ES114: Probability, Statistics, and Data Visualization

PCA, LLN and CLT 16th Apr 2025





• Recap





- Recap
 - Markov and Chebyshev's Inequality





- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors



- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors
 - ► Covariance Matrix



- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors
 - Covariance Matrix
 - Multivariate Normal



- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors
 - Covariance Matrix
 - Multivariate Normal
 - Gaussian Whitening



- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors
 - Covariance Matrix
 - Multivariate Normal
 - Gaussian Whitening
- Principal Component Analysis



- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors
 - Covariance Matrix
 - Multivariate Normal
 - Gaussian Whitening
- Principal Component Analysis
- Law of Large Numbers



- Recap
 - Markov and Chebyshev's Inequality
 - Random Vectors
 - Covariance Matrix
 - Multivariate Normal
 - Gaussian Whitening
- Principal Component Analysis
- Law of Large Numbers
- Central Limit Theorem

References



- PCA Explained
- PCA Vizualization
- WLLN and CLT
- WLLN and CLT 2

Iris Data





(attributes, measurements, dimensions)

Iris Data (1D)



Sepal length: X^1 (Random Variable)

Iris Data (1D) Sepal length: X^1 (Random Variable)





Iris Data (2D)



 X^1 : Sepal length, X^2 : Sepal width

Iris Data (2D)



 X^1 : Sepal length, X^2 : Sepal width



Iris Data (3D) X^1 : Sepal length, X^2 : Sepal width, X^3 : Petal length



Iris Data (3D) X^1 : Sepal length, X^2 : Sepal width, X^3 : Petal length

3D Scatter Plot of Iris Dataset (First Three Features)



Iris Data (4D)



	Sepal Length	Sepal Width	Petal Length	Petal Width
	X^1	X ²	<i>X</i> ³	X^4
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2
5	5.4	3.9	1.7	0.4

Iris Data (4D)



	Sepal Length	Sepal Width	Petal Length	Petal Width
	X^1	X ²	<i>X</i> ³	X^4
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2
5	5.4	3.9	1.7	0.4

How would you visualize it in 2D?



Images



Image Data as a Random Vector



X =

	\times'						
]]] []	0, 0, 0,	0, 0, 0,	0, 0, 0,	· · · , · · · ,	0, 21, 0,	0, 0, 0,	0], 0], 198],
 [[, 0, 0, 0,	0, 0, 0,	0, 0, 0,	····, ···,	85, 0, 0,	243, 0, 48,	252], 0], 242]]

Image Data as a Random Vector



]]]]	0, 0, 0,	0, 0, 0,	0, 0, 0,	••••, •••,	0, 21, 0,	0, 0, 0,	0], 0], 198],
] []	, 0, 0, 0,	0, 0, 0,	0, 0, 0,	••••,	85, 0, 0,	243, 0, 48,	252], 0], 242]]

• Are all the features useful?

Image Data as a Random Vector



]]]]	0, 0, 0,	0, 0, 0,	0, 0, 0,	····, ···,	0, 21, 0,	0, 0, 0,	0], 0], 198],
 [[, 0, 0, 0,	0, 0, 0,	0, 0, 0,	····, ····,	85, 0, 0,	243, 0, 48,	252], 0], 242]]

- Are all the features useful?
- How to capture important information?



X^1	X^2
0.41709128	5.0 € ° \
0.04221766	5.0 (- x2
0.38424179	5.0 ×3
0.90469106	5.0
0.60924091	5.0 -
0.58330889	5.0
0.56814491	5.0
0.68974537	5.0 (
0.23745621	5.0
0.70578727	5.0











"Manifold"







Projecting Data





Projecting Data





Projecting Data





 $||X \cdot v - \bar{X} \cdot v||^2 = ||(X - \bar{X}) \cdot v||^2$

Principal Component



Find the direction v which maximizes variance in the data?

Principal Component



Find the direction v which maximizes variance in the data?





