

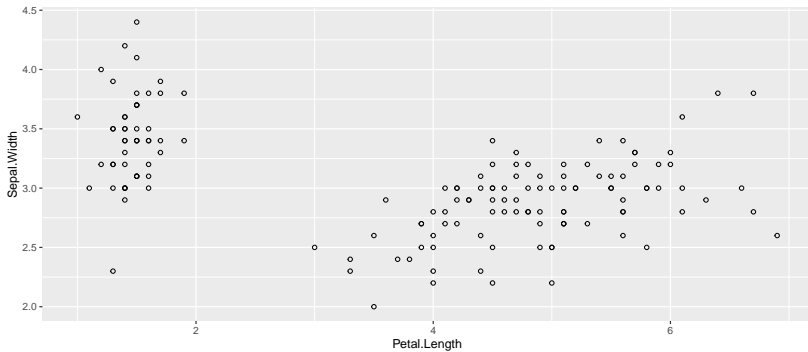
Principal Components Analysis

Toby Dylan Hocking

Background/motivation: dimensionality reduction

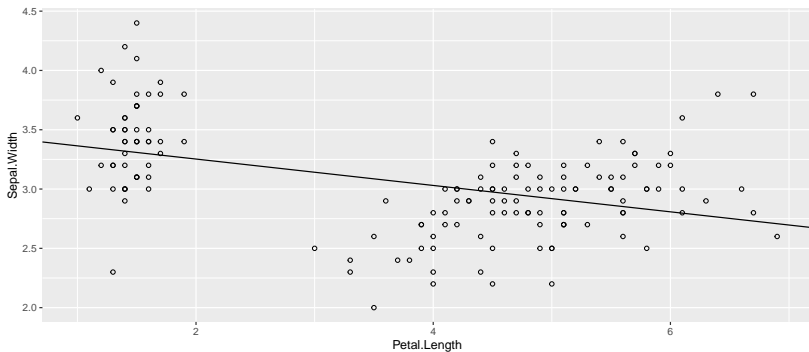
- ▶ High dimensional data are difficult to visualize.
- ▶ For example each observation/example in the zip data is of dimension $16 \times 16 = 256$ pixels.
- ▶ We would like to map each observation into a lower-dimensional space for visualization / understanding patterns in the data.

Example 1: 2d iris data



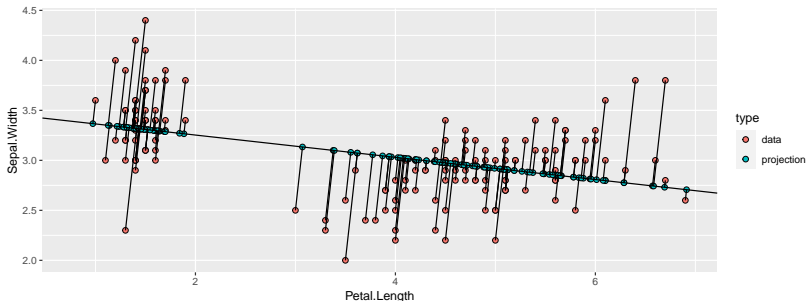
Project 2d data onto 1d subspace (line)

Why this line?

