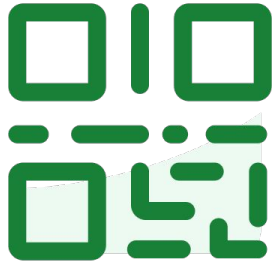




slido



Join at slido.com
#1565862



Click **Present with Slido** or install our [Chrome extension](#) to display joining instructions for participants while presenting.



LECTURE 25

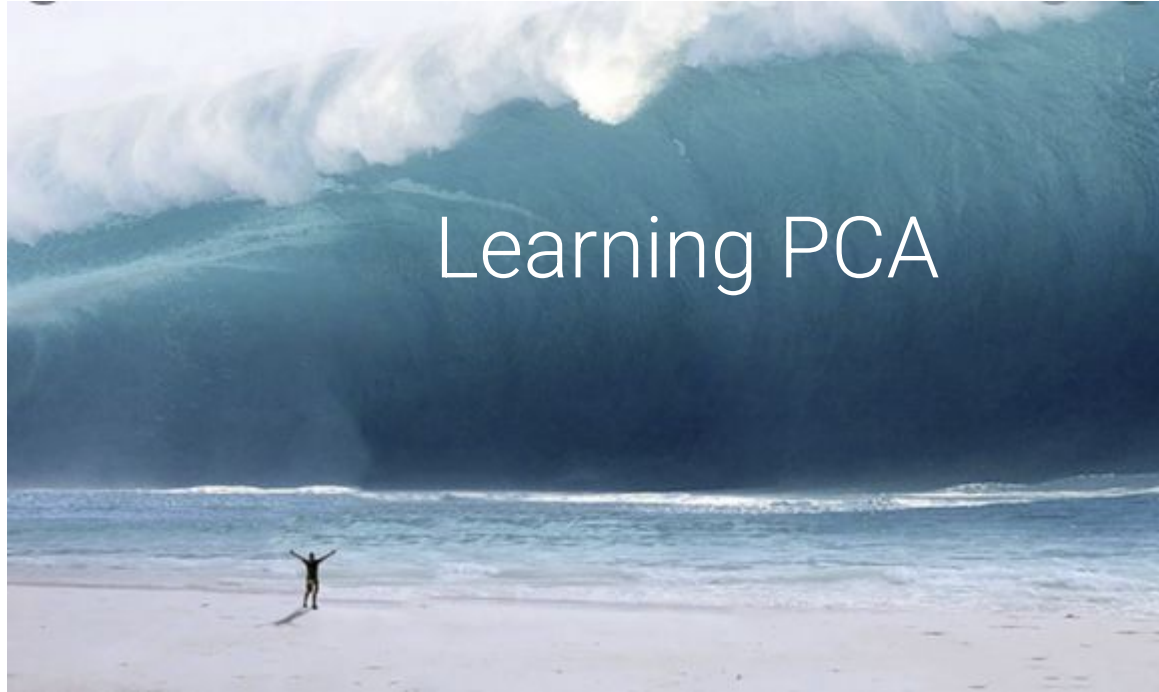
Principal Component Analysis II

PCA: A technique for EDA and feature generation.

Data 100/Data 200, Spring 2025 @ UC Berkeley

Narges Norouzi and Josh Grossman

Content credit: [Acknowledgments](#)



Learning PCA

Be kind to yourself 💙



PCA as Loss Minimization

Lecture 25, Data 100 Spring 2025

PCA as Loss Minimization

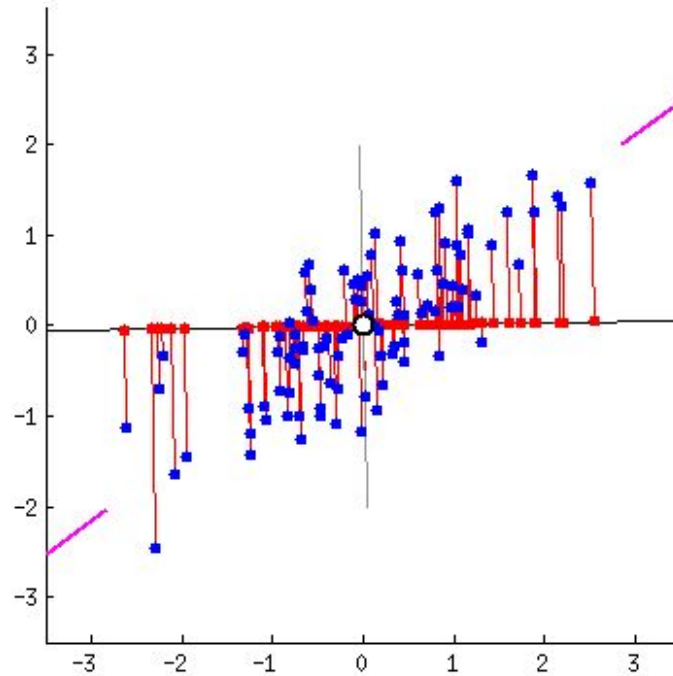
Singular Value Decomposition

PCA with SVD

Centering Data and Computing Variance

Useful reference for this lecture:

stats.stackexchange.com/questions/134282/relationship-between-svd-and-pca-how-to-use-svd-to-perform-pca

