

From Principal Component Analysis to Deep Learning with De-Noising Variational Auto-Encoders

Dr. Daniel Guterding
daniel.guterding@gmail.com

February 5th, 2020

Lecture topic

Data sets with a large number of dimensions (variables) are a challenge for machine learning algorithms with respect to computational effort and memory usage.

The lecture will answer three closely related questions in this context:

1. How to reduce the number of dimensions in a sensible way?
2. How to distinguish relevant from irrelevant dimensions?
3. How to suppress noise in high-dimensional data?

Discussed algorithms

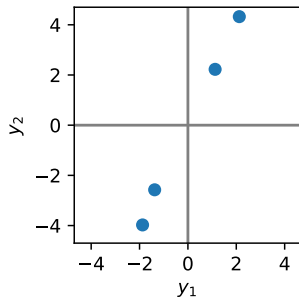
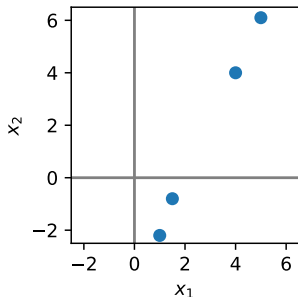
1. Principal Component Analysis (PCA)
2. De-Noising Variational Auto-Encoder (DVAE)

Principal Component Analysis - Foundations

- ▶ take **N data points** with **D dimensions**
- ▶ write single data point as column vector with \underline{x}_i ($D \times 1$)
- ▶ matrix of all data points is $\underline{X} = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)$ with ($D \times N$)
- ▶ center of all data points is $\underline{\mu} = \frac{1}{N} \sum_{i=1}^N \underline{x}_i$
- ▶ **centered data points** are $\underline{Y} = (\underline{y}_1, \underline{y}_2, \dots, \underline{y}_N)$ with $\underline{y}_i = \underline{x}_i - \underline{\mu}$

Example

$$\underline{x}_1 = \begin{pmatrix} 4.0 \\ 4.0 \end{pmatrix}, \underline{x}_2 = \begin{pmatrix} 5.0 \\ 6.1 \end{pmatrix}, \underline{x}_3 = \begin{pmatrix} 1.5 \\ -0.8 \end{pmatrix}, \underline{x}_4 = \begin{pmatrix} 1.0 \\ -2.2 \end{pmatrix}, \underline{\mu} = \begin{pmatrix} 2.875 \\ 1.775 \end{pmatrix}$$



Principal Component Analysis - Foundations

- ▶ calculate **covariance matrix** $\underline{\underline{\Sigma}} = \underline{\underline{Y}}\underline{\underline{Y}}^T$, real and symm. with $(D \times D)$
- ▶ solve eigenvalue equation $\underline{\underline{\Sigma}} = \underline{\underline{V}}\underline{\underline{\Lambda}}\underline{\underline{V}}^T$, real eigvecs and eigvals
- ▶ **diagonal variance matrix** $\underline{\underline{\Lambda}} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_D)$
- ▶ eigenvectors $\underline{\underline{V}} = (\underline{v}_1, \underline{v}_2, \dots, \underline{v}_D)$ are **orthogonal**
- ▶ \underline{v}_i is i -th principal component with variance λ_i
- ▶ **principal components can be sorted by variance**

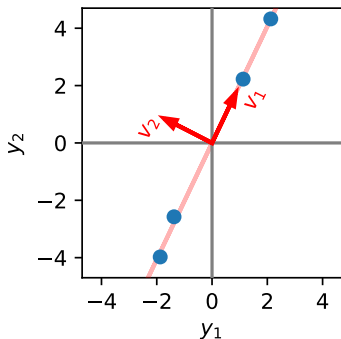
Example

$$\underline{\underline{Y}} = \begin{pmatrix} 1.125 & 2.125 & -1.375 & -1.875 \\ 2.225 & 4.325 & -2.575 & -3.975 \end{pmatrix}$$

$$\underline{\underline{\Sigma}} = \underline{\underline{Y}}\underline{\underline{Y}}^T = \begin{pmatrix} 11.1875 & 22.6875 \\ 22.6875 & 46.0875 \end{pmatrix}$$

$$\underline{\underline{V}} = (\underline{v}_1, \underline{v}_2) = \begin{pmatrix} 0.442 & -0.897 \\ 0.897 & 0.442 \end{pmatrix}$$

$$\underline{\underline{\Lambda}} = \text{diag}(\lambda_1, \lambda_2) = \text{diag}(57.3, 0.02)$$

