

PRINCIPLE COMPONENT ANALYSIS (PCA)

Motivation example:

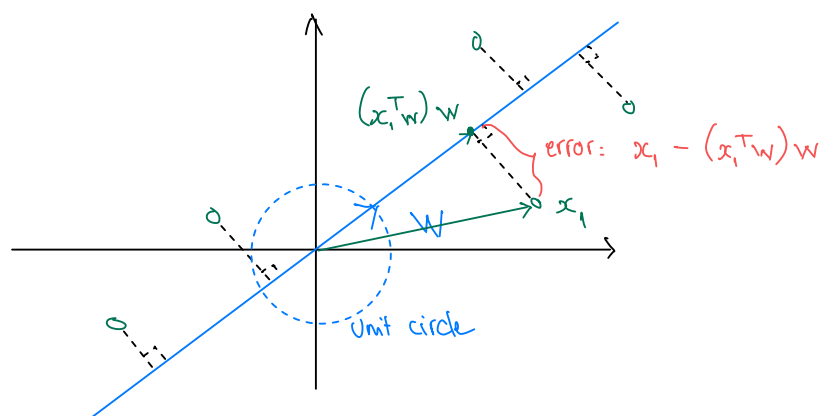
- PCA is common technique used in unsupervised learning (representative learning) for dimensions reduction.

What is unsupervised learning?

Unlike supervised learning, the dataset used in unsupervised learning doesn't have labels.

- PCA used to identify the principal vectors or hyperplanes that capture the **most variance** in the dataset. These vectors/hyperplanes can be used to **reconstruct the dataset** with the **least amount of errors**.

For example:



Goal: Find vector w of unit length

that minimize the errors:

$$\min_{\|w\|=1} \sum_{i=1}^m \|x_i - (x_i^T w)w\|^2$$

Solution: Vector w is the eigenvector correspond to the maximum eigenvalue of C

Question: Why vector w that minimize errors is also the vector that capture the most variance in the dataset?

$$\text{Given: } f(x) = \min_{\|w\|=1} \sum_{i=1}^m \|x_i - (x_i^T w)w\|^2$$

Some w minimize the sum, also minimize the average, so:

$$\begin{aligned} f'(x) &= \min_{\|w\|=1} \frac{1}{m} \sum_{i=1}^m \|x_i - (x_i^T w)w\|^2 \\ &= \min_{\|w\|=1} \frac{1}{m} \sum_{i=1}^m (x_i^T x_i - (x_i^T w)^2) \\ &= \min_{\|w\|=1} \frac{1}{m} \sum_{i=1}^m -(x_i^T w)^2 \end{aligned} \quad \begin{array}{l} \text{< minimize w.r.t. } w, \text{ remove} \\ \text{unrelated term } x_i^T x_i > \end{array}$$

This is the same as:

$$\begin{aligned} g(x) &= \max_{\|w\|=1} \frac{1}{m} \sum_{i=1}^m (x_i^T w)^2 \\ &= \max_{\|w\|=1} \frac{1}{m} \sum_{i=1}^m (x_i^T w)^T (x_i^T w) \\ &= \max_{\|w\|=1} w^T \left(\frac{1}{m} \sum_{i=1}^m x_i x_i^T \right) w \end{aligned}$$

covariance matrix C

So we conclude that:

$$\min_{\|w\|=1} \sum_{i=1}^m \|x_i - (x_i^T w) w\| \quad \text{equivalent} \quad \max_{\|w\|=1} w^T C w$$

\nwarrow
 w that minimize errors

\nearrow
 w that maximize the bilinear form of covariance matrix C

where $C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$

This leads to the solution explanation:

Recall eigenvector / eigenvalue is defined as:

$$Cw = \lambda w$$

$$\Leftrightarrow w^T C w = \lambda$$

$$\Leftrightarrow \max_{\|w\|=1} w^T C w = \max_{\|w\|=1} \lambda$$

Hence, w is the eigenvector corresponding to the maximize eigenvalue of C .

How does covariance matrix formed?

We know:

$$C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$$

$$= \frac{1}{m} (x_1 \mid \dots \mid x_m) \begin{pmatrix} x_1^T \\ \vdots \\ x_m^T \end{pmatrix}$$

$$= \frac{1}{m} X^T X$$

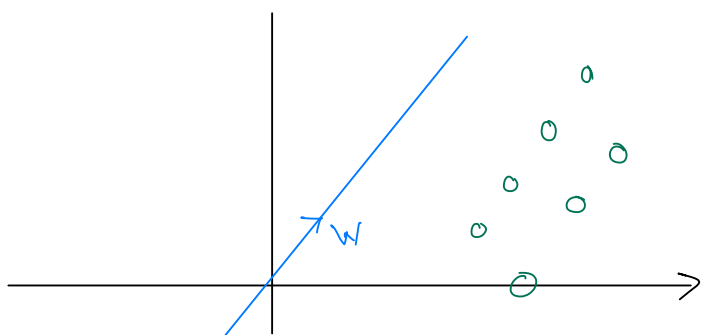
What if our training set has 2 features of length, one is measured in feet, the other is measured in centimeters?

