# Dimensionality Reduction with PCA

### Yufei Tao

#### Department of Computer Science and Engineering Chinese University of Hong Kong

Dimensionality Reduction with PCA

1/25

(日)

P = a set of points in  $\mathbb{R}^d$ , where the dimensionality d is large.

The goal of **dimensionality reduction** is to convert P into a set P' of points in a lower-dimensional subspace such that P' does not lose "too much" information about P.

We will learn a classical method called **principled component analysis** (PCA) to achieve the purpose.



Fix an integer  $k \leq d$ .

A *k*-subspace  $\Sigma$  is the span of *k* unit vectors  $\boldsymbol{u}_1, \boldsymbol{u}_2, ..., \boldsymbol{u}_k$  in  $\mathbb{R}^d$  that are mutually orthogonal. We refer to  $\{\boldsymbol{u}_1, ..., \boldsymbol{u}_k\}$  as a basis of  $\Sigma$ .

Recall that the span of  $\boldsymbol{u}_1, ..., \boldsymbol{u}_k$  is defined as

$$\Big\{ \sum_{i=1}^k c_i \boldsymbol{u}_i \mid \text{ any real values } c_1, c_2, ..., c_k \Big\}.$$

The span of two vectors is a plane:



## Subspace

A k-subspace with  $k \ge 2$  has an infinite number of bases.



Dimensionality Reduction with PCA

æ

4/25

イロト イボト イヨト イヨト

#### Projection

Let p be a point (a.k.a., a vector) in  $\mathbb{R}^d$ . Consider a k-subspace  $\Sigma$  with a basis  $\{u_1, ..., u_k\}$ . For each  $i \in [1, k]$ , the coordinate of p on  $u_i$  is  $p \cdot u_i$ . The projection of p onto  $\Sigma$  is

$$\hat{\boldsymbol{p}} = \Sigma_{i=1}^k c_i \boldsymbol{u}_i.$$

where  $c_i = \boldsymbol{p} \cdot \boldsymbol{u}_i$ . Note:  $\hat{\boldsymbol{p}}$  is a *d*-dimensional point.



**Dimensionality Reduction with PCA** 

#### Projection

Let  $\boldsymbol{p}$  be a point (a.k.a., a vector) in  $\mathbb{R}^d$ . Consider a k-subspace  $\boldsymbol{\Sigma}$  with basis  $\{\boldsymbol{u}_1, ..., \boldsymbol{u}_k\}$ . Let  $\hat{\boldsymbol{p}}$  be the projection of  $\boldsymbol{p}$  onto  $\boldsymbol{\Sigma}$ .

**Lemma 1:** 
$$|\boldsymbol{p}|^2 = |\boldsymbol{p} - \hat{\boldsymbol{p}}|^2 + \sum_{i=1}^k (\boldsymbol{p} \cdot \boldsymbol{u}_i)^2$$
.

**Proof:** Follows immediately from  $|\hat{\boldsymbol{p}}|^2 = \sum_{i=1}^{k} (\boldsymbol{p} \cdot \boldsymbol{u}_i)^2$  and  $|\boldsymbol{p}|^2 = |\boldsymbol{p} - \hat{\boldsymbol{p}}|^2 + |\hat{\boldsymbol{p}}|^2$ .



**Dimensionality Reduction with PCA** 

*P* : a set of *n* points in  $\mathbb{R}^d$  whose geometry center is the origin.

**Think:** If the geometry center of *P* is not the origin, what transformations do you need to make this happen?

**Goal:** Find a *k*-subspace  $\Sigma$  to minimize its **reconstruction error** defined as

$$rac{1}{n}\sum_{oldsymbol{p}\in P}|oldsymbol{p}-\hat{oldsymbol{p}}|^2$$

where  $\hat{\boldsymbol{p}}$  is the projection of  $\boldsymbol{p}$  onto  $\boldsymbol{\Sigma}$ .

Recall that the geometry center of our point set P is the origin.

Define the **covariance** of *P* between dimensions *i* and *j*  $(i, j \le d)$  as

$$\frac{1}{n}\sum_{p\in P}p[i]\cdot p[j].$$

where p[i] is the coordinate of p on dimension i.

**Think:** What is the coordinate mean of the points in *P* on dimension *i*?

The covariance matrix **A** of *P* is a  $d \times d$  matrix where **A**[*i*,*j*]  $(i, j \in [1, d])$  is the covariance of *P* between dimensions *i* and *j*.

Note that **A** is symmetric, namely,  $\mathbf{A} = \mathbf{A}^{T}$ .