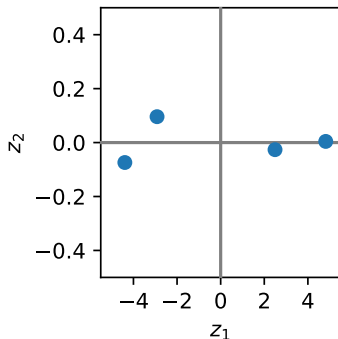
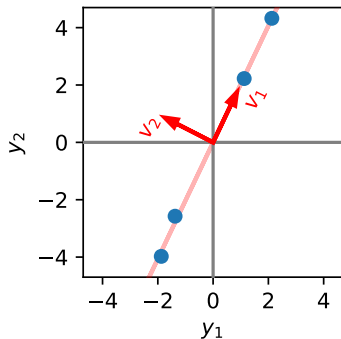


Principal Component Analysis - Foundations

- ▶ principal components form **new coordinate system**
- ▶ matrix of linear transformation is unitary $\underline{\underline{V}}^T = \underline{\underline{V}}^{-1}$
- ▶ new data matrix $\underline{\underline{Z}} = \underline{\underline{V}}^T \underline{\underline{Y}}$ with $(D \times N) = (D \times D) \cdot (D \times N)$
- ▶ variance is **diagonal**: $\underline{\underline{Z}} \underline{\underline{Z}}^T = \underline{\underline{V}}^T \underline{\underline{Y}} \underline{\underline{Y}}^T \underline{\underline{V}} = \underline{\underline{V}}^T \underline{\underline{\Sigma}} \underline{\underline{V}} = \underline{\underline{V}}^T \underline{\underline{V}} \underline{\underline{\Lambda}} \underline{\underline{V}}^T \underline{\underline{V}} = \underline{\underline{\Lambda}}$

Example

$$\underline{\underline{Z}} = \begin{pmatrix} 2.493 & 4.819 & -2.918 & -4.394 \\ -0.026 & 0.004 & 0.096 & -0.074 \end{pmatrix}$$



Principal Component Analysis - Recap

- ▶ principal component analysis finds **new coordinate system**
- ▶ **transformation** between old and new coordinates is **linear**
- ▶ **new basis vectors are called principal components**
- ▶ covariance matrix is diagonal in these coordinates
- ▶ principal components can be ordered by their contribution to the data set's total variance
- ▶ for data that almost lie on a line, the **variance can be concentrated within a few principal components**
- ▶ this allows for **dimensionality reduction by discarding principal components** with low associated variance

Principal Component Analysis - Dimensionality Reduction

- ▶ discarding principal component \underline{v}_i discards variance λ_i
- ▶ easily done by erasing the i -th column in $\underline{\underline{V}}$
- ▶ usually \underline{v}_i **with smallest λ_i discarded first**
- ▶ new transformation $\underline{\underline{\tilde{V}}} = (\underline{v}_1, \underline{v}_2, \dots, \underline{v}_{D-1})$ with $(D \times (D - 1))$
- ▶ **reduced data** $\underline{\underline{\tilde{Z}}} = \underline{\underline{\tilde{V}}}^T \underline{\underline{Y}}$ with $((D - 1) \times N)$
- ▶ discarding \underline{v}_i erases i -th row in $\underline{\underline{Z}}$

Example

$$\underline{\underline{V}} = (\underline{v}_1, \underline{v}_2) = \begin{pmatrix} 0.442 & -0.897 \\ 0.897 & 0.442 \end{pmatrix}, \underline{\underline{\Lambda}} = \text{diag}(\lambda_1, \lambda_2) = \text{diag}(57.3, 0.02)$$

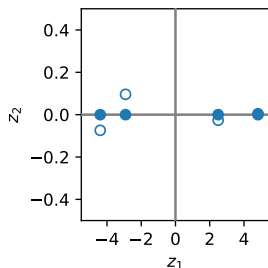
$$\underline{\underline{Z}} = \underline{\underline{V}}^T \underline{\underline{Y}} = \begin{pmatrix} 2.493 & 4.819 & -2.918 & -4.394 \\ -0.026 & 0.004 & 0.096 & -0.074 \end{pmatrix}$$

$$\underline{\underline{\tilde{V}}} = (\underline{v}_1) = \begin{pmatrix} 0.442 \\ 0.897 \end{pmatrix}$$

$$\underline{\underline{\tilde{Z}}} = \underline{\underline{\tilde{V}}}^T \underline{\underline{Y}} = (2.493 \quad 4.819 \quad -2.918 \quad -4.394)$$

