

Support vector machines unit:

- ▶ What PCA is (last week Monday)
- ▶ Applications of PCA (last week Wednesday, in lab)
- ▶ The math of PCA (last week Friday)
- ▶ PCA algorithms (**Today**)
- ▶ (Begin neural nets on Wednesday)

Today:

- ▶ Finish the math of PCA, maximum-variance view
- ▶ Point of comparison: minimum-information-loss view
- ▶ Practical parts for implementation

The most important source for all of this was Deisenroth et al, *Mathematics for Machine Learning*, 2020, 286–293.

(Available on Canvas)

Let  $\bar{\mathbf{x}} = \mathbb{E}[\mathbf{x}] = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  be the mean of the data. We *center* the data so that it has mean  $\mathbf{0}$ :  $\mathbf{x}_n^C = \mathbf{x}_n - \bar{\mathbf{x}}$ . From this point on, we assume  $X$  has been centered, and so  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ .

Let  $\mathbf{W} = \begin{pmatrix} \mathbf{v}_0 \\ \vdots \\ \mathbf{v}_{M-1} \end{pmatrix} \in \mathbb{R}^{M \times D}$  be a set of orthonormal basis vectors.

Let  $\mathbf{z}_n = \mathbf{W}\mathbf{x}_n$  be the projection of  $\mathbf{x}_n$  in the vector space defined by these basis vectors.

We can project  $\mathbf{z}_n$  back into the original space with  $\tilde{\mathbf{x}}_n = \mathbf{W}^T \mathbf{z}_n$ .

Let  $\mathbf{v}_0 \in \mathbb{R}^D$  be the zeroth (“first”) principal component (which we want to find).

Let  $z_{0,n} = \mathbf{v}_0^T \mathbf{x}_n$  be the zeroth coordinate in the projection of  $\mathbf{x}_n$ .

Let  $z_0$  be a random variable modeling the value of the zeroth coordinate. Assume data is centered, that is,  $\mathbb{E}[z_0] = 0$ .

Let  $\mathbf{C} = \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{x}_n \mathbf{x}_n^T$  be the data covariance matrix. Then

$$\begin{aligned} \text{Var}[z_0] &= \frac{1}{N} \sum_{n=0}^{N-1} (z_{0,n} - 0)^2 \\ &= \frac{1}{N} \sum_{n=0}^{N-1} (\mathbf{v}_0^T \mathbf{x}_n)^2 \\ &= \frac{1}{N} \sum_{n=0}^{N-1} (\mathbf{v}_0^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{v}_0)^2 \\ &= \mathbf{v}_0^T \left( \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{v}_0 \\ &= \mathbf{v}_0^T \mathbf{C} \mathbf{v}_0 \end{aligned}$$

which we want to maximize subject to  $\|\mathbf{v}_0\|^2 = 1$

Optimization problem as a Lagrangian:

$$\mathcal{L}(\mathbf{v}_0, \lambda_0) = \mathbf{v}_0 \mathbf{C} \mathbf{v}_0 + \lambda_0(1 - \mathbf{v}_0^T \mathbf{v}_0)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_0} = 1 - \mathbf{v}_0^T \mathbf{v}_0$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}_0} = 2\mathbf{v}_0^T \mathbf{C} - 2\lambda_0 \mathbf{v}_0^T = 0$$

$$\mathbf{v}_0^T \mathbf{C} = \lambda_0 \mathbf{v}_0^T$$

Plug that into our formula for the variance:

$$\text{Var}[z_0] = \mathbf{v}_0^T \mathbf{C} \mathbf{v}_0$$

$$= \lambda_0 \mathbf{v}_0^T \mathbf{v}_0 = \lambda_0$$

To find the next principal component, transform the data by removing the effect of the principal component  $\mathbf{v}_0$ :

$$\mathbf{X}_1 = \mathbf{X} - \mathbf{v}_0 \mathbf{v}_0^T \mathbf{X}$$

... compute the corresponding data covariance matrix  $\mathbf{C}_1$ , find the eigenvector with greatest eigenvalue. As a loop to find  $M$  principal components  $\mathbf{v}_0, \dots, \mathbf{v}_{M-1}$ :

