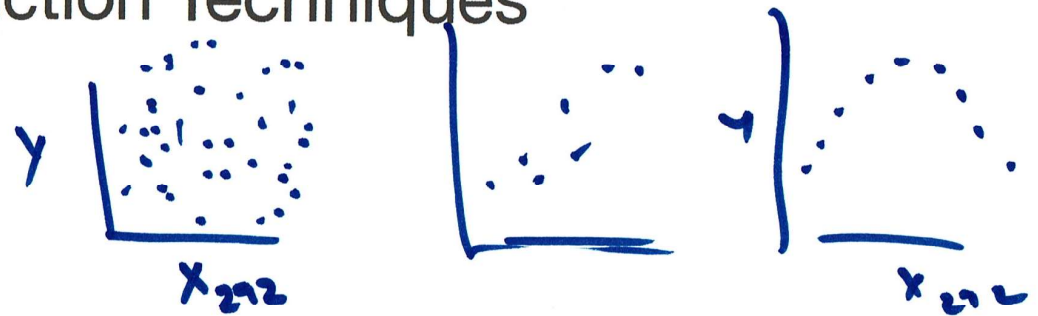$$y = a + b_1 \square + b_2 \square \ldots \ldots \ldots + b_{2000} \square$$
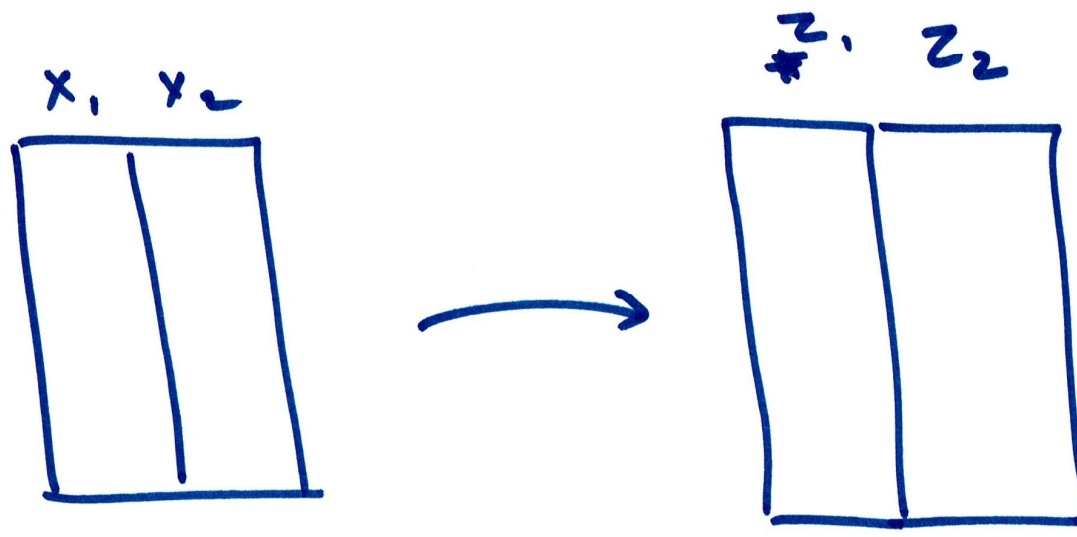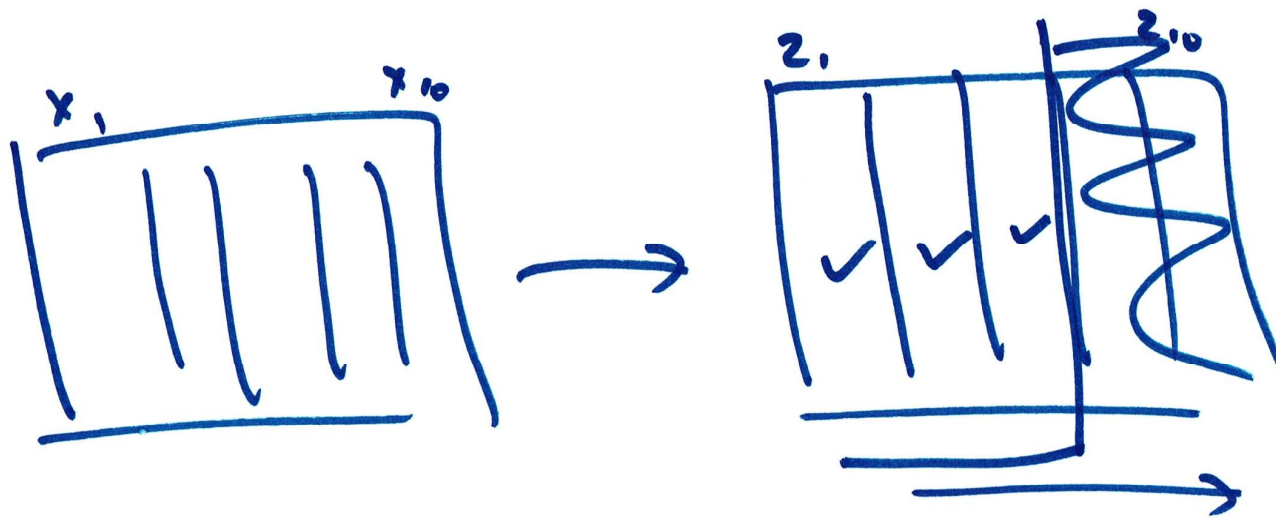
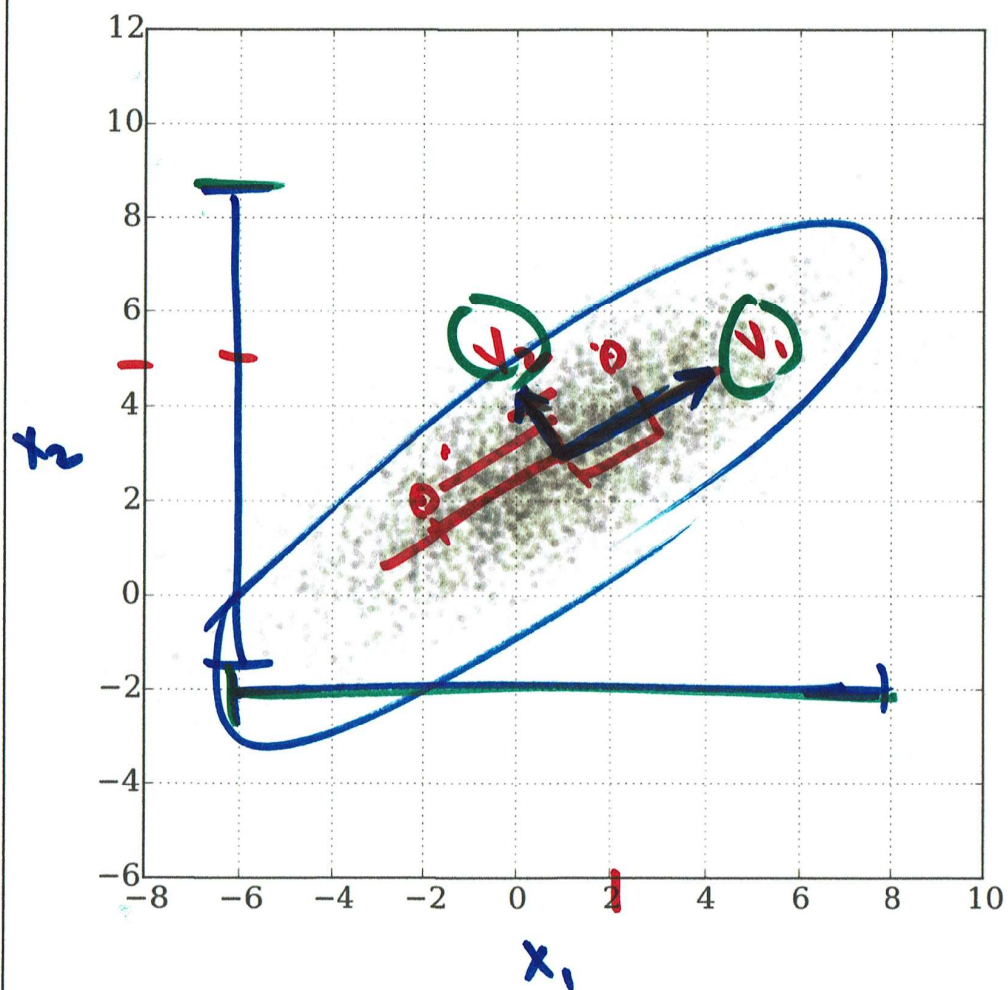# Dim. Reduction Techniques



