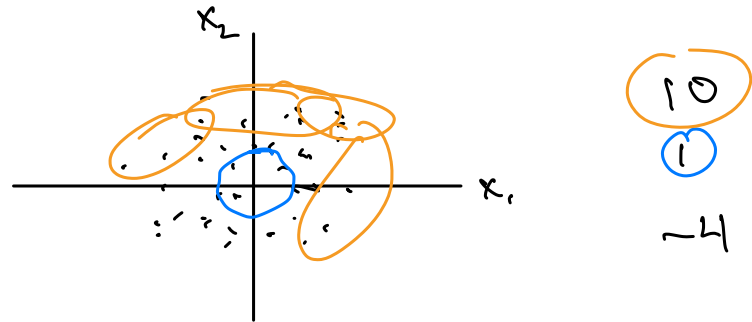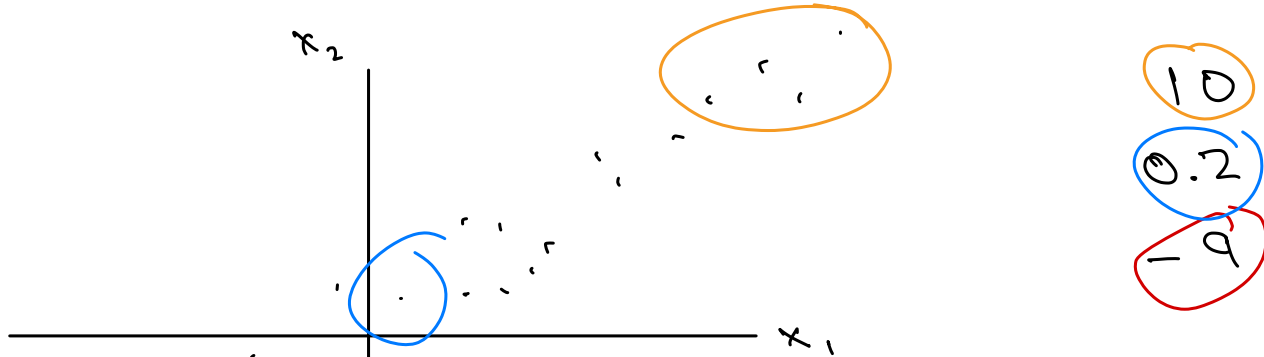# 11 Principal Components Analysis

PCA is an <u>unsupervised</u> learning technique — there is no response.
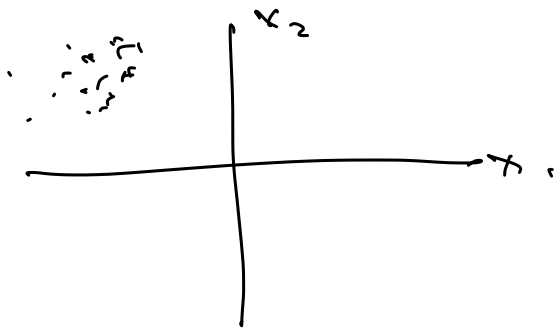
# Assumptions

Have $\underline{n}$ observations of $\underline{p}$ variables, $\underline{x}_1, \ldots, \underline{x}_n$

$$\underline{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} = \left[ \begin{pmatrix} \text{course overall} \\ \text{instr overall} \\ \text{hrs/week spent on HW} \end{pmatrix} \right]$$

Moreover, assume $\underline{x}_i \leftarrow \underline{x}_i - \overline{\underline{x}}$ are centered)

If M is a symmetric real $p \times p$ matrix its __eigen decomposition__ (or __spectral decomposition__) is

$$M = A^T D A$$

where $A^T$ is $p \times p$ matrix whose columns are eigenvectors of M & $D$ is diag matrix of eigenvalues.

$$A^T = \begin{bmatrix} a_1 & a_2 & \cdots & a_p \end{bmatrix} \implies A = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_p^T \end{bmatrix}$$

↳ eigenvects

⟶ rows are
e.vecs.

Use convention that $a_1, \ldots, a_p$ are normalized so A is orthogonal matrix $A^T A = A A^T = I$. A diagonalizes M:

$$A M A^T = D.$$

Ex  $n=3$, $p=2$

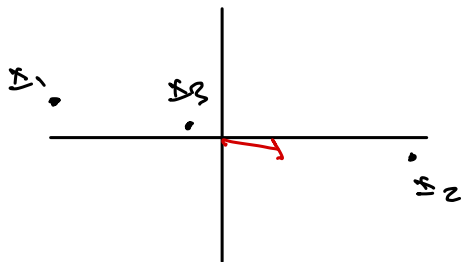standard basis vectors

$$x_1 = \begin{pmatrix} -10 \\ 1 \end{pmatrix} = -10\,\underline{e}_1 + 1\cdot\underline{e}_2 \quad \approx -10.5\,\underline{a}_1$$

$$x_2 = \begin{pmatrix} 12 \\ -1/2 \end{pmatrix} = 12\,\underline{e}_1 - \tfrac{1}{2}\,\underline{e}_2 \quad \approx 12.2\,\underline{a}_1$$

$$x_3 = \begin{pmatrix} -2 \\ 1/2 \end{pmatrix} = -2\,\underline{e}_1 + \tfrac{1}{2}\,\underline{e}_2 \quad \approx -1.8\,\underline{a}_1$$