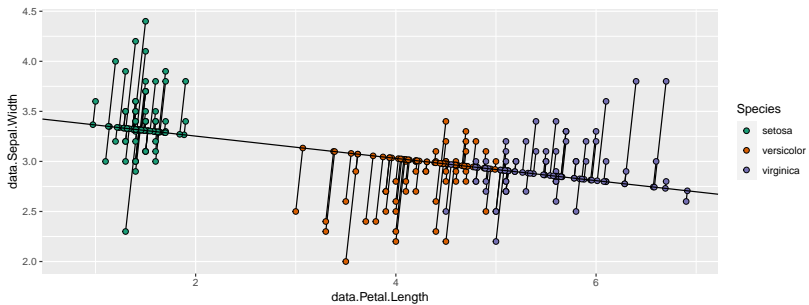


Principal Components Projection

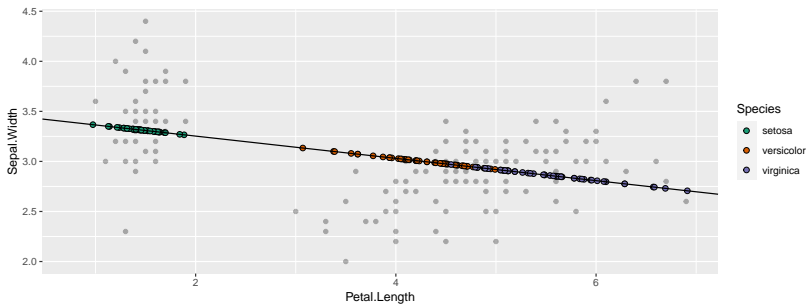
The first principal component is the line which minimizes the reconstruction error, squared distance between projection and data.



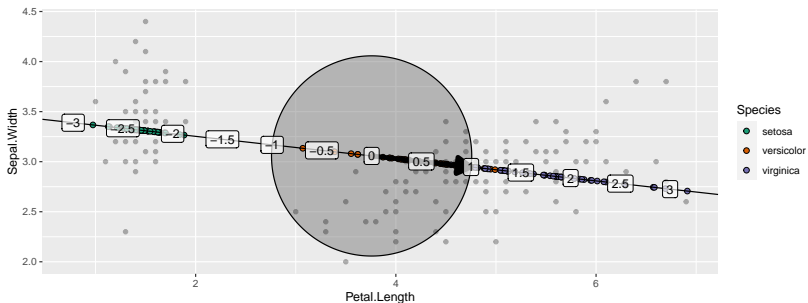
Map label onto projection



Map label onto projection

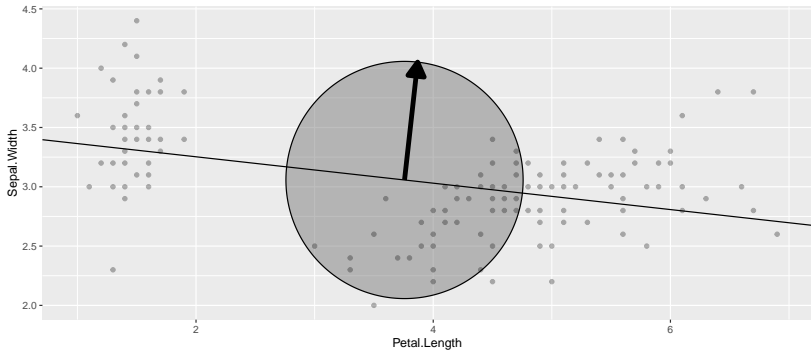


Principal component 1, amount along projection

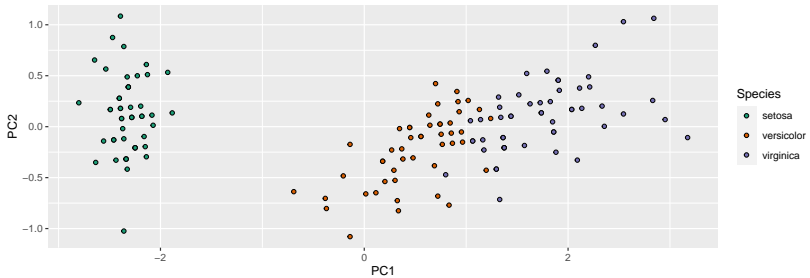


- ▶ 0 represents mean of data.
- ▶ $0 \rightarrow 1$ represents an orthogonal unit vector.