You should normalize that shit, the formula is:

$$d_j := d_j - \sigma_j \quad \forall j \in [d], \quad \sigma_j = \sqrt{\frac{1}{m} \sum_{i=1}^m x_{ij}^2}$$

So far, we know how to find "best" line w that go through $(0,0)$.

What if our data is not centered?

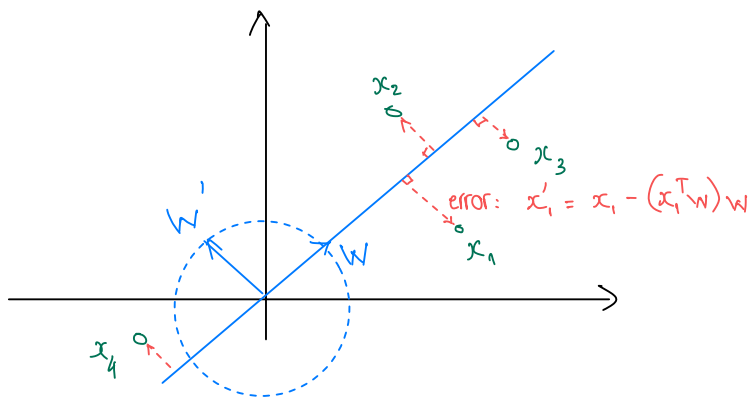


Solution: Before finding w , "center" all the data points:

- Find mean: $\mu = \frac{1}{m} \sum_{i=1}^m x_i$
- Subtract mean from data:

$$x_i = x_i - \mu \quad \forall i$$

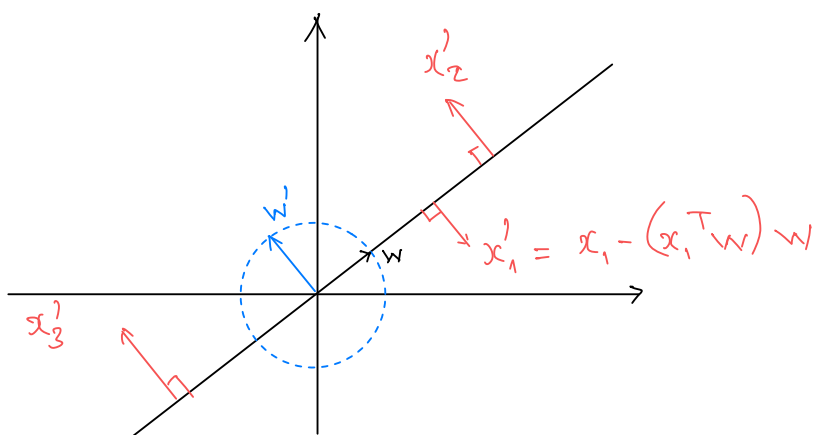
Now that we found w that can reconstruct the dataset with some errors, What if there is some information inside those errors that we also want to compress?



Solution: Find another w that can reconstruct those errors:

- Update: $S' = \{x'_1, \dots, x'_m\}$ where $x'_i = x_i - (x_i^T w)w$
- Find $w' = \arg \max_{\|w\|=1} w^T C' w$

Great, now we can find 2 vectors w and w' . What is the relationship between them?



Solution:

Since w' is the line "best" describe the errors, and all errors are perpendicular to w .

$$\Rightarrow w' \perp w$$

$$\Rightarrow w'^T w = 0$$

$$\text{for any } \begin{cases} w' = \arg \max_{\|w\|=1} w^T C' w \\ w = \arg \max_{\|w\|=1} w^T C w \\ C = \frac{1}{n} \sum_{i=1}^m x_i x_i^T, \quad C' = \frac{1}{m} \sum_{i=1}^m x'_i x'^T_i \\ x'_i = x_i - (x_i^T w)w \end{cases}$$

w' is the second relevant eigenvector

Now that we know how to iteratively find w that "compressed" the dataset, when should we stop iterating?

- The most intuitive answer is when the residuals/errors equals to 0.

Theoretically, this is always the case after d rounds:

$$x_i - \left[(x_i^T w_1)w_1 + \dots + (x_i^T w_d)w_d \right] = \vec{0} \in \mathbb{R}^d, \forall i \in [m]$$

$$\Rightarrow x_i = (x_i^T w_1) w_1 + \dots + (x_i^T w_d) w_d, \quad \forall i \in [m]$$

Since $w_1 \perp \dots \perp w_d$, we are essentially changing the basis of x_i

Why change the basis reduced dimension?

- Informally, the standard basis vectors are not always good at expressing the dataset. Some basis vector may contain more noise than information, some may be redundant.
- Formally, if the data exists in a low-dimensional space, then residues become 0 much earlier than d rounds.