Basic idea behind PCA!

$$\underline{x} = b_1\,\underline{e}_1 + b_2\,\underline{e}_2 + \cdots + b_p\,\underline{e}_p$$

Re-express as:

$$\underline{x} = \underbrace{a_1}_{\text{biggest}}\,\underline{a}_1 + \underbrace{a_2}_{\text{next biggest}}\,\underline{a}_2 + \cdots + \underbrace{a_p}_{\text{smallest}}\,\underline{a}_p$$

## 11.1 Probabilistic Approach

[Ex] Suppose $\underline{x}$ is a random vector with cov. matrix $\Sigma$, (mean zero)
want matrix $A$ that decorrelates $\underline{x}$:

$$Var(A\underline{x}) = [\text{want to be}] = D = \text{diagonal matrix.}$$

$\underline{z} = A\underline{x}$ are uncorrelated.

$$Var(A\underline{x}) = A\Sigma A^T = D$$

$[\Sigma$ is real, $p \times p$, symmetric $] \Longrightarrow$ Using $A$ from the
eigen decomp of $\Sigma$ will be the solution. If
rows of $A$ are eigvecs of $\Sigma$, then

$$\Sigma = A^T D A \quad + \quad A\Sigma A^T = D \quad + \quad D \text{ holds eigenvalues of } \Sigma.$$

**Note** Same argument for finite sample version.

$\underline{x}_1, \ldots, \underline{x}_n$ over $p$ features, need a convention

for storing in a matrix:

$$X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

$n \times p$

(usual design matrix
w/ast column of 1s)

$$X^T = \begin{bmatrix} \underline{x}_1 & \underline{x}_2 & \cdots & \underline{x}_n \end{bmatrix} \quad p \times n.$$

**Def** The <u>sample covariance matrix</u> of $X$ is

$$\hat{\Sigma} = \frac{1}{n-1} X^T X \quad (\text{psp matrix})$$

$= $ covariance matrix of "old features"

Want: A set of "new features" $z_i = A x_i$
that are uncorrelated & sorted in order of decreasing
variability.

Store $z_1, \ldots, z_n$ in matrix

$$Z^T = [z_1 \; z_2 \cdots z_n]$$

$$Z = \begin{bmatrix} z_{11} & z_{12} & \cdots & z_{1p} \\ z_{21} & z_{22} & & z_{2p} \\ \vdots & \vdots & & \vdots \\ z_{n1} & z_{n2} & & z_{np} \end{bmatrix} = n \times p \text{ design matrix of}$$

$p$ "new features"

$$Z^T = [\underline{z}_1, \underline{z}_2 \cdots \underline{z}_n] = [A\underline{x}_1, A\underline{x}_2 \cdots A\underline{x}_n] = A X^T$$

want : cov matrix of $\underline{z}$ to be diagonal

$$Diag = \hat{\Sigma}_z = \frac{1}{n-1} Z^T Z = \frac{1}{n-1}(A X^T)(X A^T)$$

$$= \frac{1}{n-1} A X^T X A^T = A \left(\frac{1}{n-1} X^T X\right) A^T = A \hat{\Sigma} A^T$$

$\Rightarrow$ Set rows of $A$ to be eigenvectors $\hat{\Sigma}$

If you do this, then $\hat{\Sigma}_z = $ diag with eigenvalues of $\hat{\Sigma}$

$$\hat{\Sigma}_z = diag(\lambda_1, \lambda_2, \cdots, \lambda_p)$$

where $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0$

○ $\lambda_i$ is sample variance of "new features"

$$z_{11}, z_{21}, \cdots, z_{n1}$$

- $d_j$ is sample variance of $j$th "new feature"

$$z_{1j}, z_{2j}, \ldots, z_{nj}$$

---

**PCA**

$$\underline{x}_i = x_{i1}\,\underline{e}_1 + x_{i2}\,\underline{e}_2 + \cdots + x_{ip}\,\underline{e}_p$$

$$= z_{i1}\,\underline{a}_1 + z_{i2}\,\underline{a}_2 + \cdots + z_{ip}\,\underline{a}_p$$

eigenvectors

principal components ("new features")
= scores = component scores
= "loadings"

$(z_{i1}\,\underline{a}_1)$
= loading