Principal component 2



Re-plot using PC units



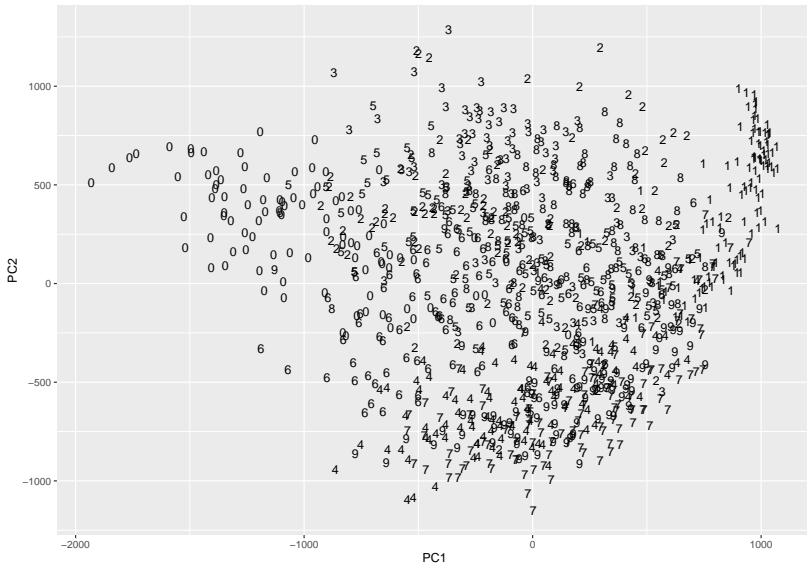
Mathematical representation

Each of the n inputs $x_i \in \mathbb{R}^p$ where p is the input dimension, $p = 2$ for iris in previous slides, $p = 784$ for images in next slides.

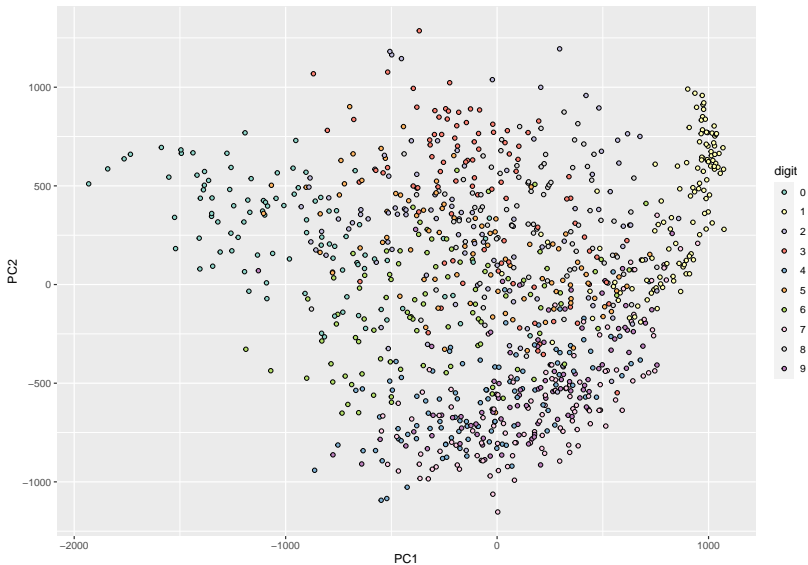
$$\min \sum_{i=1}^n \|x_i - \mu - V_q \lambda_i\|^2.$$

- ▶ $\mu \in \mathbb{R}^p$ is mean vector.
- ▶ $V_q \in \mathbb{R}^{p \times q}$ is an orthogonal matrix (each column is an orthogonal unit vector).
- ▶ $\lambda_i \in \mathbb{R}^q$ is a vector of principal components (contribution of each unit vector).

