

## Dimensionality Reduction

- Some of the most interesting applications of machine learning involve data that lives in a **high-dimensional space**.
  - \* Ex: An image that is 1000 pixels wide and 800 pixels tall lives in **800000-dimensional space**.
- Moreover, the number of examples  $n$  in our dataset may be significantly lower than the **dimension  $d$** .
  - As a result, classifiers can end up overfitting the data.
- Ideally, we would like to find a **lower-dimensional representation** of the dataset that preserves the important relationships between examples.
  - This is an unsupervised learning problem.
- We will focus on **Principal Component Analysis**, which is one of the simplest and most popular methods.

• It will be useful to recall a few important properties of the covariance matrix  $\Sigma_{\underline{x}} = \mathbb{E}[(\underline{X} - \mathbb{E}[\underline{X}])(\underline{X} - \mathbb{E}[\underline{X}])^T]$ .

→ All of the **eigenvalues** are real and non-negative.

Using this fact, we can sort them (along with the corresponding eigenvectors) into descending order  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ .

→ All of the **eigenvectors** are orthonormal to each other.

Specifically, let  $\underline{v}_1, \dots, \underline{v}_d$  be the eigenvectors (sorted as above),

$$\text{then } \underline{v}_i^T \underline{v}_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

→ Collecting the eigenvectors into a matrix  $\mathbf{V} = [\underline{v}_1 \dots \underline{v}_d]$  and the eigenvalues into a diagonal matrix  $\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_d \end{bmatrix}$  we can write the eigendecomposition

$$\Sigma_{\underline{x}} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \quad \text{where } \mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I} \quad \text{Orthogonal Matrix}$$

→ These properties hold for sample covariance matrices too.

- Given a random vector  $\underline{X}$  with mean vector  $\underline{\mu}_x = \mathbb{E}[\underline{X}]$  and covariance matrix  $\Sigma_x = V \Lambda V^T$ , we can carry out the following coordinate transformation:

→ Center the distribution at the origin by subtracting the mean:

$$\underline{X}_{\text{centered}} = \underline{X} - \underline{\mu}_x$$

→ Rotate the coordinate system by multiplying by the orthogonal matrix of eigenvectors:

$$\underline{Z} = V^T \underline{X}_{\text{centered}} = V^T (\underline{X} - \underline{\mu}_x)$$

No redundancy between transformed features.

- The entries of  $\underline{Z}$  are uncorrelated with each other,  $\text{Cov}[z_i, z_j] = 0$  and their variances are equal to the eigenvalues  $\text{Var}[z_i] = \lambda_i$ ,  $i \neq j$

The covariance matrix is  $\Sigma_z = \underbrace{V^T}_{\text{Covariance of Linear Transform}} \Sigma_x \underbrace{V}_{\text{I}} = \underbrace{\Lambda}_{\text{diagonal}}$

- To reduce the dimension to  $k$ , we simply throw out all but the top  $k$  entries of  $\underline{Z}$ .

- **Principal component analysis** follows the same steps, except that we also need to estimate the mean vector and covariance matrix from the data  $\underline{x}_1, \dots, \underline{x}_n$ . ← Usually training data.

→ Collect the  $n$   $d$ -dimensional examples into an  $n \times d$  data matrix:

$$X = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix}$$

All-ones vector  
 $\underline{1}_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$   $\uparrow$   $n$

→ Compute the sample mean vector:  $\hat{\underline{\mu}}_x = \frac{1}{n} \sum_{i=1}^n \underline{x}_i = \frac{1}{n} X^T \underline{1}_n$

→ Center the data  $\underline{x}_{\text{centered},i} = \underline{x}_i - \hat{\underline{\mu}}_x$ ,  $X_{\text{centered}} = X - \underline{1}_n \hat{\underline{\mu}}_x^T$

→ Compute the sample covariance matrix:

$$\hat{\underline{\Sigma}}_x = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \hat{\underline{\mu}}_x)(\underline{x}_i - \hat{\underline{\mu}}_x)^T = \frac{1}{n-1} \sum_{i=1}^n \underline{x}_{\text{centered},i} \underline{x}_{\text{centered},i}^T$$

$$= \frac{1}{n-1} (X - \underline{1}_n \hat{\underline{\mu}}_x^T)^T (X - \underline{1}_n \hat{\underline{\mu}}_x^T) = \frac{1}{n-1} X_{\text{centered}}^T X_{\text{centered}}$$

→ Compute the eigendecomposition  $\hat{\underline{\Sigma}}_x = V \Lambda V^T$ .

• Principal Component Analysis:

→ Given a dataset  $X$ , compute the sample mean vector  $\hat{\underline{\mu}}_x$  and sample covariance matrix  $\hat{\underline{\Sigma}}_x$ .