Principal components are the directions which capture the maximum variance of the data



Principal components are the directions which capture the maximum variance of the data

• Data is *d*-dimensional



Principal components are the directions which capture the maximum variance of the data

- Data is *d*-dimensional
- Will need at most *d* directions to capture entire variance


Principal components are the directions which capture the maximum variance of the data

- Data is *d*-dimensional
- Will need at most d directions to capture entire variance
- Top k directions that preserve the maximum variance are the principal components



$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \\ x_{51} & x_{52} & x_{53} & x_{54} \end{bmatrix}$$



$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \\ x_{51} & x_{52} & x_{53} & x_{54} \end{bmatrix}$$

$$Cov(x_1, x_2) = \frac{1}{5} \sum_{i=1}^{5} (x_{i1} - \mu_1)(x_{i2} - \mu_2)$$



$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \\ x_{51} & x_{52} & x_{53} & x_{54} \end{bmatrix}$$

$$Cov(x_1, x_2) = \frac{1}{5} \sum_{i=1}^{5} (x_{i1} - \mu_1) (x_{i2} - \mu_2)$$
$$\frac{1}{N} (X - \mu)^T (X - \mu) = \frac{1}{N} \bar{X}^T \bar{X}$$
$$= \frac{1}{N} \begin{bmatrix} \bar{x}_1^T \bar{x}_1 & \bar{x}_1^T \bar{x}_2 & \bar{x}_1^T \bar{x}_3 & \bar{x}_1^T \bar{x}_4 \\ & \vdots \\ \bar{x}_5^T \bar{x}_1 & \bar{x}_5^T \bar{x}_2 & \bar{x}_5^T \bar{x}_3 & \bar{x}_5^T \bar{x}_4 \end{bmatrix}$$



Assume \bar{X} , $n \times d$ is mean normalized

• What can you say about $\bar{X}^T \bar{X}$?



- What can you say about $\bar{X}^T \bar{X}$?
 - **Square** $(d \times d)$



- What can you say about $\bar{X}^T \bar{X}$?

 - ► Square $(d \times d)$ ► Symmetric $(\bar{X}^T \bar{X})^T = \bar{X}^T (\bar{X}^T)^T = \bar{X}^T \bar{X}$



- What can you say about $\bar{X}^T \bar{X}$?
 - Square $(d \times d)$
 - Symmetric $(\bar{X}^T \bar{X})^T = \bar{X}^T (\bar{X}^T)^T = \bar{X}^T \bar{X}$
 - $\bar{X}^T \bar{X}$ is positive semi-definite so non-negative eigen values



- What can you say about $\bar{X}^T \bar{X}$?
 - Square $(d \times d)$
 - Symmetric $(\bar{X}^T\bar{X})^T = \bar{X}^T(\bar{X}^T)^T = \bar{X}^T\bar{X}$
 - ▶ $\bar{X}^T \bar{X}$ is positive semi-definite so non-negative eigen values
 - Since symmetric one can chose eigen vectors to be orthonormal

$$V = \left[\begin{array}{ccc} | & | & | \\ v_1 & v_2 & \dots & v_d \\ | & | & | \end{array} \right]$$



The direction which preserves the most variance is,

$$\bar{X}^T \bar{X} v_1 = \lambda_1 v_1$$



The direction which preserves the most variance is,

• Eigen vector of covariance matrix with largest eigen value

$$ar{X}^Tar{X}v_1=\lambda_1v_1$$

• Variance preserved is equal to the eigen value



The direction which preserves the most variance is,

$$\bar{X}^T \bar{X} v_1 = \lambda_1 v_1$$

- Variance preserved is equal to the eigen value
- $d \times d$ matrix can have atmost d eigen vectors



The direction which preserves the most variance is,

$$\bar{X}^T \bar{X} v_1 = \lambda_1 v_1$$

- Variance preserved is equal to the eigen value
- $d \times d$ matrix can have atmost d eigen vectors
- Total variance is $\sum_{i=1}^{d} \lambda_i$