8/25

・ロト ・ 一 マ ・ コ ト ・ 日 ト

Let A be the covariance matrix defined on the previous slide. If

$$Av = \lambda v$$

for a  $d \times 1$  <u>non-zero</u> vector **v** and some real value  $\lambda$ , we call

- $\lambda$  an **eigenvalue** of **A**, and
- v an eigenvector of A.

We also say that  $\mathbf{v}$  is an **eigenvector for**  $\lambda$ .

**Lemma 2:** Matrix **A** has *d* eigenvalues  $\lambda_1, ..., \lambda_d$ . For each  $i \in [1, d]$ , denote by  $\mathbf{v}_i$  an eigenvector for  $\lambda_i$ . Then:

- $\lambda_1, ..., \lambda_d$  are all real values.
- For any distinct *i*, *j* ∈ [1, *d*], the vectors *v<sub>i</sub>* and *v<sub>j</sub>* are orthogonal to each other.

The lemma is a well-known result from linear algebra (proof omitted).

9/25

ロト (得) (ヨト (ヨト

Principle Component Analysis (PCA)

**PCA** (P, k)

/\* The geometry center of P is at the origin \*/

- 1.  $\mathbf{A} \leftarrow$  the covariance matrix of P
- 2. compute the eigenvalues of **A** and sort them in descending order:  $\lambda_1 \ge \lambda_2 \ge ... \ge \lambda_d$
- 3. compute a **unit** eigenvector  $\mathbf{v}_i$  for  $\lambda_i$  for each  $i \in [1, d]$
- 4. return  $v_1, ..., v_k$

10/25

(日本)

The rest of the slides will prove:

**Theorem 1:** The *k*-subspace with  $\{v_1, ..., v_k\}$  as the basis has the smallest the reconstruction error among all *k*-subspaces.

Dimensionality Reduction with PCA

11/25

▲ 同 ▶ → 三 ▶

Variance along a Direction

Fix an arbitrary unit vector  $\boldsymbol{u}$ . Define the variance of  $\boldsymbol{u}$  as

$$\frac{1}{n}\sum_{\boldsymbol{p}\in P}|\boldsymbol{p}\cdot\boldsymbol{u}|^2.$$

Indeed, the above is precisely the variance of the set  $\{ \boldsymbol{p} \cdot \boldsymbol{u} \mid \boldsymbol{p} \in P \}$ .

**Think:** What is the mean of  $\{\boldsymbol{p} \cdot \boldsymbol{u} \mid \boldsymbol{p} \in P\}$ ?

Dimensionality Reduction with PCA

12/25

< ロ > < 同 > < 回 > < 回 >

#### Variance along a Direction

**Lemma 3:** For any unit vector **u**, the variance of **u** equals

### u<sup>⊤</sup>Au

where  $\mathbf{A}$  is the covariance matrix of P.

**Proof:** Define **P** as the  $n \times d$  matrix whose row  $i \in [1, n]$  is  $p^T$ , where p is the *i*-th point of P (note: p is a  $d \times 1$  vector and hence  $p^T$  is a  $1 \times d$  vector). Rudimentary calculation shows

variance of 
$$\boldsymbol{u} = \frac{1}{n} (\mathbf{P}\boldsymbol{u})^T (\mathbf{P}\boldsymbol{u})$$
  
$$= \boldsymbol{u}^T \left(\frac{1}{n} \mathbf{P}^T \mathbf{P}\right) \boldsymbol{u}$$
$$= \boldsymbol{u}^T \mathbf{A} \boldsymbol{u}$$

as claimed

Variance along a Direction

**Lemma 4:** If  $\lambda$  is an eigenvalue of **A** and **v** is a unit eigenvector for  $\lambda$ , then the variance of **v** is  $\lambda$ .

**Proof:** 

variance of 
$$\mathbf{v} = \mathbf{v}^T \mathbf{A} \mathbf{v} = \mathbf{v}^T (\mathbf{A} \mathbf{v}) = \lambda \mathbf{v}^T \mathbf{v} = \lambda$$
.