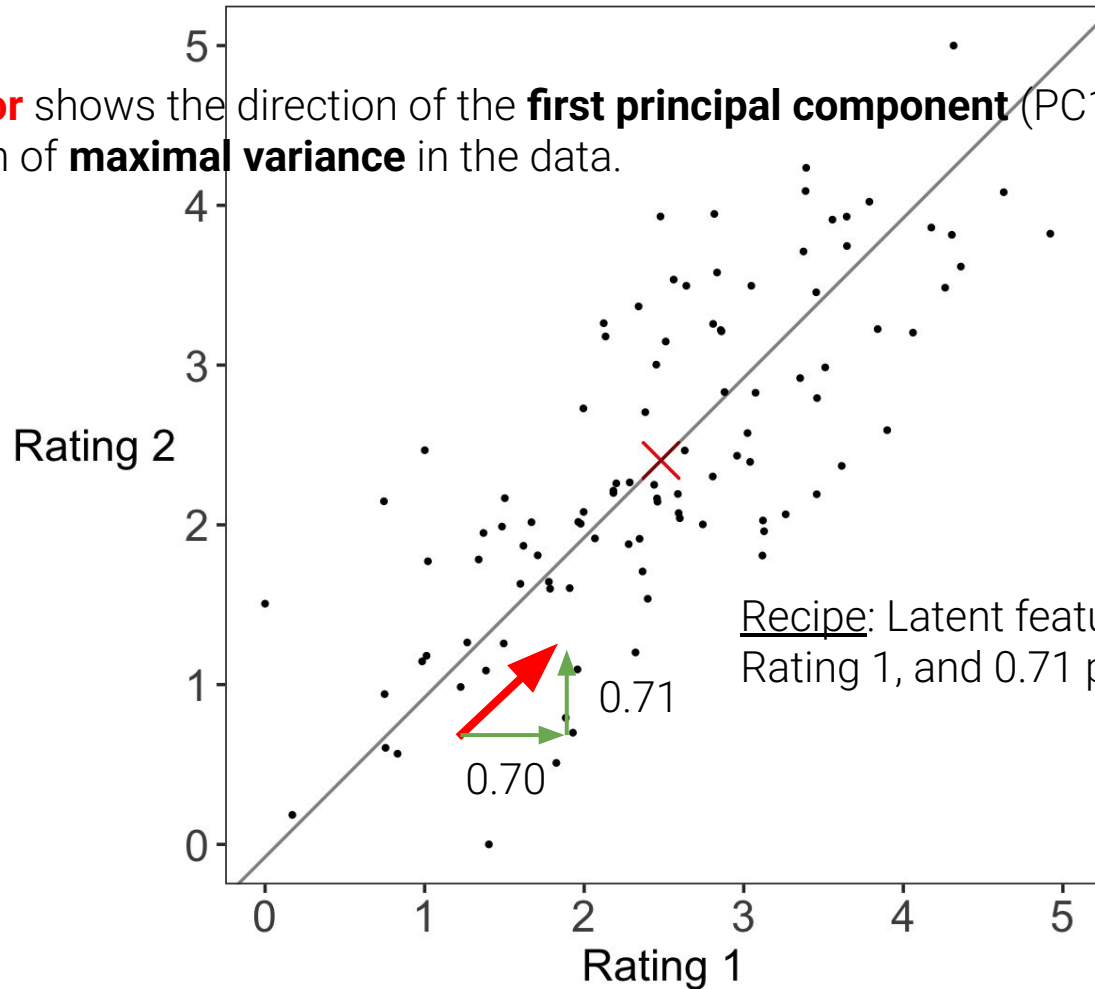


Maximizing variance = **Spreading out red dots**

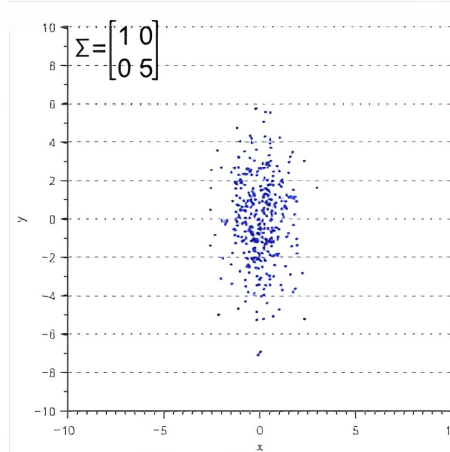
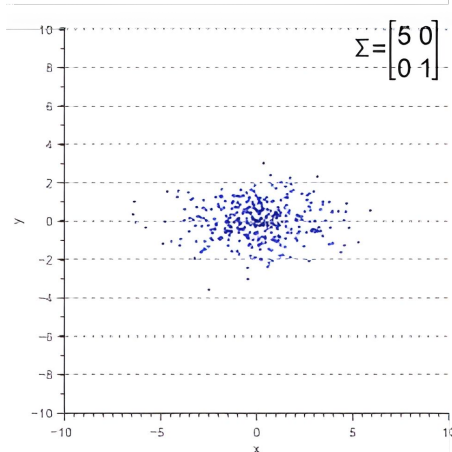
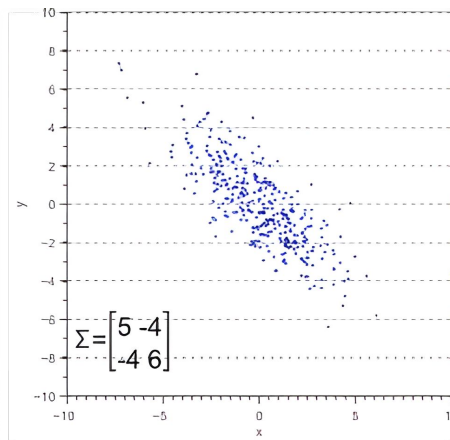
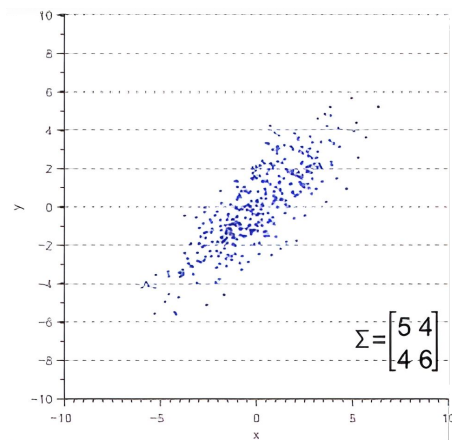
Equivalent: Minimize sum of squared **perpendicular** distances from points to projected point



This **length 1 vector** shows the direction of the **first principal component (PC1)**.
This is the direction of **maximal variance** in the data.



Recipe: Latent feature 1 is 0.70 parts
Rating 1, and 0.71 parts Rating 2.



Diagonal elements:
Variance of data along each dimension.

Off-diagonal elements:
Covariance of different dimensions. Think "unscaled" correlation!



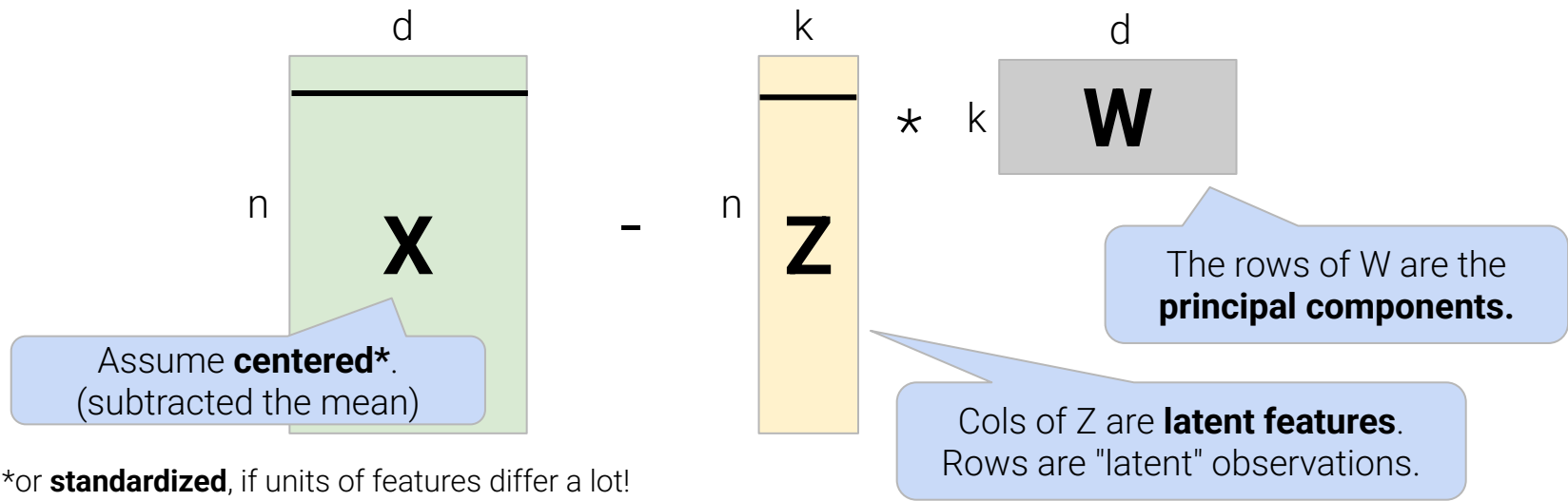
Derive PCA using Loss Minimization

Goal: Minimize the **reconstruction loss** of our **matrix factorization model**:

$$L(Z, W) = \frac{1}{n} \sum_{i=1}^n \|X_i - Z_i W\|^2$$

Diagram illustrating the matrix factorization model components:

