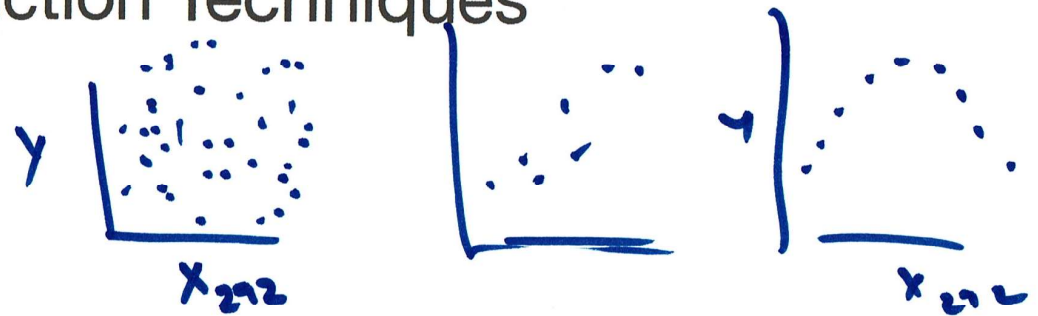


$$y = a + b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_{2000} \cdot x_{2000}$$

Dim. Reduction Techniques

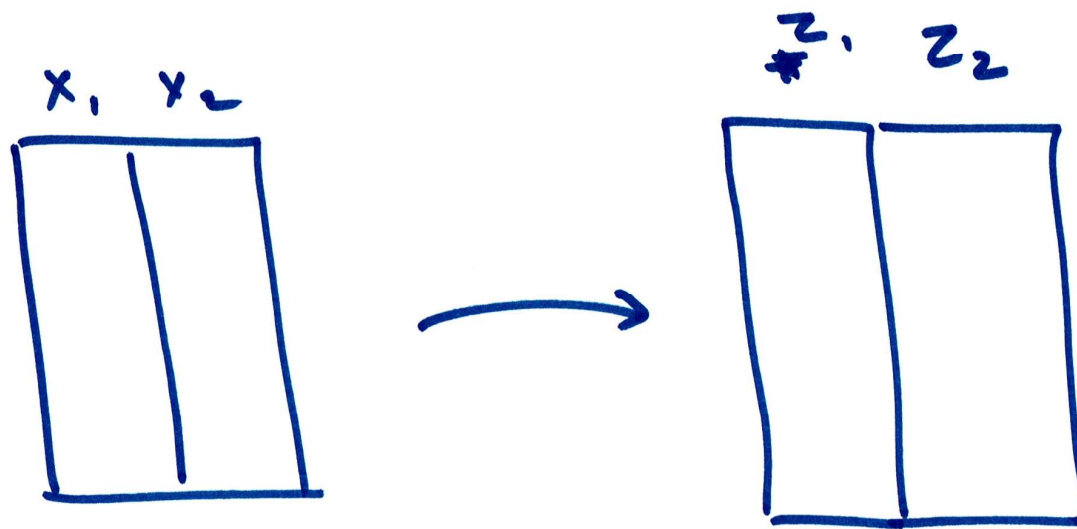


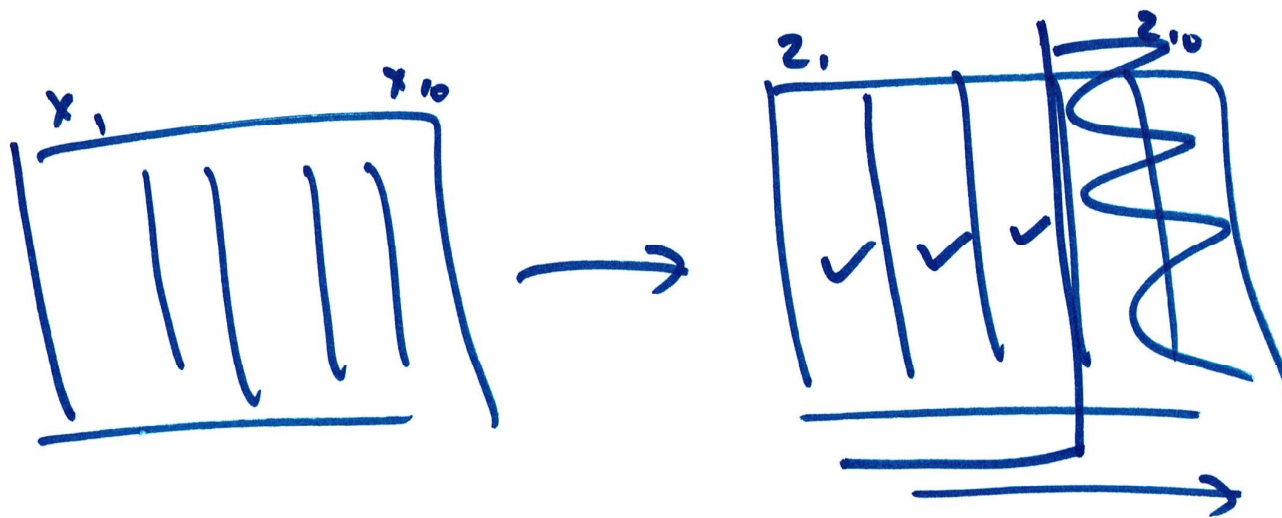
- Feature elimination

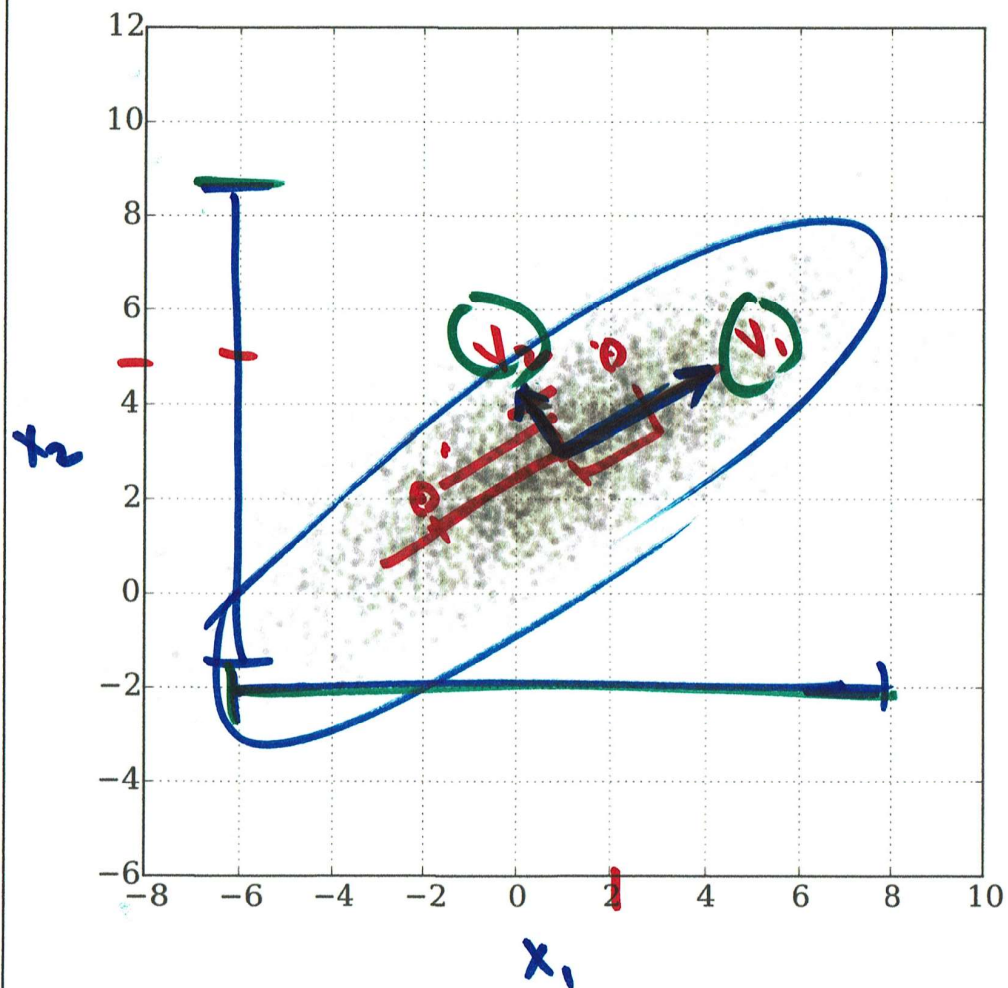
- Simply identify and remove variables (columns) that are not important
- The disadvantage is that we would gain no insight from those dropped variables and lose any information they contain

- Feature extraction

- Create a few new variables from the old variables
- **PCA** Principal Component Analysis: is the most popular feature extraction technique (linear)
- t-SNE (non-linear) ←







Handwritten matrix calculations and the resulting principal component equations:

$$\begin{bmatrix} 2 & 5 \\ -2 & 2 \end{bmatrix} \begin{matrix} x_1 \\ x_2 \end{matrix} \rightarrow \begin{bmatrix} 1 & 0.8 \\ -2 & 0.5 \end{bmatrix}$$