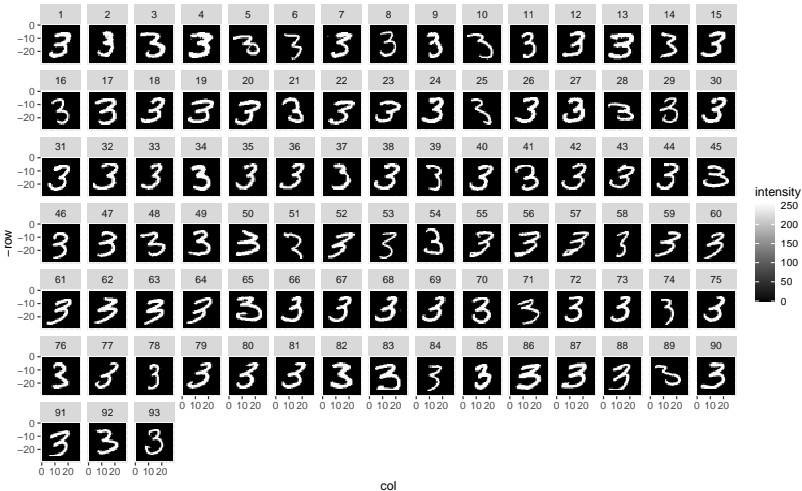
Same analysis with MNIST digit data



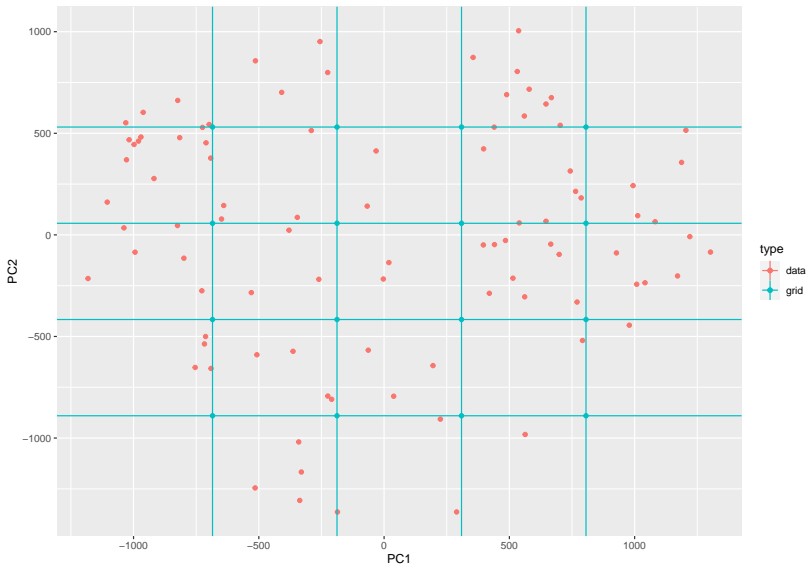
Alternate visualization



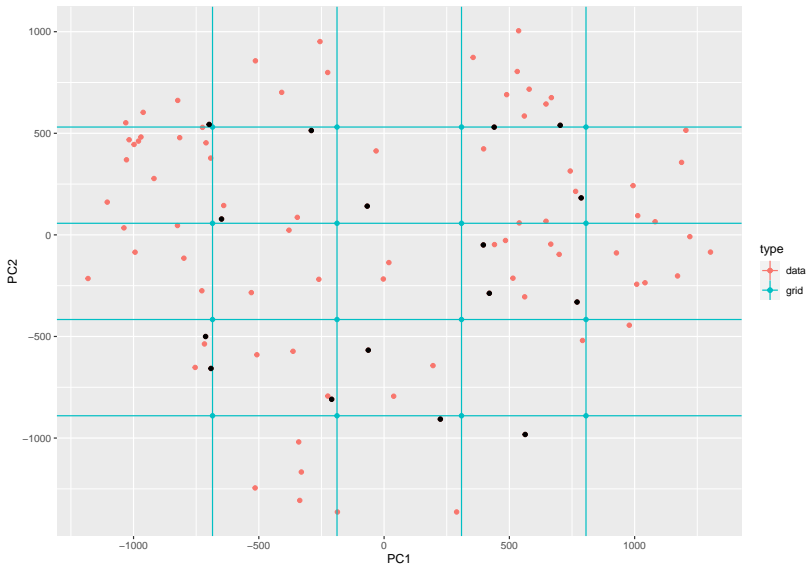