- X_i (Row Vector): Dimensions $1 \times d$
- Z_i (Row Vector): Dimensions $1 \times k$
- W (Matrix): Dimensions $k \times d$

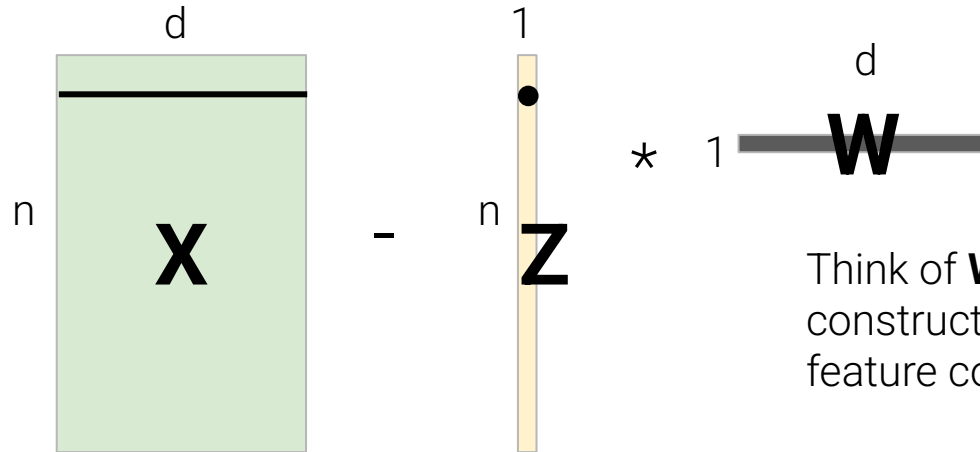




Simplified Derivation: consider (k=1)

Let consider the situation when $k=1$ (i.e., we construct only **one latent feature**):

$$L(z, w) = \frac{1}{n} \sum_{i=1}^n (X_i - z_i w) (X_i - z_i w)^T$$



Think of \mathbf{w} as the recipe for constructing \mathbf{z} from the \mathbf{d} feature cols in \mathbf{X} .

Simplified Derivation: Optimizing for w

Minimize the loss with respect to w :

$$L(w) = -w \Sigma w^T$$

i.e., make w length 1

Make **w really big** (to infinity) ... but we have the **orthonormality constraint $ww^T = 1$**

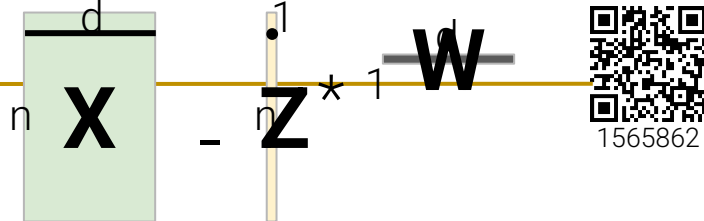
Use [Lagrange multiplier \$\lambda\$](#) to introduce the constraint **$ww^T = 1$** to our optimization problem:

$$L(w, \lambda) = -w \Sigma w^T + \lambda (ww^T - 1)$$

Intuition for Lagrange multiplier: When we take the partial derivative with respect to λ and set equal to 0 to minimize L , we will recover the constraint:

$$\frac{\partial}{\partial \lambda} L(w, \lambda) = ww^T - 1 = 0$$

See skipped slides of [Lecture 24](#) to see how the covariance matrix (Σ) shows up!¹⁰





$$L(w, \lambda) = -w\Sigma w^T + \lambda (ww^T - 1)$$

Take **derivative with respect to w** (vector calculus – out of scope for Data 100):

$$\frac{\partial}{\partial w} (-w\Sigma w^T + \lambda (ww^T - 1)) = -2\Sigma w^T + 2\lambda w^T$$

Think of ww^T like squaring \rightarrow Derivative is $2w^T$

Setting equal to zero: $-2\Sigma w^T + 2\lambda w^T = 0$

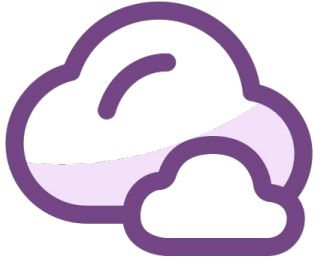
$$\Sigma w^T = \lambda w^T$$

What is this?

Remember that λ is a scalar!



slido



What is this?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.



$$L(w, \lambda) = -w\Sigma w^T + \lambda (ww^T - 1)$$

Eigenvalue of Σ

$$\Sigma w^T = \lambda w^T$$

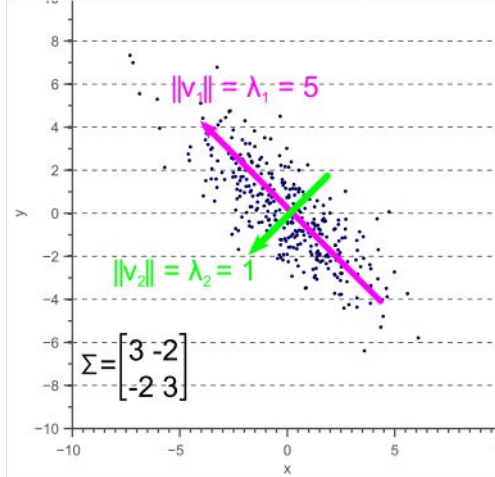
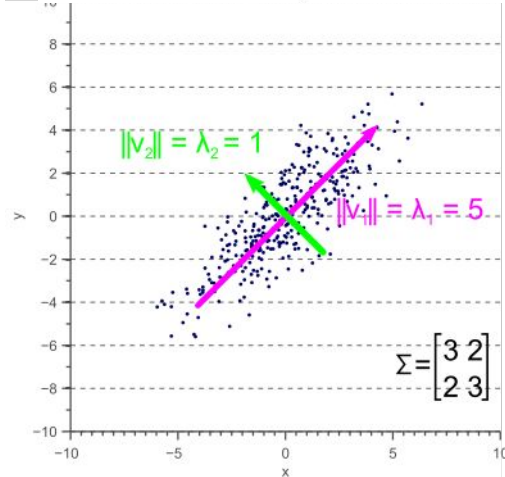
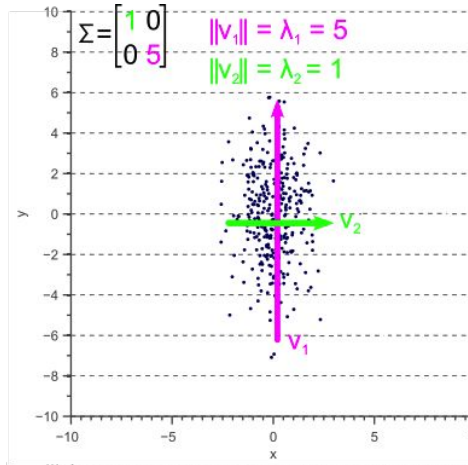
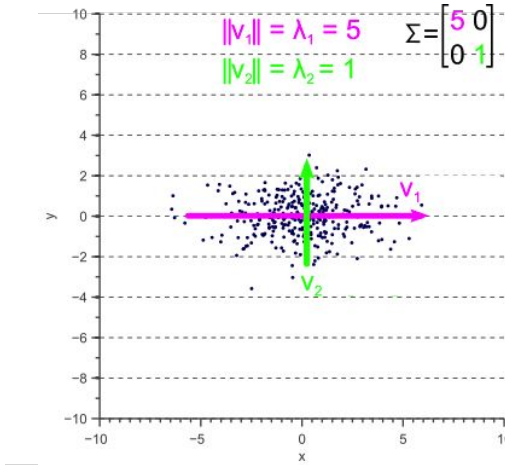
Eigenvector of Σ