Dimensionality Reduction with PCA

э

14/25

- 4 同 6 4 日 6 4 日 6

Total Variance of a Subspace

Consider a *k*-subspace  $\Sigma$  with a basis  $\{u_1, ..., u_k\}$ . We define the total variance of  $\Sigma$  as

$$\sum_{i=1}^k \text{ variance of } \boldsymbol{u}_i.$$

15/25

(日)

Total Variance vs. Construction Error

**Lemma 5:** Consider a *k*-subspace  $\Sigma$  with a basis  $\{u_1, ..., u_k\}$ . It holds that

$$\frac{1}{n}\sum_{\boldsymbol{p}\in P}|\boldsymbol{p}|^2 = \text{total variance of }\Sigma + \text{ construction error of }\Sigma.$$

**Proof:** Follows immediately from Lemma 1 and the definitions of total variance (Slide 15) and construction error (Slide 7).

**Implication:** We need to maximize total variance to minimize construction error.

Dimensionality Reduction with PCA

-

16/25

・ 同 ト ・ ヨ ト ・ ヨ ト

The rest of the slides will prove:

**Theorem 2:** The *k*-subspace with  $\{v_1, ..., v_k\}$  as the basis has the largest total variance among all *k*-subspaces.

Theorem 1 is then an immediate corollary of Theorem 2 and Lemma 5.

Dimensionality Reduction with PCA

17/25

**Lemma 6:** Among all unit vectors, the vector  $\mathbf{v}_1$  (the unit eigenvector for the largest eigenvalue  $\lambda_1$ ) has the largest variance.

**Proof:** By Lemma 3, finding a unit vector of the largest variance can be modeled as the following optimization problem:

find a vector  $\boldsymbol{w}$  to maximize  $\boldsymbol{w}^T \boldsymbol{A} \boldsymbol{w}$  subject to  $\boldsymbol{w}^T \boldsymbol{w} = 1$ .

We can solve the problem using the method of Lagrange multipliers. Introduce a real value  $\lambda$ , and define the Lagrangian function:

$$f(\boldsymbol{w}, \lambda) = \boldsymbol{w}^{T} \boldsymbol{A} \boldsymbol{w} - \lambda (\boldsymbol{w}^{T} \boldsymbol{w} - 1)$$

which yields

$$\frac{\partial f}{\partial \boldsymbol{w}} = 2\boldsymbol{A}\boldsymbol{w} - 2\lambda\boldsymbol{w}.$$

Setting  $\frac{\partial f}{\partial w} = 0$  gives a condition that an optimal w must satisfy:

$$\mathbf{A}\mathbf{w} = \lambda \mathbf{w}.$$

In other words,  $\lambda$  is an eigenvalue and **w** is a unit eigenvector of  $\lambda$ .

Lemma 4 shows that the variance of  $\boldsymbol{w}$  is exactly  $\lambda$ . Thus, the variance of  $\boldsymbol{w}$  cannot exceed the maximum eigenvalue  $\lambda_1$  of  $\boldsymbol{A}$ . Setting  $\boldsymbol{w}$  to  $\boldsymbol{v}_1$  achieves this variance.

19/25

< ロ > < 同 > < 回 > < 回 >

**Lemma 7:** Among all the unit vectors orthogoal to  $v_1$ , the vector  $v_2$  (the unit eigenvector for the second largest eigenvalue  $\lambda_2$ ) has the largest variance.

**Proof:** By Lemma 3, the goal is to solve the following optimization problem:

find a vector  $\boldsymbol{w}$  to maximize  $\boldsymbol{w}^T \boldsymbol{A} \boldsymbol{w}$  subj. to  $\boldsymbol{w}^T \boldsymbol{w} = 1$  and  $\boldsymbol{w}^T \boldsymbol{v}_1 = 0$ .

Again, we can solve the problem using the method of Lagrange multipliers. Introduce real values  $\lambda$  and  $\phi$ , and define the Lagrangian function:

$$f(\boldsymbol{w}, \lambda, \phi) = \boldsymbol{w}^T \mathbf{A} \boldsymbol{w} - \lambda (\boldsymbol{w}^T \boldsymbol{w} - 1) - \phi \boldsymbol{w}^T \boldsymbol{v}_1$$

which yields

$$\frac{\partial f}{\partial \boldsymbol{w}} = 2\boldsymbol{A}\boldsymbol{w} - 2\lambda\boldsymbol{w} - \phi\boldsymbol{v}_1.$$

The optimal **w** needs to satisfy  $\frac{\partial f}{\partial w} = 0$ , namely:

$$2\mathbf{A}\boldsymbol{w} - 2\lambda\boldsymbol{w} - \phi\boldsymbol{v}_1 = 0. \tag{1}$$

Next, we argue that  $\phi$  must be 0. The above equation leads to:

$$2\boldsymbol{v}_1^T \mathbf{A} \boldsymbol{w} - 2\lambda \boldsymbol{v}_1^T \boldsymbol{w} - \phi \boldsymbol{v}_1^T \boldsymbol{v}_1 = 0.$$
 (2)