Another PCA on just one digit class



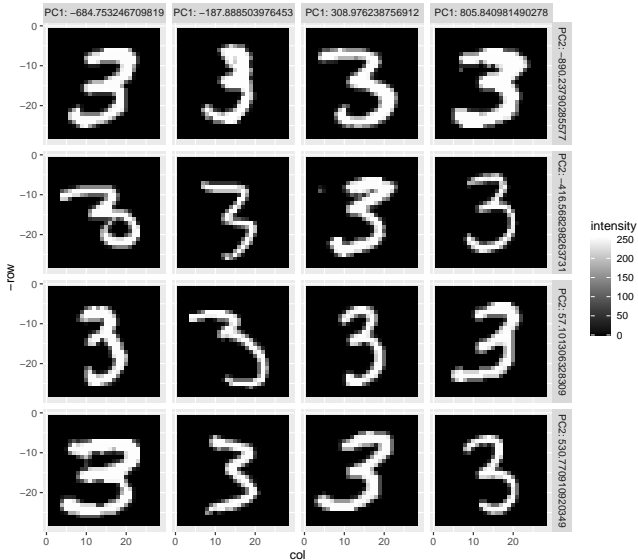
Mapping onto first two PCs



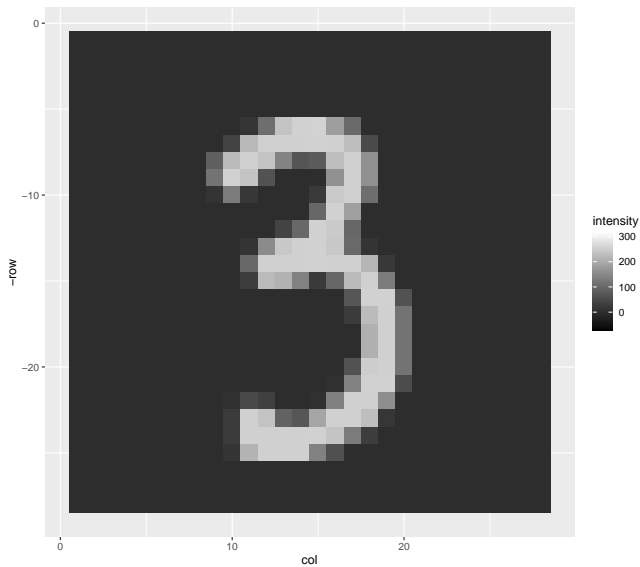
Highlight closest data point to each grid point