Principal Component Analysis - Dimensionality Reduction

- ▶ **transformation back to original coordinates**
with $\underline{\tilde{Y}} = \underline{\tilde{V}} \underline{\tilde{Z}} = \underline{\tilde{V}} \underline{\tilde{V}}^T \underline{Y}$
- ▶ although $\underline{V} \underline{V}^T = \underline{1}$, we have $\underline{\tilde{V}} \underline{\tilde{V}}^T \neq \underline{1}$
- ▶ **information loss** due to reduced dimension
- ▶ back-transformed data lie in subspace of lower dimension



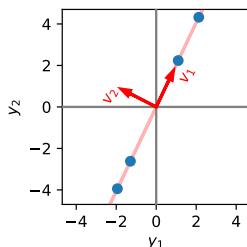
Example

$$\underline{Y} = \begin{pmatrix} 1.125 & 2.125 & -1.375 & -1.875 \\ 2.225 & 4.325 & -2.575 & -3.975 \end{pmatrix}$$

$$\underline{\tilde{V}} = (\underline{v}_1) = \begin{pmatrix} 0.442 \\ 0.897 \end{pmatrix}$$

$$\underline{\tilde{Z}} = \underline{\tilde{V}}^T \underline{Y} = \begin{pmatrix} 2.493 & 4.819 & -2.918 & -4.394 \end{pmatrix}$$

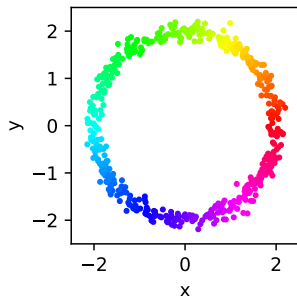
$$\underline{\tilde{Y}} = \underline{\tilde{V}} \underline{\tilde{Z}} = \underline{\tilde{V}} \underline{\tilde{V}}^T \underline{Y} = \begin{pmatrix} 1.101 & 2.129 & -1.289 & -1.941 \\ 2.237 & 4.323 & -2.617 & -3.942 \end{pmatrix}$$



Principal Component Analysis - Problems

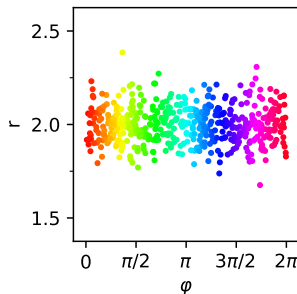
Non-Linear Data

- ▶ principal component analysis only works well for **linear data**
- ▶ introduction of non-linear variables would further increase dimensionality
- ▶ kernel methods are solution (kernel PCA)

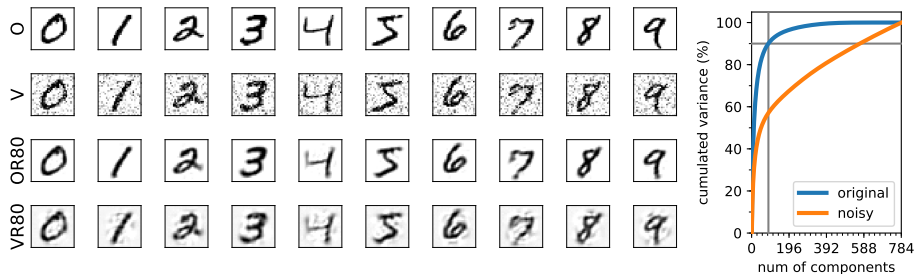


Algorithm

- ▶ diagonalization of covariance matrix does not scale well with dimensionality D
- ▶ implementations usually use **singular value decomposition**
- ▶ same result, but relation to variance harder to understand