→ Compute the eigendecomposition

$$\hat{\underline{\Sigma}}_x = V \Lambda V^T$$

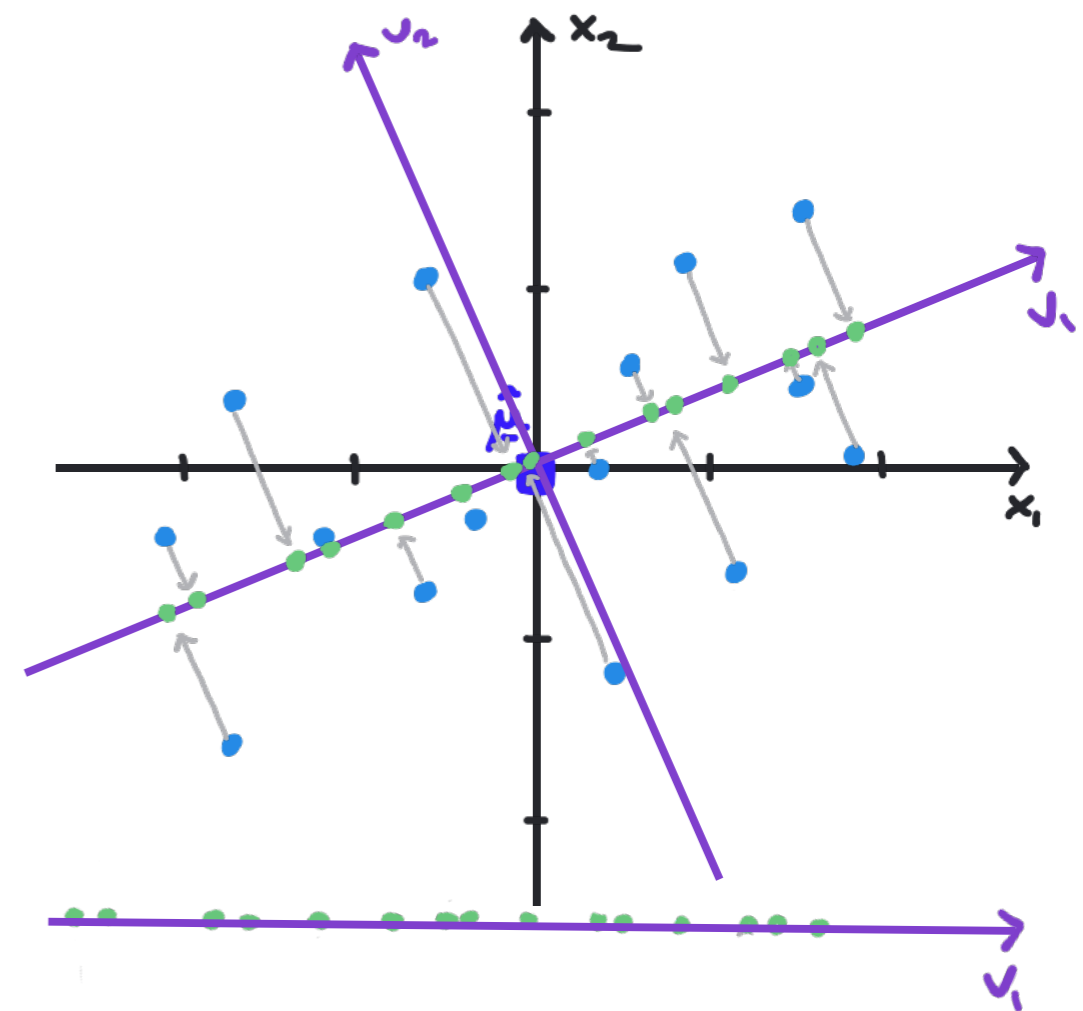
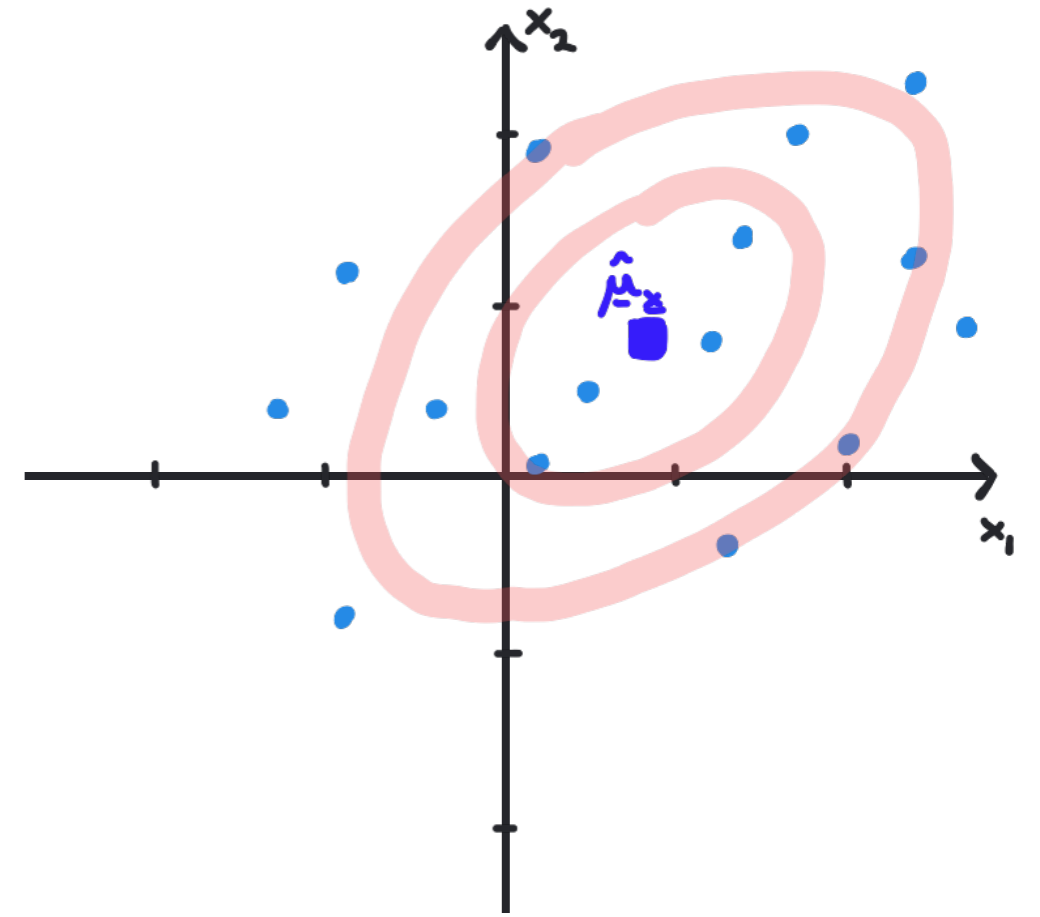
→ Keep only the first  $k$  eigenvectors