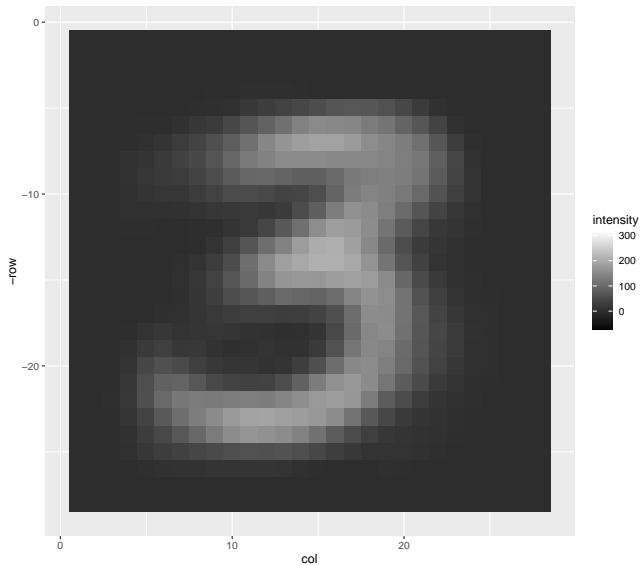
Digits highlighted



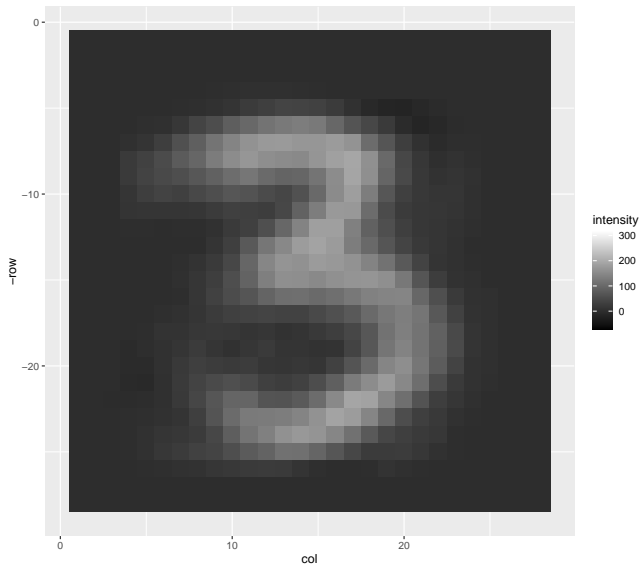
One digit



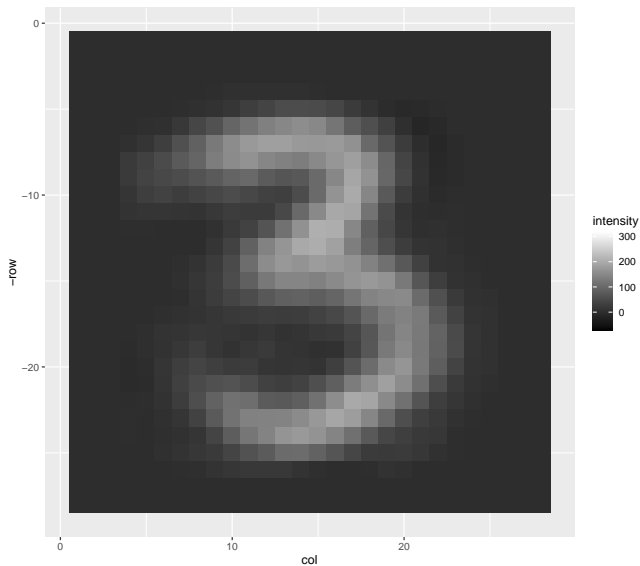
Reconstruction with no components (mean)