w is a **unit** (i.e., length 1) **eigenvector** of the **covariance matrix**. When we multiply a matrix by one of its eigenvectors, it's equivalent to multiplying the eigenvector by its (scalar) eigenvalue.

$L(w, \lambda)$ is **minimized** when $w\Sigma w^T = w\lambda w^T = \lambda ww^T = \lambda$ is **big**. So, the **optimal w** is the eigenvector with the **largest eigenvalue λ** .

In other words, the **optimal w** points in the direction of the **greatest variance** of the data (PC1!), and λ is the variance in that direction (so long as X is centered!).

Eigenvectors and eigenvalues of covariance matrix



Same data rotated four ways.

$v_1/5$ is PC1. **Unit-length** eigenvector with largest eigenvalue points in direction of maximum variance. Eigenvalue ($\lambda_1=5$) is variance.

$v_2/1$ is PC2. Unit-length eigenvector with second largest eigenvalue points in direction of maximum variance that is **orthogonal** to first eigenvector.

visiondummy.com/2014/04/geometric-interpretation-covariance-matrix



Singular Value Decomposition

Lecture 25, Data 100 Spring 2025

PCA as Loss Minimization

Singular Value Decomposition

PCA with SVD

Centering Data and Computing Variance

Useful reference for this lecture:

stats.stackexchange.com/questions/134282/relationship-between-svd-and-pca-how-to-use-svd-to-perform-pca



Singular value decomposition (SVD) decomposes a matrix into a product of three matrices.

- We assume you have taken (or are taking) a linear algebra course.
- We will not explain SVD in its entirety—in particular, we will not prove:
 - **How** the SVD is computed
 - **Why** SVD is a valid decomposition of rectangular matrices

Today, we will only cover **what** is needed for PCA.

Checkout the EE16 notes for a deeper study of SVD:

<https://eecs16b.org/notes/sp24/note15.pdf>



Singular value decomposition (SVD) describes a matrix decomposition into three matrices:

$$X = U S V^T$$

$$X \in \mathbb{R}^{n \times d}$$

rank $r \leq d$

Assume $d < n$.

$$U \in \mathbb{R}^{n \times d}$$

columns of U are
orthonormal (i.e., unit
length and orthogonal)

Columns of U are
eigenvectors of XX^T

$$S$$

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

diagonal matrix
of **singular values**,
ordered from
largest to smallest.
 r positive singular
values. Remaining **$d-r$**
singular values are 0.

$$V^T$$

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

columns+rows of V
are **orthonormal**

Columns of V
(rows of V^T) are
eigenvectors of $X^T X$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n X_i^T X_i$$

There are technically infinite possible factorizations, but if singular values are distinct + in decreasing order then solution space is constrained to shuffling columns of U and V .

Looks like covariance matrix without $1/n!$ 17



$$U \in \mathbb{R}^{n \times d}$$

- Columns of U are **orthonormal**: $\vec{u}_i^\top \vec{u}_j = 0$ for all i, j and all vectors \vec{u}_i are unit vectors
- Columns of U are called the **left singular vectors**
- $UU^T = I_n$ and $U^T U = I_d$
- Can think of U as a rotation

Eigenvectors of XX^T :

$$\begin{aligned} XX^T &= USV^T (USV^T)^T = USV^T V S^T U^T \\ &= USS^T U^T \end{aligned}$$

Right multiply by U : $\Rightarrow XX^T U = US^T S U^T U$

$$= USS^T$$

SVD One-by-One: S

S

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

- Diagonal values (**singular values**), are ordered from **greatest to least**
- r **non-zero** singular values

$$\underbrace{\left[\begin{array}{ccccccc} s_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & s_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_r & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{array} \right]}_d$$

- r is the rank of X
- The majority of the matrix is zero
- The **singular values** s_1, s_2, \dots, s_r are **non-negative** and **sorted** $s_1 \geq s_2 \geq \dots \geq s_r$
- Can think of S as a scaling operation



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

- Columns of V are orthonormal \rightarrow rows of V^T are orthonormal
- Rows of V are also orthonormal
- Columns of V are called the **right singular vectors**
- $VV^T = V^TV = I_d$
- Can think of V as a rotation

Eigenvectors of X^TX :

$$X^TX = (USV^T)^T USV^T = VS^T U^T USV^T$$

$$X^TX = VS^T SV^T$$

Right multiply by V : $\Rightarrow X^TXV = VS^T SV^TV$

$$= VS^T S$$



1565862

NumPy SVD

```
U, S, Vt = np.linalg.svd(X, full_matrices = False)
```

[\[documentation\]](#)

"economy" or "compact" SVD

$$U \in \mathbb{R}^{n \times d}$$

$$X = U S V^T$$

width	height	area	Perim.
2.97	1.35	24.78	8.64
-3.03	-0.65	-15.22	-7.36
-4.03	-1.65	-20.22	-11.36
3.97	-1.65	3.78	4.64
3.97	3.35	48.78	14.64
-2.03	-3.65	-20.22	-11.36
-1.03	-2.65	-15.22	-7.36
0.97	0.35	6.78	2.64
1.97	-3.65	-16.22	-3.36
2.97	-2.65	-7.22	0.64
...

 $n \times d$

-0.13	0.01	0.03	-0.21
0.09	-0.08	0.01	0.56
0.12	-0.13	0.09	-0.07
-0.03	0.18	0.01	-0.05
-0.26	-0.09	0.09	-0.06
0.12	-0.05	0.17	-0.05
0.09	0	0.1	-0.08
-0.04	0.01	0	-0.08
0.08	0.18	0.04	-0.05
0.03	0.19	0.02	-0.05
...

 $n \times d$

197.39			
	27.43		
		23.26	
			0

 $d \times d$

-0.1	-0.07	-0.93	-0.34
0.67	-0.37	-0.26	0.59
0.31	-0.64	0.26	-0.65
0.67	0.67	0	-0.33

 $d \times d$ X is therefore **rank 3**.