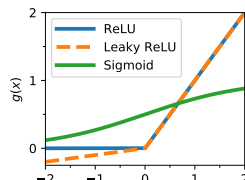
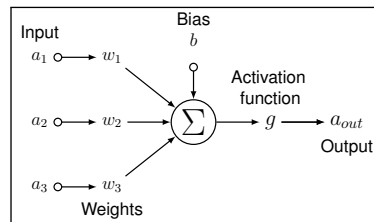
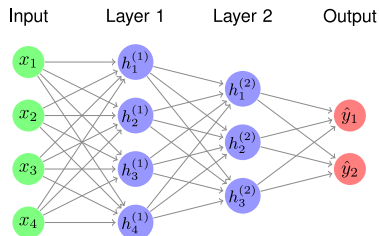
Principal Component Analysis - De-Noising of Images



- ▶ MNIST Numbers contains 70000 hand-written numbers from 0 to 9
- ▶ $28 \times 28 = 784$ pixel per image, values from 0 (white) to 255 (black)
- ▶ **transform** every **image to column vector** with dimension (784×1)
- ▶ 80 (about 10%) of principal components contain 90% of variance
- ▶ we introduce random noise into a derived data set
- ▶ contained variance converges slower than for original data
- ▶ noise leads to small reconstruction errors and lowered contrast

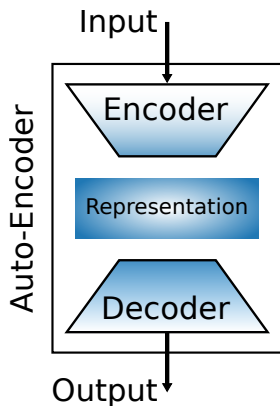
Artificial Neural Networks - Foundations

- ▶ inspired by the structure of the brain
- ▶ layers of **simple units** with **complex connections**
- ▶ allows for **non-linear transformations**
- ▶ complexity of network can be controlled
- ▶ unit performs **weighted sum** and applies **activation function**
- ▶ output of unit is $a_{\text{out}} = g(b + \sum_i a_i w_i)$
- ▶ activation function ReLU
 $g(x) = \max(0, x)$ popular
- ▶ **weights w_i must be learned**
- ▶ high computational effort, training best done on GPUs



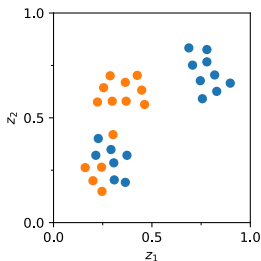
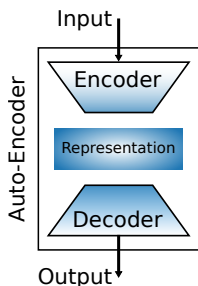
Auto-Encoder - Foundations

- ▶ encoder compresses high-dim. input data into low-dim. representation vector \underline{z}
- ▶ decoder decompresses \underline{z} into high-dim. output data
- ▶ compression is **lossy**, depends on structure of encoder/decoder and dimension of \underline{z}
- ▶ principal component analysis is primitive Auto-Encoder; $\underline{\tilde{V}}^T$ as encoder, $\underline{\tilde{V}}$ as decoder
- ▶ **real Auto-Encoder uses neural network as encoder/decoder**
- ▶ allows for better compression and reconstruction because of more complex structure



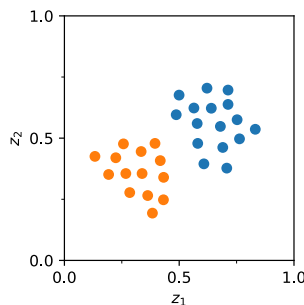
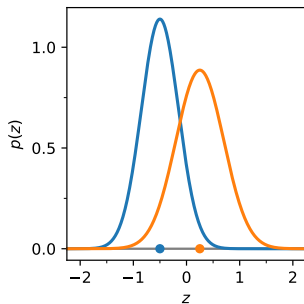
Auto-Encoder - Foundations

- ▶ **training** of encoder/decoder **minimizes quadratic deviation** between input and output
- ▶ representation vector \underline{z} lies in so-called **latent space**
- ▶ possible to generate and decode new \underline{z} , so-called **generative model**
- ▶ **transformation is point-wise**, points that are close in latent space not necessarily related
- ▶ no guarantee that similar \underline{z} lead to similar outputs
- ▶ **insufficient reconstruction between training data**
- ▶ Variational Auto-Encoder solves these problems

