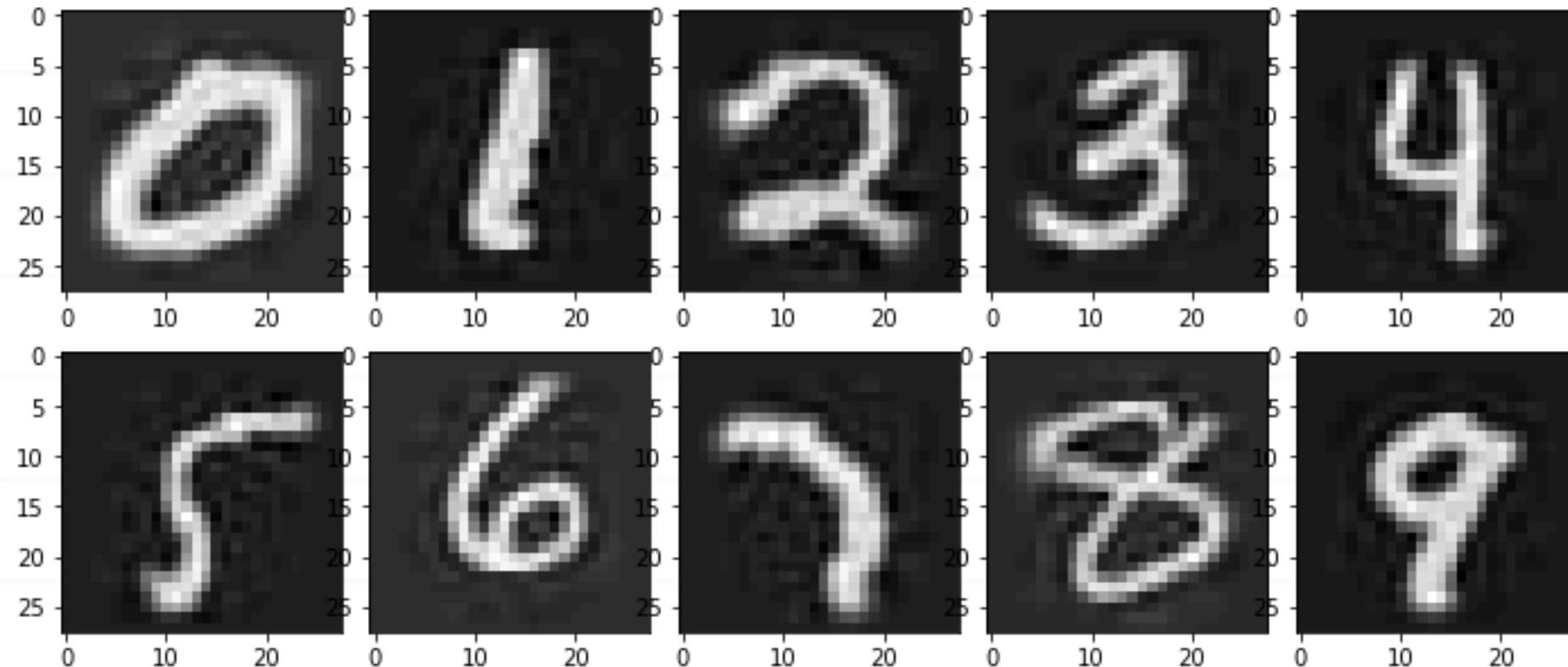
PCA with M=10



PCA with M=50



PCA with M=150



PCA as transformation

$$\tilde{x}_n = \sum_{i=1}^M \left(u_i^\top x_n \right) u_i$$

$$y_n = \begin{bmatrix} u_1^\top x_n \\ u_2^\top x_n \\ \vdots \\ u_M^\top x_n \end{bmatrix} = [u_1 \ u_2 \ \dots \ u_M]^\top x_n = U^\top x_n$$