1. Center the data matrix by subtracting the mean of each attribute column.
2. To find \mathbf{v}_i , the i -th **principal component**:
 - \mathbf{v} is a **unit vector** that linearly combines the attributes.
 - \mathbf{v} gives a one-dimensional projection of the data.
 - \mathbf{v} is chosen to minimize the sum of squared distances between each point and its projection onto \mathbf{v} .
 - Choose \mathbf{v} such that it is orthogonal to all previous principal components.

Let's now use SVD to compute the **principal components**.



PCA with SVD

Lecture 25, Data 100 Spring 2025

PCA as Loss Minimization

Singular Value Decomposition

PCA with SVD

Centering Data and Computing Variance



Assume we have constructed the Singular Value Decomposition (SVD) of X :

$$X = USV^T$$

If X is **centered**, the **covariance matrix** (Σ) of X is:



Assume we have constructed the Singular Value Decomposition (SVD) of X :

$$X = USV^T$$

If X is **centered**, the **covariance matrix** (Σ) of X is:

$$\Sigma = (1/n) X^T X$$

$$= (1/n) (USV^T)^T USV^T$$

Substitution.

$$= (1/n) V S^T U^T U S V^T$$

Matrix transpose rule. $(AB)^T = B^T A^T$

$$= (1/n) V S^T S V^T$$

U is orthonormal. $U^T U = \mathbf{I}$

$$= (1/n) V S^2 V^T$$

S is diagonal. Square diagonal elements.



Assume we have constructed the Singular Value Decomposition (SVD) of X :

$$X = USV^T$$

If X is **centered**, the **covariance matrix** (Σ) of X is:



Assume we have constructed the Singular Value Decomposition (SVD) of X :

$$X = USV^T$$

If X is **centered**, the **covariance matrix** (Σ) of X is:

$$\Sigma = (1/n) VS^2V^T$$

Where we left off.

$$\Sigma V = (1/n) VS^2V^TV$$

Right multiply by V .

$$\Sigma V = \frac{S^2}{n} V$$

V is orthonormal. $V^TV = \mathbf{I}$. S^2 can be moved.

$$\Sigma w^T = \lambda w^T$$

Looks like our loss minimization solution! 



1565862

Principal Components (PCs) are the Eigenvectors of the Covariance Matrix

Assume we have constructed the Singular Value Decomposition (SVD) of X :

$$X = USV^T$$

$$\Sigma V = \frac{S^2}{n} V$$

The **columns of V** (rows of V^T) are the **eigenvectors** of the **covariance matrix** (Σ) and therefore the **PCs**. Orthogonal **directions** of greatest variance of the data.

The **(singular values)² / n** are the **eigenvalues** of Σ . **Variance** of the data in each of the directions given by the PCs.

$$\Sigma w^T = \lambda w^T$$

← Same as definitions of corresponding terms ([earlier slide](#))



1565862

Computing Latent Vectors Using $X * V$

Constructing a 2 PC approximation ($k=2$). PCs as a "recipe" to make Z from features in X .

$$X * V = Z$$

width	height	area	Perim.
2.97	1.35	24.78	8.64
-3.03	-0.65	-15.22	-7.36
-4.03	-1.65	-20.22	-11.36
3.97	-1.65	3.78	4.64
3.97	3.35	48.78	14.64
-2.03	-3.65	-20.22	-11.36
...

PC1	PC2		
-0.1	0.67	0.31	0.67
-0.07	-0.37	-0.64	0.67
-0.93	-0.26	0.26	0
-0.34	0.59	-0.65	-0.33

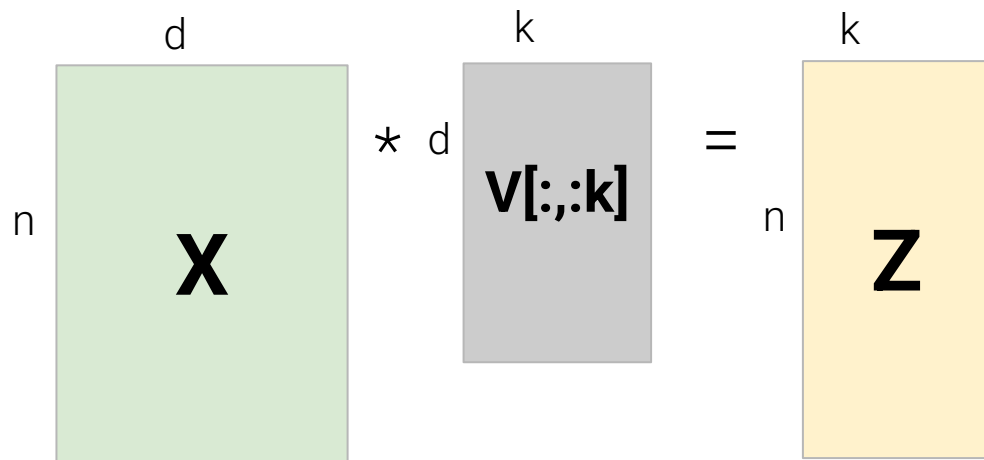
-26.43	0.16
17.05	-2.18
23.25	-3.54
-5.38	5.03
-51.09	-2.59
23.19	-1.45
...	...



We have now shown that if we construct the singular value decomposition of X :

$$X = USV^T$$

The **first k** columns of V (rows of V^T) are the **first k PCs**. We construct the **latent features (Z)** of X by projecting X onto the **k PCs**. We choose k ! Often, $k=2$.



There is another way to compute Z :

PC "Recipe" of features

$$Z = XV = USV^T V$$

$$= US$$

U contains normalized columns of Z . Scale up by S .



1565862

Computing Latent Vectors Using $U * S$

Equivalent construction of Z ! **U cols are normalized Z cols.** S "scales up" U cols to Z cols.

$$U * S = Z$$

-0.13	0.01	0.03	-0.21
0.09	-0.08	0.01	0.56
0.12	-0.13	0.09	-0.07
-0.03	0.18	0.01	-0.05
-0.26	-0.09	0.09	-0.06
0.12	-0.05	0.17	-0.05
...

197.39			
	27.43		
		23.26	
			0

-26.43	0.16
17.05	-2.18
23.25	-3.54
-5.38	5.03
-51.09	-2.59
23.19	-1.45
...	...



Given the entire Z matrix we can recover the **centered X** by **multiplying by V^T** :