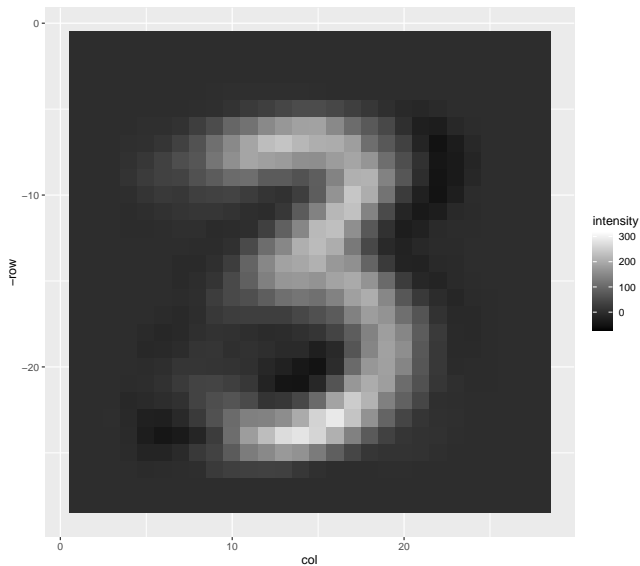
Reconstruction with one PC



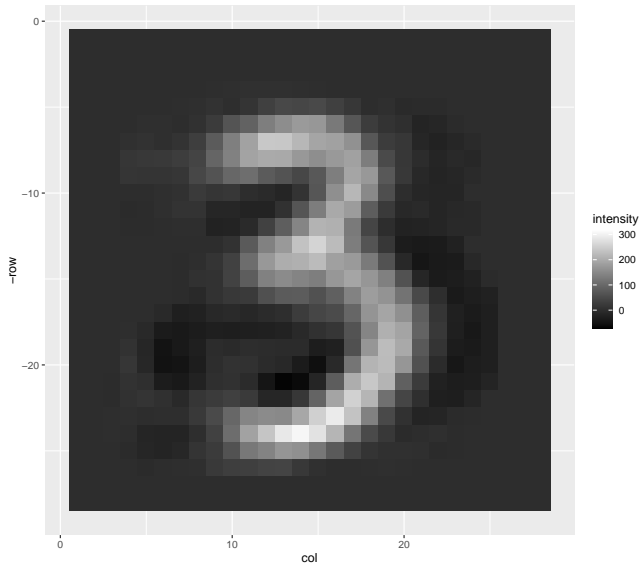
Reconstruction with 2 PCs



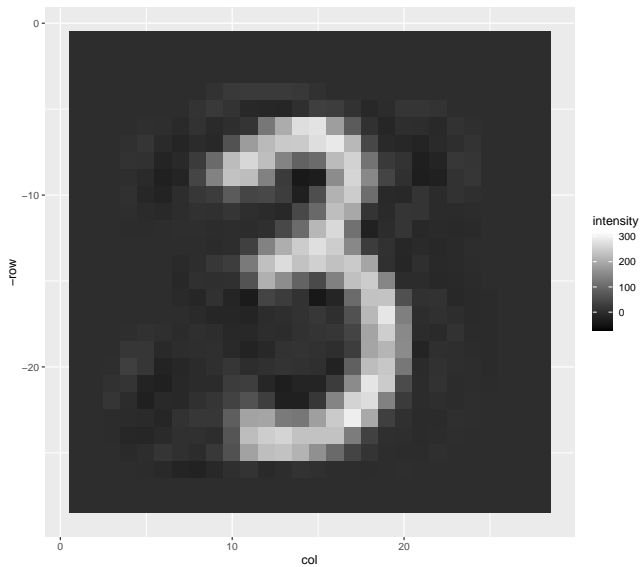
Reconstruction with 5 PCs



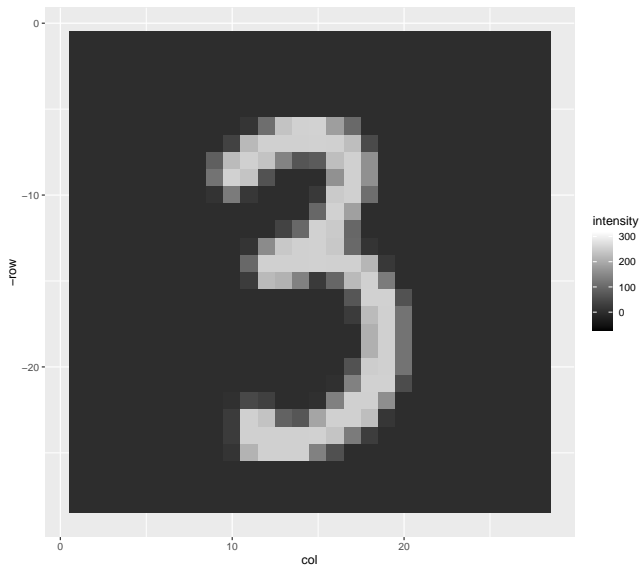
Reconstruction with 10 PCs



Reconstruction with 50 PCs



Reconstruction with all PCs



How to compute PCA?

SVD = Singular Value Decomposition (many algorithms available to compute).

$$X = UDV^T$$

- ▶ $X \in \mathbb{R}^{n \times p}$ data matrix.
- ▶ $U \in \mathbb{R}^{n \times p}$ orthogonal matrix.
- ▶ $D \in \mathbb{R}^{p \times p}$ diagonal matrix.
- ▶ $V \in \mathbb{R}^{p \times p}$ orthogonal matrix.
- ▶ The V_q we want for PCA is the first q columns of V .
- ▶ The columns of UD are the principal components, λ_i values.

Possible exam questions

► TODO