- 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 loose 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)

$X_1 \quad Y_2$



$Z_1 \quad Z_2$

$x_1$ $x_{10}$

2. $3_{10}$

$$x_1 \quad x_2$$

$$\begin{bmatrix} 2 & 5 \\ -2 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0.3 \\ -2 & 0.5 \end{bmatrix}$$

1.2   0.8

1.9   0.1

$$Z_1 = \boxed{\phantom{x}}\, x_1 + \boxed{\phantom{x}}\, x_2$$

$$Z_2 = \boxed{\phantom{x}}\, x_1 + \boxed{\phantom{x}}\, x_2$$

6

# 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

$$X_{now} = \frac{X - \mu}{\sigma}$$

$$C \Sigma_m = Cov(X_{now})$$

Eigen decomposition

$$\begin{array}{|cccc|}
\hline
e_1, & e_2, & \ldots & e_{10} \\
\hline
v_1, & v_2 & \ldots & v_{10} \\
\hline
\end{array}$$

$$\boxed{Z}$$

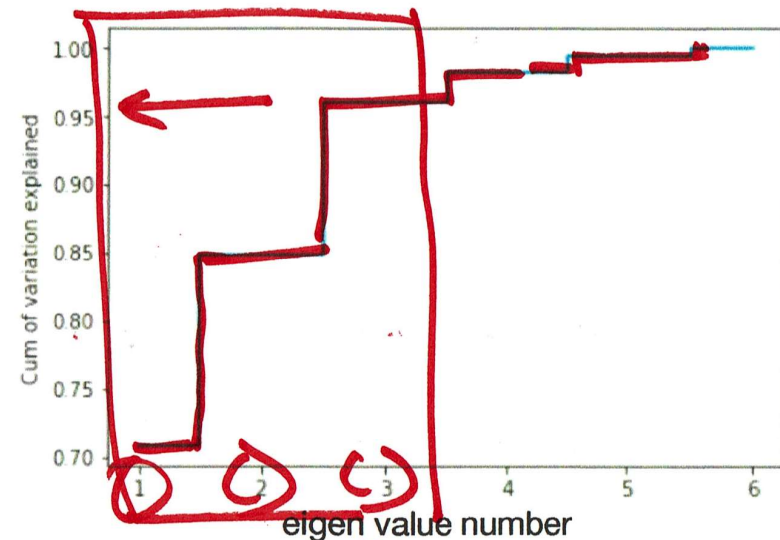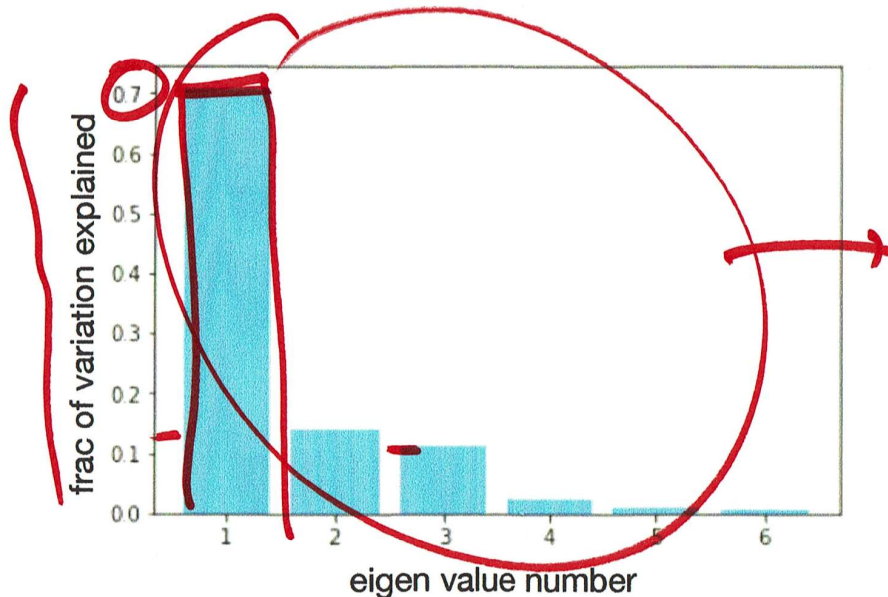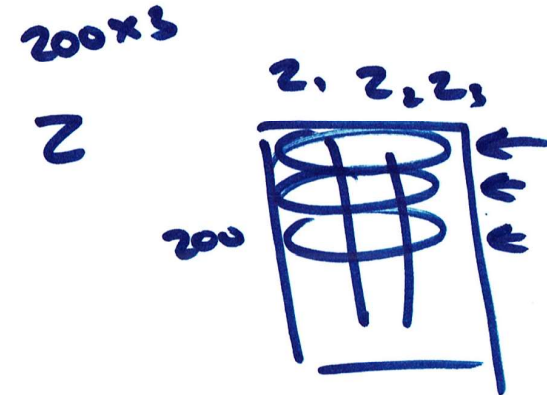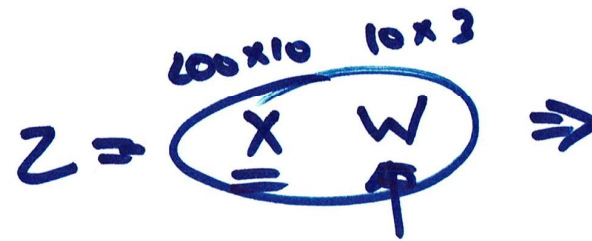$$\boxed{10 \times 3} \leftarrow W = \underset{10}{\Bigg\{} \left( v_1 \ \Big| \ v_2 \ \Big| \ v_3 \right)$$

- 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

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

- Transform to new variables

$$X = \boxed{\begin{array}{c} \phantom{xxxxxx} \\ \phantom{xxxxxx} \\ \phantom{xxxxxx} \\ \phantom{xxxxxx} \end{array}}$$

10 (top), 200 (right)

$$W = \boxed{\phantom{xx}}$$

3 (top), 10 (right)

$$Z = \underset{\substack{200 \times 10 \quad 10 \times 3}}{\left( X \quad W \right)} \Rightarrow \underset{200 \times 3}{Z}$$

$z_1, z_2, z_3$

200 (left of grid)

$$z_1 = \boxed{\phantom{x}} \, x_1 + \boxed{\phantom{x}} \, x_2 + \boxed{\phantom{x}} \, x_3$$

$\cdots \quad x_{10}$

$$z_2 = \quad \cdots$$

$$z_3 = \quad \cdots$$

$$S = 100 + 20 \, \widehat{z_1}$$
$$+ 15 \, z_2$$
$$- 10 \, z_3$$