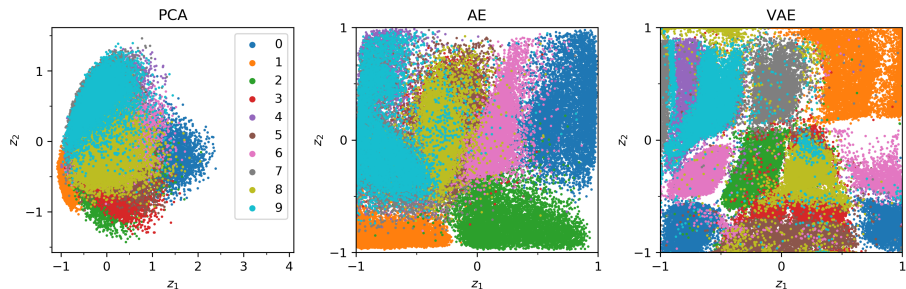


Variational Auto-Encoder - Foundations

- ▶ **learns probability distribution** $p(\underline{z})$
- ▶ input data are mapped to normal distribution with $\underline{\mu}$ and $\ln \underline{\Sigma}$, no correlation
- ▶ $\underline{z} = \underline{\mu} + \underline{\sigma}\varepsilon$, $\underline{\sigma} = \exp\left(\frac{1}{2} \ln \underline{\Sigma}\right)$, $\varepsilon \sim \mathcal{N}(0, 1)$
- ▶ **Kullback-Leibler divergence** measures deviation of distribution from $\mathcal{N}(0, 1)$
- ▶ $D_{\text{KL}} = \frac{1}{2} \sum_i (-1 - \ln(\sigma_i^2) + \mu_i^2 + \sigma_i^2) \geq 0$
- ▶ **KLD is added to quadratic error**, regularizes $\underline{\mu}$ and $\ln \underline{\Sigma}$
- ▶ KLD centers \underline{z} of training data around $\underline{0}$
- ▶ **similar input data, similar \underline{z}**
- ▶ algorithm seems ad-hoc, but mathematically well-established

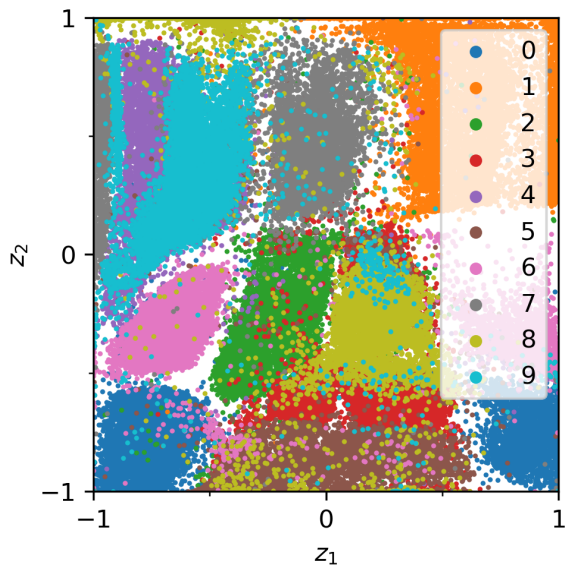


Algorithms Comparison - Latent Space



- ▶ use MNIST data set
- ▶ plots use classification info, model does not
- ▶ in all cases two-dimensional representation is learned
- ▶ numbers 4, 5, 6 particularly problematic
- ▶ AE and VAE also not perfect
- ▶ grouping of numbers is done best by VAE

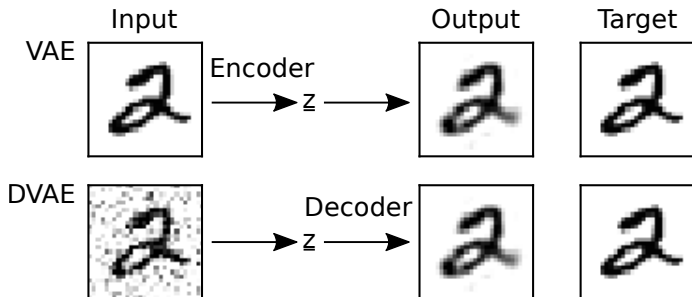
Variational Auto-Encoder - Generative Learning