$$ZV^T = XVV^T = USV^T = X$$

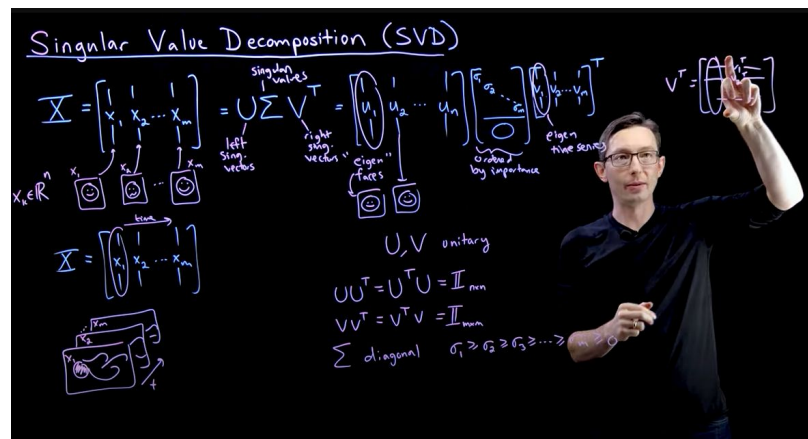
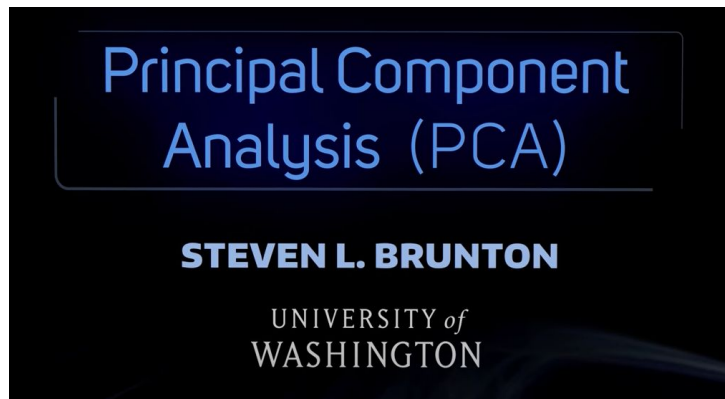
This is like **inverting** our PC recipe \rightarrow How do we combine our latent features (Z) to get back our original features (X)?

If you choose **$k < r$** , where r is $\text{rank}(X)$, you will only recover X **approximately**.

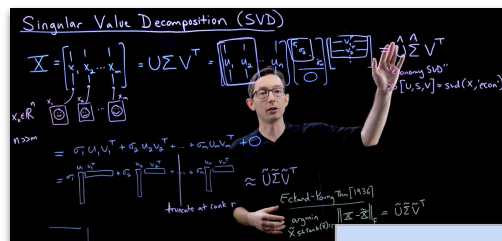
To recover the original (uncentered) X we would also need to add back the mean.



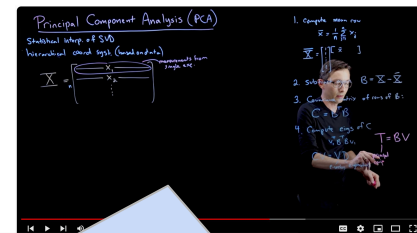
Singular Value Decomposition Playlist



S. Brunton's video on SVD



S. Brunton's video on PCA



- Extraordinarily clear.
- Videos cover the foundations, applications, linear algebra, stats and code.
- A masterpiece of teaching.
- Note: our d is his m (# of features).

Unfortunately, his definition of the **Principal Components** is wrong (~8:30).



Demo

Computing

PCA using SVD



Centering Data and Computing Variance

Lecture 25, Data 100 Spring 2025

PCA as Loss Minimization

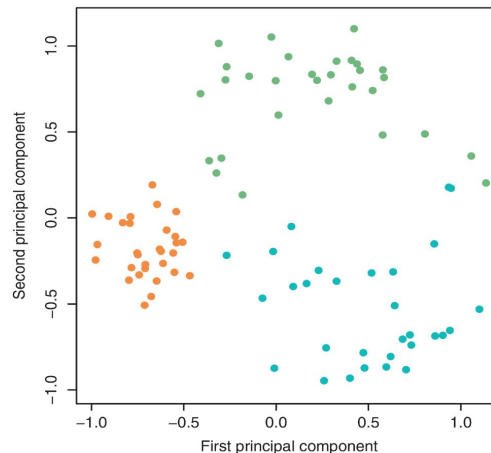
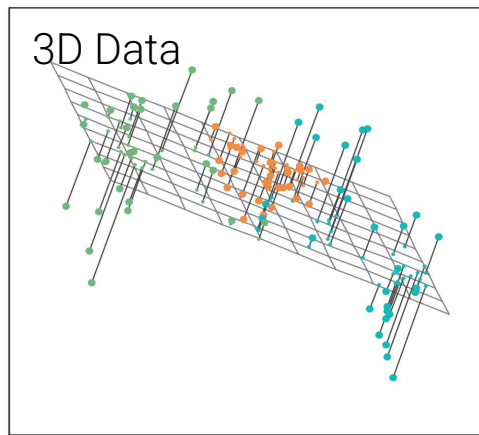
Singular Value Decomposition

PCA with SVD

Centering Data and Computing Variance

We often construct a scatter plot of the data projected onto the **first two principal components**. This is often called a **PCA plot**.

- PCA plots allow us to visually assess similarities between our data points and if there are any clusters in our dataset.



PCA Plot

If PC1+PC2 explain a large % of the **total variance**, then the PCA plot is a good representation.
If not, then a **PCA plot** is omitting lots of information.



1565862

We define the **total variance** of a data matrix as the sum of variances of attributes.

width	length	area	perimeter
20	20	400	80
16	12	192	56
...
24	12	288	72

Total Variance: **402.56** = 7.69 5.35 50.79 338.73



1565862

Variance and Singular Values

We define the **total variance** of a data matrix as the sum of variances of attributes.

$$\Sigma V = \frac{S^2}{n} V$$

Total Variance: **402.56** =

width	length	area	perimeter
20	20	400	80
16	12	192	56
...
24	12	288	72
7.69	5.35	50.79	338.73

The **component score** tells us the **variance** captured by the i^{th} principal component. The component scores are the **eigenvalues of the covariance matrix**. N is # of datapoints.

$$i^{\text{th}} \text{ component score} = \frac{(i^{\text{th}} \text{ singular value})^2}{N}$$

Variance captured by **PC1**

197.4	0	0	0
0	27.43	0	0
0	0	23.26	0
0	0	0	0

$$\rightarrow 197.39^2/100 = \mathbf{389.63}$$

$$\rightarrow 27.43^2/100 = 7.52$$

$$\rightarrow 23.26^2/100 = 5.41$$

Sum = **402.56**.

How do we compute an array of **variance ratios**, where each element is the **fraction** that each PC contributes to total data variance?

$$X = USV^T$$

```
u, s, vt = np.linalg.svd(X, full_matrices = False)
```

- A. `s / n` # `n` is `len(X)`, num features
- B. `s ** 2 / n`
- C. `s / sum(s)`
- D. `s**2 / sum(s**2)`
- E. Something else

$$i\text{-th component score} = \text{Variance captured by } i\text{-th PC} = \frac{(i\text{-th singular value})^2}{n}$$



slido



How do we compute an array of variance ratios?

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.



$$i\text{-th component score} = \frac{\text{Variance captured by } i\text{-th PC}}{n} = \frac{(i\text{-th singular value})^2}{n}$$

$$\text{total variance} = \text{sum of all the component scores} = \sum_{i=1}^k \frac{s_i^2}{N}$$

$$X = USV^T$$

$$\text{variance ratio of principal component } j = \frac{\text{component score } j}{\text{total variance}} = \frac{s_j^2/N}{\sum_{i=1}^k s_i^2/N} = \frac{s_j^2}{\sum_{i=1}^k s_i^2}$$

How do we compute an array of **variance ratios**, where each element is the **fraction** that each principal component contributes to total data variance?