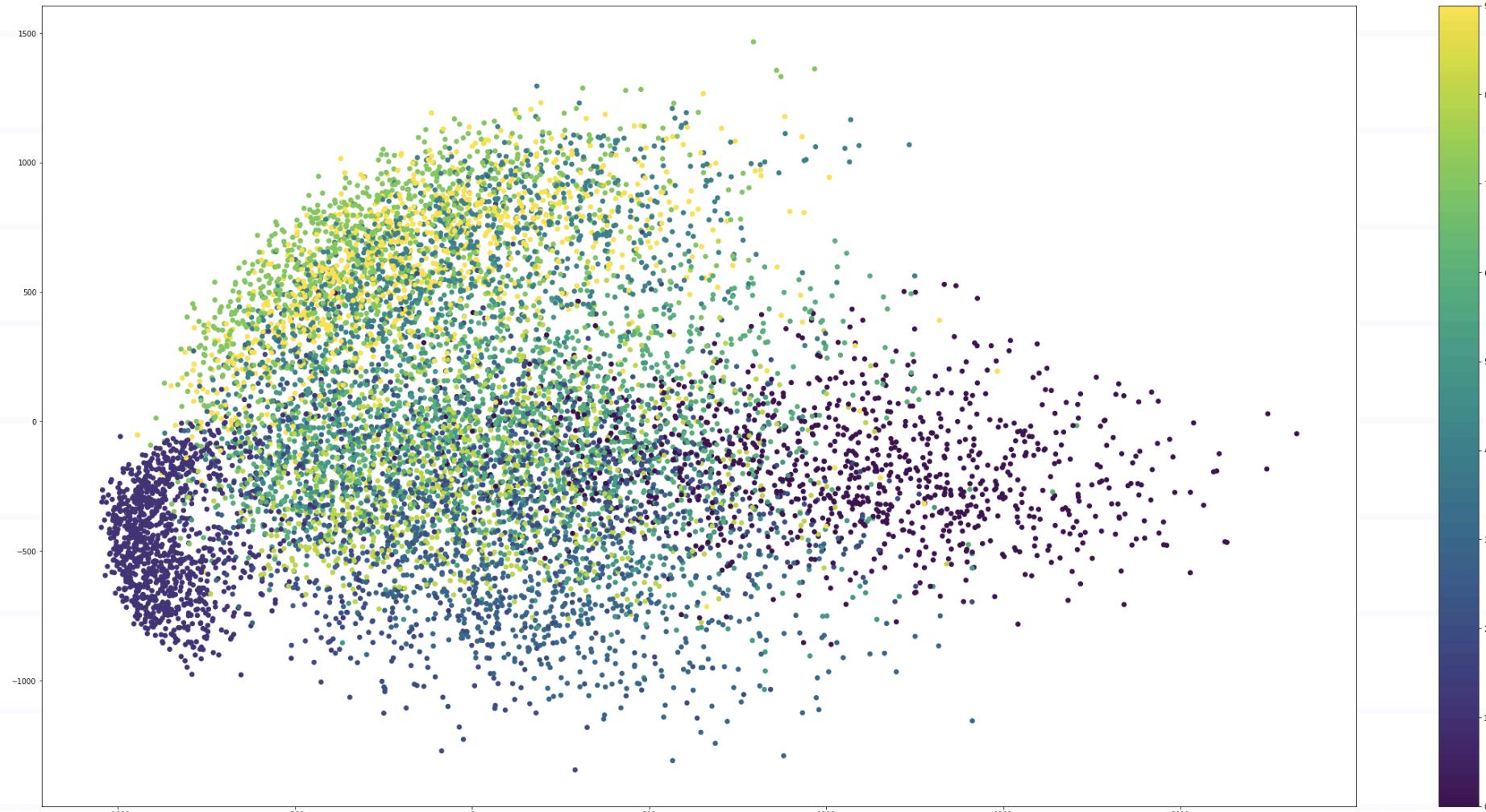
$x_n \in \mathbb{R}^D$

$y_n \in \mathbb{R}^M$

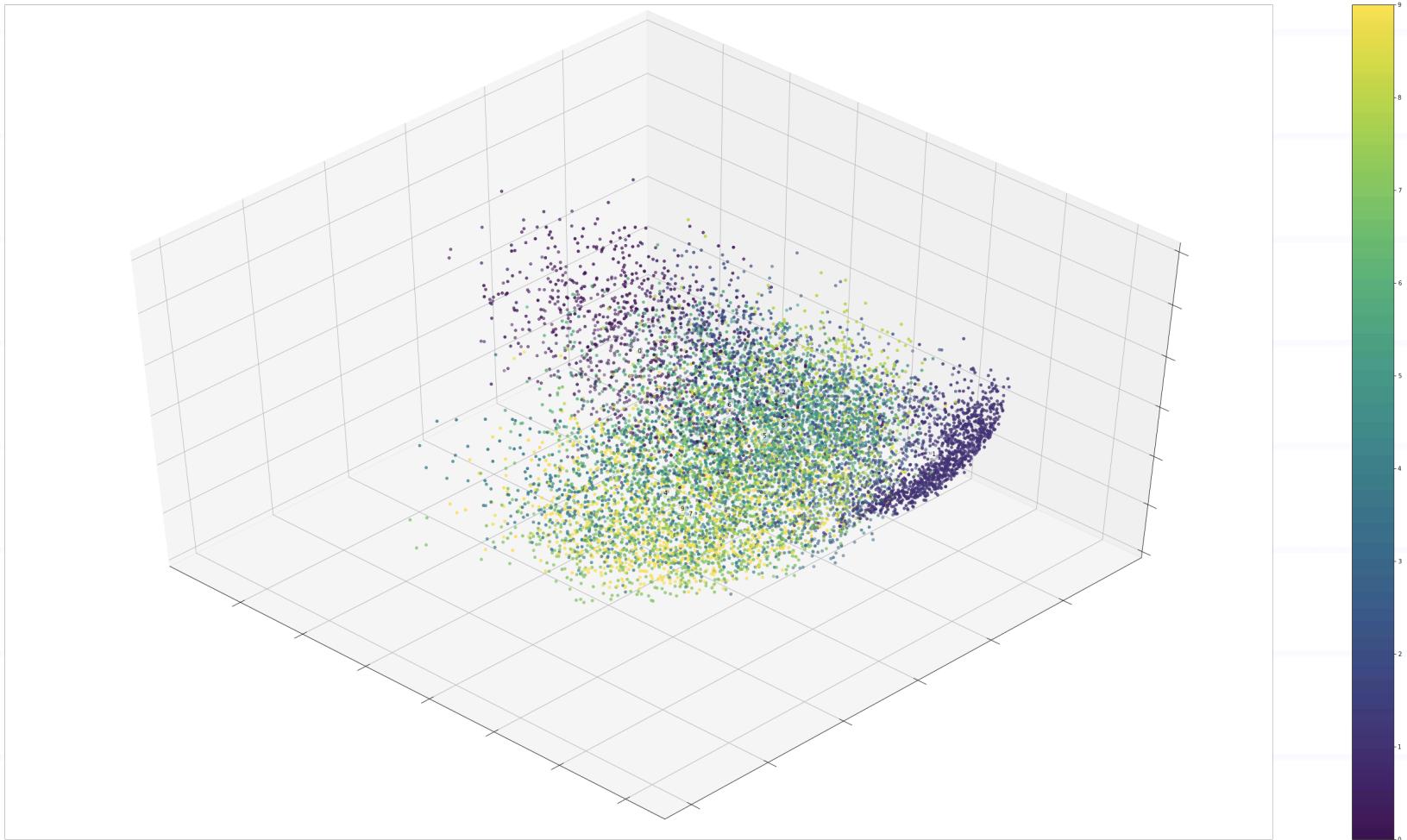
U is $D \times M$

- Prove that y_n is zero-mean
- What is the variance of y_n ?

PCA on MNIST images



PCA on MNIST images



PCA as transformation

- Is PCA invariant to

- Transversal shift
 - in images no
 - in vectors yes

- Rotation

- in images no
- in vectors yes

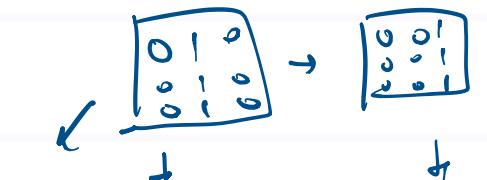
$$x_n' = x_n + a$$

$$\forall n \in \{1, \dots, N\}$$

$$y_n' \stackrel{?}{=} y_n$$

yes

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$



$$x_n' = A x_n$$

$$y_n' \stackrel{?}{=} y_n \quad \text{yes}$$

$$\begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

PCA as normalization

$$\frac{1}{N} \sum_n y_n y_n^T = \frac{1}{N} \sum_n U^T x_n x_n^T U = \underbrace{U^T S U}_{U L = L}$$

Recall $SU = UL$, u_i were eigen vectors of S , and $L = \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \dots \\ 0 & & \lambda_M \end{bmatrix}$

and $U^T U = I$ \because orthonormal

$$\therefore \frac{1}{N} \sum_n y_n y_n^T = \underbrace{U}_L \underbrace{L}_{V} V$$

Can we make it unit variance?

$$y_n = L^{-1} V U x_n$$

PCA as Matrix Factorization

$$y_n = U^T x_n \quad \forall n$$

$$U = [u_1 \ u_2 \dots \ u_m]$$

In matrix form: $[y_1 \ y_2 \ \dots \ y_N] = U^T [x_1 \ x_2 \ \dots \ x_N]$

$$\underbrace{Y}_{M \times N} = U^T \underbrace{X}_{784 \times N}$$
$$M \times 784$$

- This is a way of matrix factorization
- Here, Y, U, X can have any real values
- But for images, negative values do not make sense.
- Can we constrain U to have only non-negative values? Then Y will also be non-negative

Discrete vs Continuous Latent Variables

- In discrete
 - We map each x to a discrete cluster
- In continuous
 - We map each x to a weighted sum of some vectors (i.e., \mathbf{u}_i)