$$V_k = [\underline{v}_1 \dots \underline{v}_k]$$

→ Center the data:  $X_{\text{centered}} = X - \mathbf{1}_n \hat{\underline{\mu}}_x^T$

→ Project onto  $V_k$ :  $Z = X_{\text{centered}} V_k$

→ We can use  $Z$  as a  $k$ -dimensional representation of our original  $d$ -dimensional dataset  $X$  where  $k < d$ .



- One important application of PCA is visualizing high-dimensional datasets in 2 or 3 dimensions.
- Another application is as a pre-processing step for supervised learning to avoid overfitting:

→ Estimate the sample mean vector and covariance matrix from the training data only:

$$\hat{\underline{\mu}}_{\underline{x}} = \frac{1}{n_{\text{train}}} \underline{X}_{\text{train}} \quad \hat{\underline{\Sigma}}_{\underline{x}} = \frac{1}{n_{\text{train}} - 1} (\underline{X}_{\text{train}} - \underline{1}_{n_{\text{train}}} \hat{\underline{\mu}}_{\underline{x}}^T)^T (\underline{X}_{\text{train}} - \underline{1}_{n_{\text{train}}} \hat{\underline{\mu}}_{\underline{x}}^T)$$

→ Compute the eigendecomposition  $\hat{\underline{\Sigma}}_{\underline{x}} = \underline{V} \underline{\Lambda} \underline{V}^T$  and  $\underline{V}_k = [\underline{v}_1, \dots, \underline{v}_k]$ .

→ Center and project both the training and test data:

$$\underline{X}_{\text{train, reduced}} = (\underline{X}_{\text{train}} - \underline{1}_{n_{\text{train}}} \hat{\underline{\mu}}_{\underline{x}}^T) \underline{V}_k \quad \underline{X}_{\text{test, reduced}} = (\underline{X}_{\text{test}} - \underline{1}_{n_{\text{test}}} \hat{\underline{\mu}}_{\underline{x}}^T) \underline{V}_k$$

→ Use the reduced training data  $\underline{X}_{\text{train, reduced}}$  along with the original training labels  $\underline{Y}_{\text{train}}$  to train the algorithm (such as a classifier). Test it using the reduced test data  $\underline{X}_{\text{test, reduced}}$  along with the original test labels  $\underline{Y}_{\text{test}}$ .