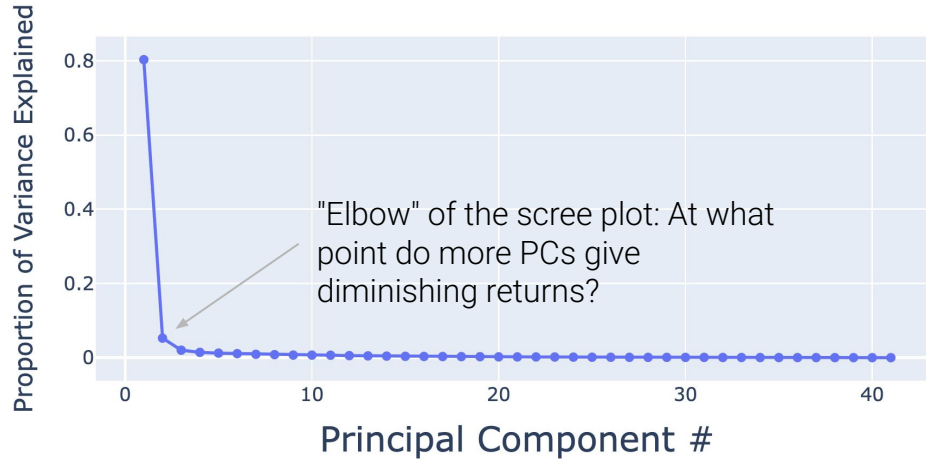
```
u, s, vt = np.linalg.svd(X, full_matrices = False)
```

- A. s / n # n is $\text{len}(X)$, num features
- B. $s ** 2 / n$
- C. $s / \text{sum}(s)$
- ☒ D. $s ** 2 / \text{sum}(s ** 2)$
- E. Something else

Scree Plot

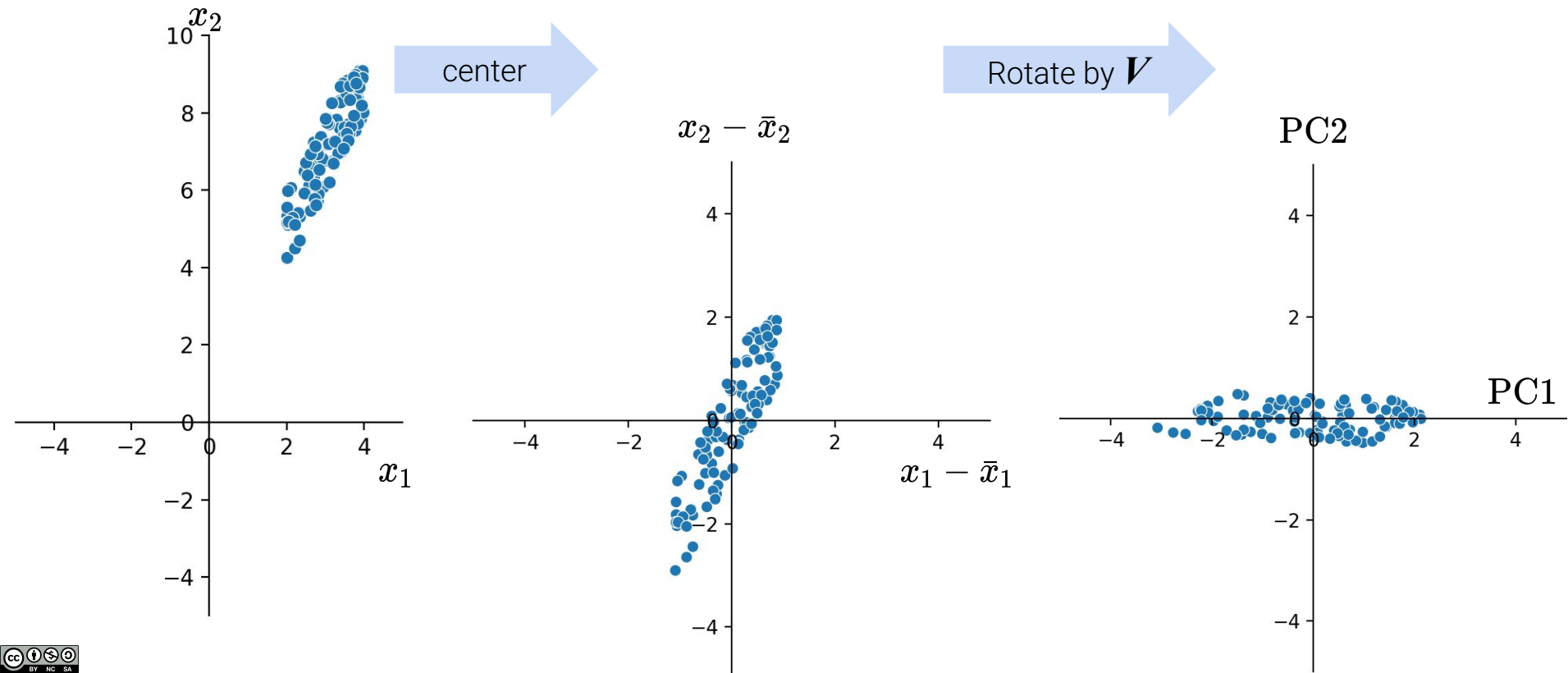
If PC1+PC2 explain a large % of the total variance, then the PCA plot is a good representation. If not, then a **PCA plot** is omitting lots of information.

A **scree plot** shows the variance ratio captured by each principal component, largest first.

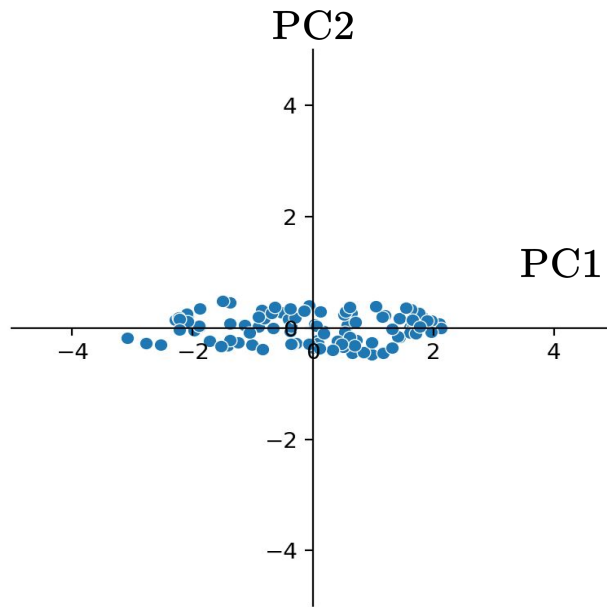


Scree [\[wikipedia\]](https://en.wikipedia.org/wiki/Scree)

PCA first **centers the data matrix**, then rotates it so that the direction with the **greatest variance** is aligned with the x-axis. [We saw this "changing frame of reference" last lecture.]