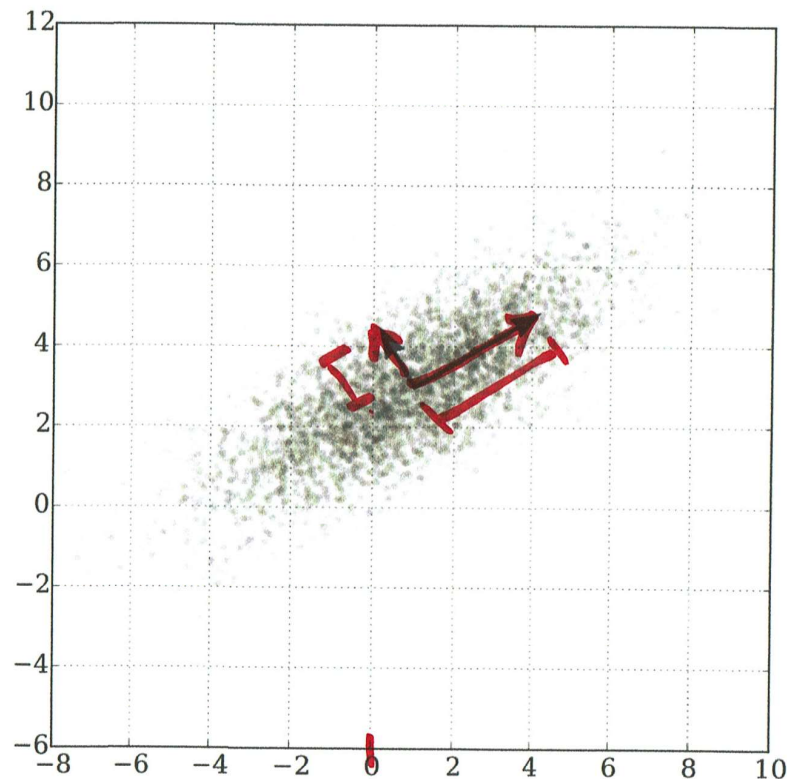
Below the first matrix, the values 1.2 and 0.9 are written. Below the second matrix, the values 1.9 and 0.1 are written, with a green bracket underneath.

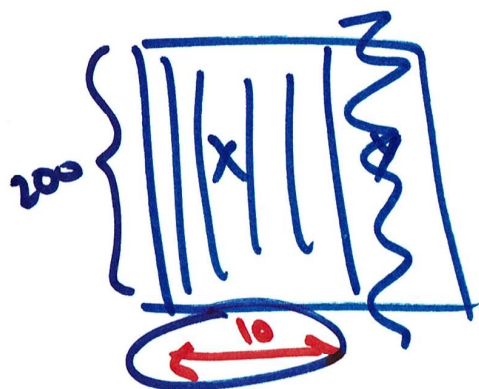
Below the matrices, a diagram shows a horizontal band of data points with a red line and a green line. Arrows point from the equations below to the lines in the diagram.

$$\begin{cases} Z_1 = \boxed{1.2} x_1 + \boxed{0.9} x_2 \\ Z_2 = \boxed{1.9} x_1 + \boxed{0.1} x_2 \end{cases}$$

PCA

- creates new variables using linear combinations of old variables
- is designed to create variables that are independent of one another
- also manages to tell us how important each of these new variables are
- this “importance”, helps us to choose how many variables we will use





$$X_{\text{new}} = \frac{X - \mu}{\sigma}$$

$$C_{\text{cov}} = \text{Cov}(X_{\text{new}})$$

Eigen decomposition

$$\begin{array}{c} \rightarrow \begin{array}{|c|c|c|c|} \hline e_1 & e_2 & \dots & e_{10} \\ \hline v_1 & v_2 & \dots & v_{10} \\ \hline \end{array} \end{array}$$