$\mathbf{C}$  = the data covariance matrix of  $\mathbf{X}$

$\mathbf{v}_0$  = the eigenvector of  $\mathbf{C}$  with greatest eigenvalue

for  $m \in [1, M)$  :

$$\mathbf{W}_m = \sum_{i=0}^{m-1} \mathbf{v}_i \mathbf{v}_i^T \quad (\text{projection matrix into subspace})$$

$$\mathbf{X}_m = \mathbf{X} - \mathbf{W}_m \mathbf{X}$$

$\mathbf{C}_m$  = the data covariance matrix of  $\mathbf{X}_m$

$\mathbf{v}_m$  = the eigenvector of  $\mathbf{C}_m$  with greatest eigenvalue

## Theorem 1 (Invariant)

*For all  $m$ , every eigenvector of  $C$  is an eigenvector of  $C_m$ .*

For all  $m$ , every eigenvector of  $C$  is an eigenvector of  $C_m$ .

**Proof.** Suppose  $\mathbf{v}_j$  is an eigenvector of  $C$  with eigenvalue  $\lambda_j$ , that is,  $C\mathbf{v}_j = \lambda_j\mathbf{v}_j$ . Then

$$\begin{aligned} \mathbf{C}_m\mathbf{v}_j &= \frac{1}{N}\mathbf{X}_m\mathbf{X}_m^T\mathbf{v}_j \\ &= \frac{1}{N}(\mathbf{X} - \mathbf{W}_m\mathbf{X})(\mathbf{X} - \mathbf{W}_m\mathbf{X})^T\mathbf{v}_j && \text{do FOIL with } \frac{1}{N}\mathbf{X}\mathbf{X}^T = \mathbf{C} \\ &= (\mathbf{C} - \mathbf{C}\mathbf{W}_m - \mathbf{W}_m\mathbf{C} + \mathbf{W}_m\mathbf{C}\mathbf{W}_m)\mathbf{v}_j \\ &= \begin{cases} \mathbf{C}\mathbf{v}_j = \lambda_j\mathbf{v}_j & \text{if } j \geq m \\ \mathbf{C}\mathbf{v}_j - \mathbf{C}\mathbf{v}_j - \mathbf{C}\mathbf{v}_j + \mathbf{C}\mathbf{v}_j = 0 & \text{if } j < m \end{cases} \end{aligned}$$

In the  $j \geq m$  case,  $\mathbf{v}_j$  is orthogonal to all the vectors in  $\mathbf{W}_m$ , so  $\mathbf{W}_m\mathbf{v}_j = 0$ .

In the  $j < m$  case,  $\mathbf{v}_j$  is a basis vector of the subspace into which  $\mathbf{W}_m$  projects, so  $\mathbf{W}_m\mathbf{v}_j = \mathbf{v}_j$ . □

*To compute the principal components:*

Given **data** and  $M$  (number of desired components),

Center the data (and store the mean  $\bar{\mathbf{x}}$ )

Compute the covariance matrix

Compute the eigenvectors and corresponding eigenvalues

Sort the eigenvectors by eigenvalues

Return the  $M$  eigenvectors with greatest eigenvalues

*To transform a data point using principal components:*

Given data point  $\mathbf{x}$  and principal components  $\mathbf{v}_0, \dots, \mathbf{v}_{M-1}$ ,

Shift the data point based on the centering,  $\hat{\mathbf{x}} = \mathbf{x} - \bar{\mathbf{x}}$

Compute the dot products of  $\hat{\mathbf{x}}$  and each principal component

Assemble the results as a new vector, and return it

## Coming up:

### **Due Fri, Mar 28:**

*Textbook reading from Chapter 10 (see Canvas)*

### **Due Mon, Mar 31:**

*Take PCA quiz*

### **Due Fri, Apr 4:**

*Implement PCA*

### **Due Wed, Apr 9:**

*Read and respond to two articles about bias in algorithms  
(See Canvas)*

### **Sometime between Mar 31 and Apr 17:**

*Make an office-hours appointment for project check-in  
(Originally the deadline was Apr 11)*