Recall that  $\mathbf{v}_1^T \mathbf{w} = 0$  and  $\mathbf{v}_1^T \mathbf{v}_1 = 1$ . Furthermore:

$$\boldsymbol{\nu}_1^T \boldsymbol{A} \boldsymbol{w} = \boldsymbol{w}^T \boldsymbol{A}^T \boldsymbol{\nu}_1 \quad \text{(transposing a } 1 \times 1 \text{ matrix})$$
  
=  $\boldsymbol{w}^T \boldsymbol{A} \boldsymbol{\nu}_1 = \boldsymbol{w}^T (\boldsymbol{A} \boldsymbol{\nu}_1) = \lambda_1 \boldsymbol{w}^T \boldsymbol{\nu}_1 = 0.$ 

Hence, from (2), we get  $\phi = 0$ . Thus, (1) now becomes:

$$2\mathbf{A}\boldsymbol{w}-2\lambda\boldsymbol{w} = 0$$

namely, the optimal  $\boldsymbol{w}$  must also be an eigenvector.

It thus follows from Lemma 4 that the optimal  $\boldsymbol{w}$  must be  $\boldsymbol{v}_2$ .

The following theorem generalizes Lemma 7.

**Theorem 3:** For any  $t \in [2, d]$ , the vector  $\mathbf{v}_t$  (the unit eigenvector for the *t*-th largest eigenvalue  $\lambda_t$ ) has the largest variance among all the unit vectors orthogonal to  $\mathbf{v}_1, ..., \mathbf{v}_{t-1}$ .

The proof is a straightforward extension of the proof of Lemma 7 and is left to you as an exercise.

We are now ready to prove Theorem 2 using mathematical induction.

**Base case.** For k = 1, the theorem's correctness follows directly from Lemma 6.

**Inductive case.** Assuming the theorem's correctness on k = t - 1 for some  $t \in [2, d]$ , we will show that it also holds on k = t.

Let  $\Sigma^*$  be the *t*-subspace that has  $\{v_1, ..., v_t\}$  as a basis. Let  $\Sigma$  be an arbitrary *t*-subspace.

Our goal is to show that the total variance of  $\Sigma$  is at most that of  $\Sigma^*$ .

23/25

イロト イポト イラト イラト

**Lemma 8:** We can find a vector  $\boldsymbol{w}$  in the *t*-subspace  $\boldsymbol{\Sigma}$  such that  $\boldsymbol{w}$  is orthogoal to  $\boldsymbol{v}_i$  for every  $i \in [1, t-1]$ .

**Proof:** Let  $\{u_1, ..., u_t\}$  be an arbitrary basis of  $\Sigma$ . Any vector w in  $\Sigma$  is in the span of  $u_1, ..., u_t$ . This means

$$\boldsymbol{w} = \Sigma_{i=1}^t c_i \boldsymbol{u}_i$$

for some real values  $c_1, c_2, ..., c_k$ . As **w** needs to be orthogonal to each of  $v_1, ..., v_{t-1}$ , we have:

This yields a set of t - 1 linear equations for the t variables  $c_1, ..., c_t$ . The linear system has infinitely many solutions.

Lemma 8 implies that  $\Sigma$  has a basis  $\{w_1, ..., w_t\}$  such that  $w_t$  is orthogonal to all of  $v_1, ..., v_{t-1}$  (think: why?). We can then assert by Theorem 3 that

variance of  $\boldsymbol{v}_t \geq$  variance of  $\boldsymbol{w}_t$ .

Let  $\sum_{t=1}^{*}$  be the subspace defined by the basis  $\{v_1, ..., v_{t-1}\}$ . Let  $\sum_{t=1}^{*}$  be the subspace defined by the basis  $\{w_1, ..., w_{t-1}\}$ .

By the inductive assumption, we have

total variance of  $\Sigma_{t-1}^* \ge$  total variance of  $\Sigma_{t-1}$ .

Therefore

total variance of  $\Sigma^*$  = total variance of  $\Sigma^*_{t-1}$  + variance of  $\boldsymbol{v}_t$   $\geq$  total variance of  $\Sigma_{t-1}$  + variance of  $\boldsymbol{v}_t$   $\geq$  total variance of  $\Sigma_{t-1}$  + variance of  $\boldsymbol{w}_t$ = total variance of  $\Sigma$ .

This completes the proof of Theorem 2.

25/25