- Principal components are all **orthogonal** to each other
 - Why? Recall that the **columns of V are orthonormal!**
- Principal Components are **axis-aligned**
 - If we plot two PCs on a 2D plane, one will lie on the x-axis, the other on the y-axis
- The **latent factors** are obtained by **projecting X** onto the principal components





When running PCA it is important to **center the data** (ensure data is centered at 0):

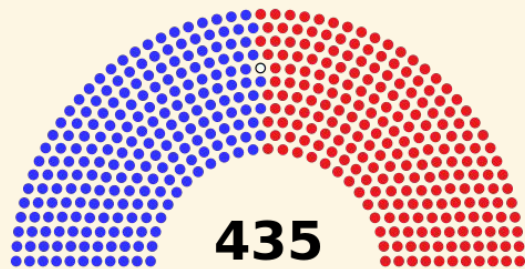
```
X = X - np.mean(X, axis = 0)
```

- The **sklearn** PCA code does this **automatically**!

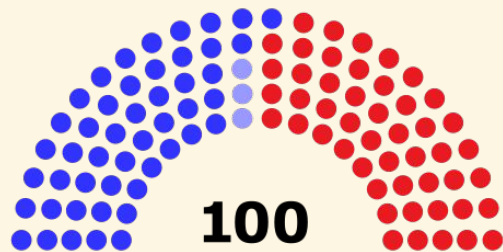
You should also *consider* **standardizing** your data:

```
X = (X - np.mean(X, axis = 0)) / np.std(X, axis = 0)
```

- The **sklearn** PCA code does not do this automatically
- Standardizing puts all cols on same scale (e.g., $X=0.5$ has same meaning across cols)
- You should **not** standardize your data if the units are **all the same** (e.g., all 0/1 cols)



House of representatives



Senate

Demo

Congressional Vote Data



1565862

Let's examine how the **House of Representatives** (of the 116th Congress, 1st session) voted in the month of **September 2019**.

Specifically, we'll look at the records of Roll call votes. From the U.S. Senate ([link](#)):

Roll call votes occur when a representative or senator votes "yea" or "nay," so that the names of members voting on each side are recorded.

Do legislators' roll call votes show a relationship with their political party?

$$\begin{array}{c} d \\ \text{---} \\ \boxed{\mathbf{X}} \\ n \end{array} * \begin{array}{c} k \\ \text{---} \\ \boxed{\mathbf{V}} \end{array} = \begin{array}{c} k \\ \boxed{\mathbf{Z}} \\ n \end{array}$$

The i -th row of \mathbf{V} indicates **how much feature i contributes** to each PC. Cols of \mathbf{V} are the PCs ("recipes").

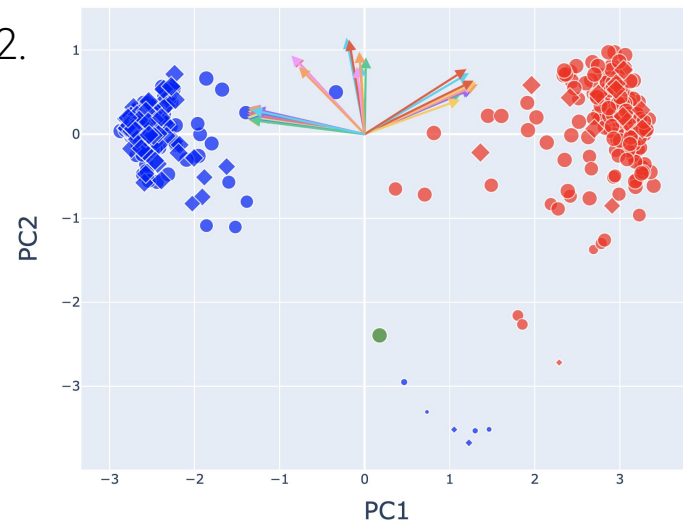
"First row of $\mathbf{V} = [0.9, -0.44] \rightarrow$ PC1 linear combo is 0.9 parts feature 1, and PC2 has -0.44 parts feature 1."

Biplots superimpose **feature influence** on plot of PC1 vs. PC2.

Biplots help us interpret how features influence the PCs: positively, negatively, or not much at all.

Simplest biplot: Plot the rows of \mathbf{V} with no scaling.

- For other scalings, which can lead to more interpretable directions/loadings, see [\[SAS biplots\]](#).





Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. Han Xiao, Kashif Rasul, Roland Vollgraf. arXiv:1708.07747

<https://github.com/zalandoresearch/fashion-mnist>



Demo

Time permitting!



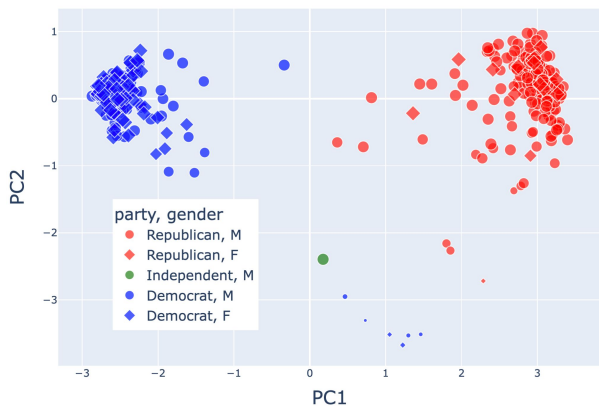
1565862

Summary: Plots Based on PCA

PCA Plot

Scatter plot of PC1 against PC2

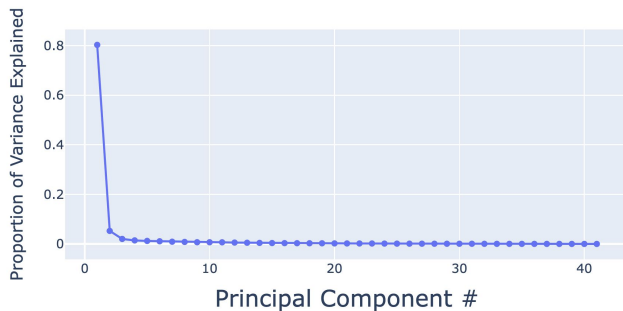
- Similarities of data points; identifying clusters.



Scree Plot

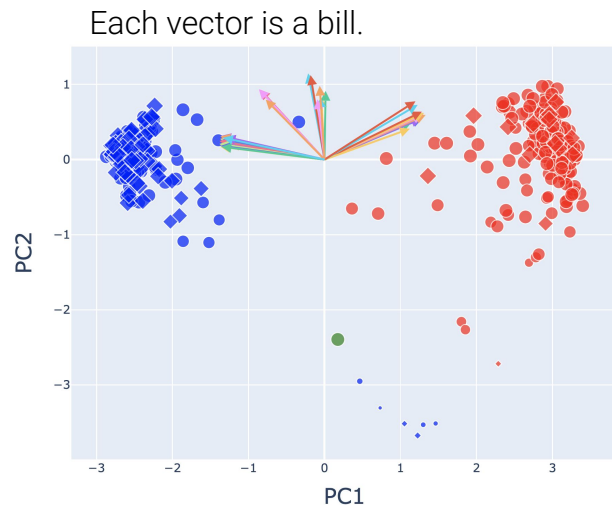
Variance ratio captured by each PC, with largest first.

- If first two ratios large, PCA plot is good representation
- “Elbow” method to assess how many PCs to use



Biplot

PCA plot + **influence of features** on PC1 and PC2.





LECTURE 25

Principal Component Analysis II

Data 100/Data 200, Spring 2025 @ UC Berkeley

Narges Norouzi and Josh Grossman

Content credit: [Acknowledgments](#)