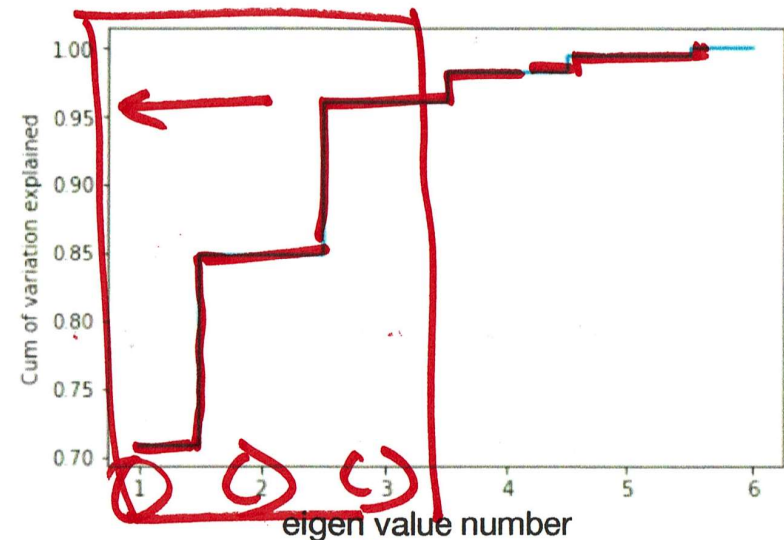
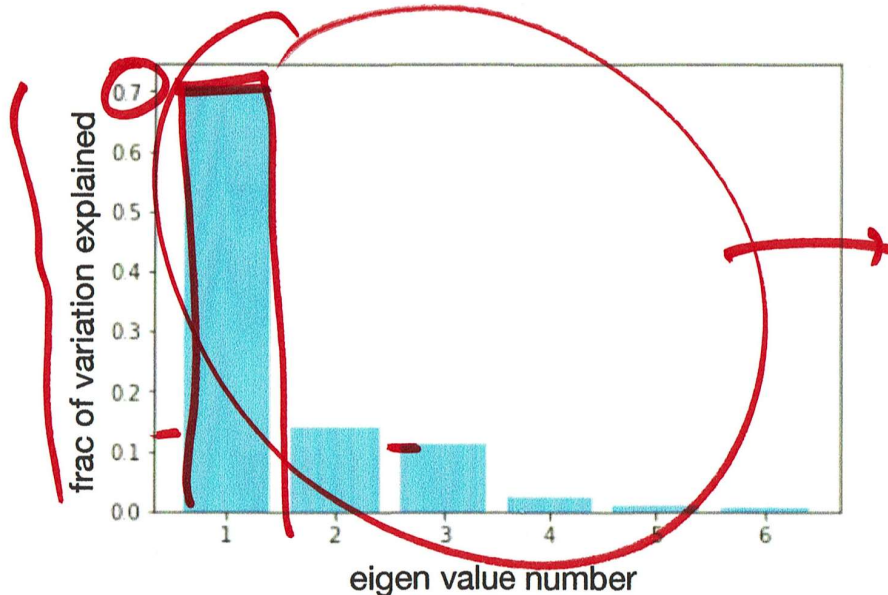
Z

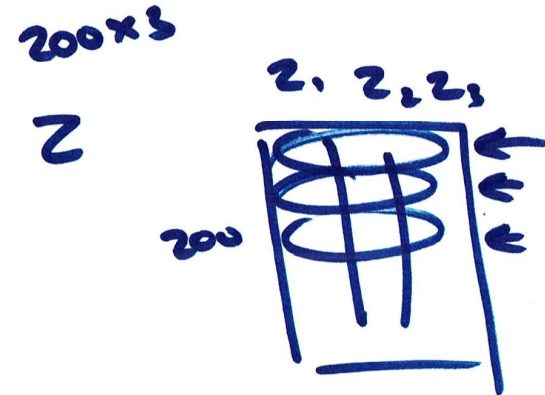
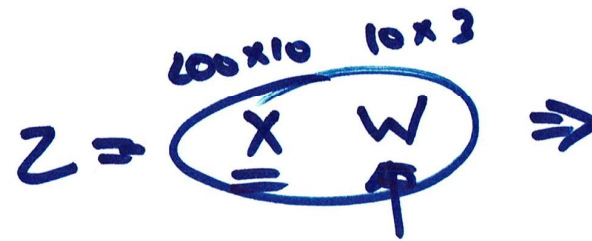
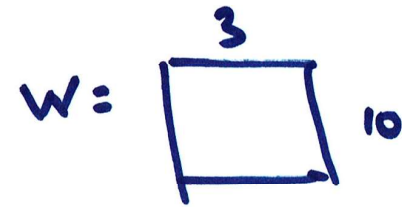
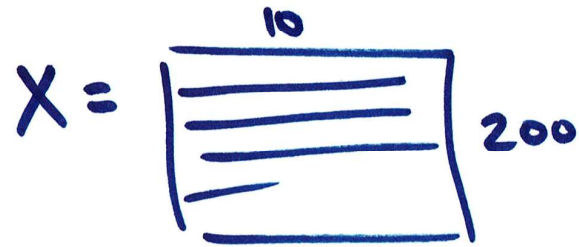
10×3

$$\leftarrow W = \begin{pmatrix} v_1 & v_2 & v_3 \end{pmatrix}$$

- Scale the data and compute the covariance matrix
- Break the covariance matrix into magnitude and direction. Eigen Vectors and the Eigen Values of the covariance matrix can be thought of as the natural axis/directions and magnitudes along those axis, of the data
- The eigen values also can be used to calculate the percentage of variation explained by each component
- Sort in the eigen values in desending order and calculate the cumulative percentage of variation explained
- Pick the number of principal components you will use
- Transform to new variables

$$\frac{e_1}{\sum e_i}, \frac{e_2}{\sum e_i}, \frac{e_3}{\sum e_i}$$





$$S = 100 + 20 \underline{Z_1} + 15 Z_2 - 10 Z_3$$

$$Z_1 = \square x_1 + \square x_2 + \square x_3 + \dots + x_{10}$$

$$Z_2 = \dots$$

$$Z_3 = \dots$$