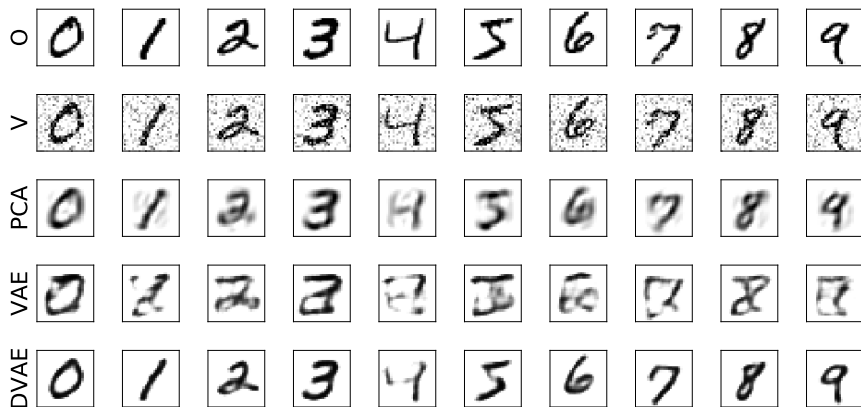
9 4 8 8 7 7 1 1 1 1
9 4 9 9 7 7 1 1 1 1
9 4 9 9 7 7 1 1 1 1
9 4 9 9 7 7 1 1 1 1
9 4 9 9 7 7 1 1 1 1
9 4 9 4 7 7 1 1 1 1
9 9 9 4 7 7 1 1 1 1
9 9 9 4 2 3 3 8 6 6
9 9 4 2 2 3 3 6 6 6
9 6 6 2 2 8 8 6 6 6
9 6 6 2 2 8 8 8 6 6
6 6 6 2 2 8 8 8 6 6
6 6 2 2 2 8 8 8 6 6
0 0 2 2 3 3 3 3 0 0
0 0 0 3 3 3 3 3 0 0
0 0 0 5 5 5 5 5 0 0
0 0 0 5 5 5 5 5 0 0

De-Noising Variational Auto-Encoder - Foundations

- ▶ train VAE with **artificial noisy input data**
- ▶ calculate reconstruction error w.r.t. **original data**
- ▶ encoder learns de-noising
- ▶ latent space in DVAE similar to VAE
- ▶ **decoder unchanged w.r.t. VAE**
- ▶ rigorous mathematical justification exists



De-Noising VAE - De-Noising of Images



- ▶ 32 dimensions in latent space
- ▶ **PCA and VAE trained with original data**
- ▶ DVAE shows best image reconstruction
- ▶ **VAE very sensitive to noise**

Summary

Principal Component Analysis

- ▶ linear transformation diagonalizes variance matrix
- ▶ principal comp. can be sorted by variance
- ▶ dimensionality reduction by discarding principal components with low variance

De-Noising Variational Auto-Encoder

- ▶ combination of neural networks as encoder/decoder
- ▶ learns representation of data as low-dim. probability distribution
- ▶ trained with artificial noisy data and known clean data as target
- ▶ generative model



Literature

- ▶ Hastie, Tibshirani, Friedman: The Elements of Statistical Learning
- ▶ Geron: Hands-On Machine Learning with Scikit-Learn and TensorFlow
- ▶ Goodfellow, Bengio, Courville: Deep Learning
- ▶ Foster: Generative Deep Learning
- ▶ Kingma, Welling: Auto-Encoding Variational Bayes, arXiv:1312.6114
- ▶ Im *et al.*: Denoising Criterion for Variational Auto-Encoding Framework, arXiv:1511.06406
- ▶ Doersch: Tutorial on Variational Autoencoders, arXiv:1606.05908
- ▶ Rolínek, Zietlow, Martius: Variational Autoencoders pursue PCA directions (by accident), arXiv:1812.06775

Training complexities of ML algorithms

- ▶ all listed complexities are provable upper bounds
- ▶ more efficient implementations may and do exist
- ▶ N : number of training samples
- ▶ D : number of features
- ▶ M : number of trees
- ▶ **Linear regression**: $\mathcal{O}(ND^2 + D^3)$, because $\underline{\beta} = (\underline{X}\underline{X}^T)^{-1}\underline{X}\underline{Y}$
- ▶ **Support Vector Machine**: $\mathcal{O}(N^2D + N^3)$
- ▶ **Decision tree**: $\mathcal{O}(N^2D)$
- ▶ **Random forest**: $\mathcal{O}(N^2DM)$
- ▶ **Artificial Neural Network**: no general proof available

MNIST Principal Components