What if the data "approximately" in a low-dimensional space?

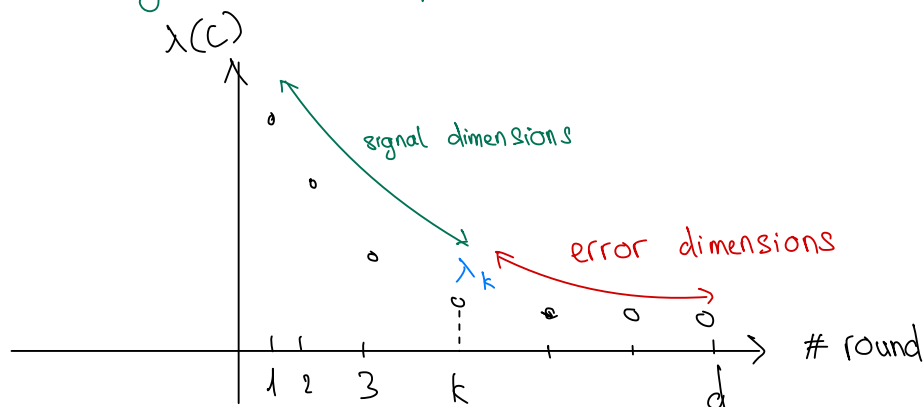
Consider this, using pythagoruous, we can say that:

$$\|x_i\|^2 = \|x_i - (x_i^T w) w\|^2 + \|(x_i^T w) w\|^2$$

$$\Leftrightarrow \frac{1}{m} \sum_{i=1}^m \|x_i\|^2 = \underbrace{\frac{1}{m} \sum_{i=1}^m \|x_i - (x_i^T w) w\|^2}_{\text{avg residual/error, as small as possible}} + \underbrace{\frac{1}{m} \sum_{i=1}^m \|(x_i^T w) w\|^2}_{\text{avg projection } x_i \text{ onto } w, \text{ as large as possible}}$$

We can prove that eigenvalue: $\lambda(C) = \frac{1}{m} \sum_{i=1}^m \|(x_i^T w) w\|^2$, C is cov matrix

"For each round, the larger $\lambda(C) = \frac{1}{m} \sum_{i=1}^m \|(x_i^T w) w\|^2$, the better the fit"



stop iterating here

Rule of thumb for # dimensions

We can say after compression, we retain 95% of information if:

$$\frac{\sum_{i=1}^k \lambda_i(C)}{\sum_{i=1}^d \lambda_i(C)} \geq 0.95$$

signal dimensions \rightarrow $\sum_{i=1}^k \lambda_i(C)$ \leftarrow total dimensions $\sum_{i=1}^d \lambda_i(C)$

PCA algorithm

Given $S = \{x_1, \dots, x_m\}$ $x_i \in \mathbb{R}^d$, "compressed" S .

Let $X \in \mathbb{R}^{m \times d}$ be the corresponding dataset matrix.

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_m^T \end{pmatrix}$$

1. Center X :

$$\bullet u = \frac{1}{m} \sum_{i=1}^m x_i \quad \forall i \in [m] \quad (u \in \mathbb{R}^d) \text{ mean vector}$$

$$\bullet x_i := x_i - u \quad \forall i \in [m]$$

$$\text{or } X := X - \mathbf{1} \cdot u^T$$

vector of ones size $m \times 1$

2. Normalize X :

$$\bullet \sigma_j = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2} \quad \forall j \in [d]$$

standard deviation feature j

$$\bullet x_{ij} := x_{ij} / \sigma_j \quad \begin{matrix} \forall i \in [m] \\ j \in [d] \end{matrix}$$

Step 1 and 2 can be thought of as Standardization:

$$\forall \begin{matrix} i \in [m] \\ j \in [d] \end{matrix}, \quad x_{ij} := \frac{x_{ij} - \mu_j}{\sigma_j}, \quad \text{where } \begin{cases} \mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij} \\ \sigma_j = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2} \end{cases}$$

3. Construct covariance matrix C :

$$\bullet C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T \quad \forall i \in [m]$$

$$C \in \mathbb{R}^{d \times d}$$

$$\text{or } C = \frac{1}{m} X^T X$$

4. Eigen decompose covariance matrix C :

$$\bullet C = Q \cdot D \cdot Q^T$$

eigenvectors

eigenvalues

$$= (Q_L \mid Q_R) \cdot \left(\begin{array}{c|c} D_{TL} & \\ \hline & D_{BR} \end{array} \right) \cdot \begin{pmatrix} Q_L^T \\ Q_R^T \end{pmatrix}$$

first k eigenvectors

5. Get the principal components:

The principal components matrix is: $X \cdot Q_L$

If you want to reconstruct the dataset using principle components, then:

$$X := (X Q_L) \cdot Q_L^T$$

Interpret the eigenvector chosen

Given any vector, covariance matrix will rotate it to roughly the same direction of the chosen eigenvector.