The direction which preserves the most variance is,

$$\bar{X}^T \bar{X} v_1 = \lambda_1 v_1$$

- Variance preserved is equal to the eigen value
- $d \times d$ matrix can have atmost d eigen vectors
- Total variance is $\sum_{i=1}^{d} \lambda_i$
- Project data onto the top k eigen vectors yields the "most informative" k-dimensional data







- Data Compression
- Data Visualization
- Data Denoising



Algorithm Principal Component Analysis (PCA)

Require: Data matrix $X \in \mathbb{R}^{n \times d}$ with *n* samples and *d* features, number of components *k*

Ensure: Top k principal components



Algorithm Principal Component Analysis (PCA)

Require: Data matrix $X \in \mathbb{R}^{n \times d}$ with *n* samples and *d* features, number of components *k*

Ensure: Top k principal components

1: Center the data:

Let
$$\mu = \frac{1}{n} \sum_{i=1}^{n} X_i$$
, $\overline{X} = X - \mu$



Algorithm Principal Component Analysis (PCA)

Require: Data matrix $X \in \mathbb{R}^{n \times d}$ with *n* samples and *d* features, number of components *k*

Ensure: Top k principal components

1: Center the data:

Let
$$\mu = \frac{1}{n} \sum_{i=1}^{n} X_i$$
, $\tilde{X} = X - \mu$

2: Compute the covariance matrix:

$$C = \frac{1}{n} \tilde{X}^{\top} \tilde{X}$$



Algorithm Principal Component Analysis (PCA) - Continued

3: Compute the eigenvectors and eigenvalues of C:

 $Cv_i = \lambda_i v_i, \quad i = 1, \dots, d$



Algorithm Principal Component Analysis (PCA) - Continued

3: Compute the eigenvectors and eigenvalues of C:

 $Cv_i = \lambda_i v_i, \quad i = 1, \dots, d$

4: Sort the eigenvectors by decreasing eigenvalues:

 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$



Algorithm Principal Component Analysis (PCA) - Continued

3: Compute the eigenvectors and eigenvalues of C:

 $Cv_i = \lambda_i v_i, \quad i = 1, \dots, d$

4: Sort the eigenvectors by decreasing eigenvalues:

 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$

5: Select the top *k* eigenvectors:

$$W = [v_1, v_2, \ldots, v_k]$$



Algorithm Principal Component Analysis (PCA) - Continued

6: **Project the data onto the top** *k* **components:**

 $Z = \tilde{X}W$



Algorithm Principal Component Analysis (PCA) - Continued

6: **Project the data onto the top** *k* **components:**



Algorithm Principal Component Analysis (PCA) - Continued

6: **Project the data onto the top** *k* **components:**

 $Z = \tilde{X}W$

7: return Projected data Z and components W

Viz

Implementation of PCA



- Notebook
- Blog





PCA may not always work!



Applying PCA for classification







 PCA























• What is the mean?



- What is the mean?
- What if the samples are Bernoulli draws from p = 0.5?




Given *n* samples x_1, x_2, \ldots, x_n

- What is the mean?
- What if the samples are Bernoulli draws from p = 0.5?
- What if the samples are from exponential distribution $\lambda = 5$?

$$\left(\mu = \frac{1}{2} - \frac{2}{2} \times i\right)$$

Sample Mean



Sample Mean:



Sample Mean



Sample Mean:



True Mean: μ Viz



Let $X_1, X_2, ..., X_n$ be a sequence of independent and identically distributed (i.i.d.) random variables with finite mean $\mu = \mathbb{E}[X_i]$ and finite variance $\sigma^2 = \text{Var}(X_i)$.



Let X_1, X_2, \ldots, X_n be a sequence of independent and identically distributed (i.i.d.) random variables with finite mean $\mu = \mathbb{E}[X_i]$ and finite variance $\sigma^2 = \operatorname{Var}(X_i)$. Define the sample mean as:

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$



Let X_1, X_2, \ldots, X_n be a sequence of independent and identically distributed (i.i.d.) random variables with finite mean $\mu = \mathbb{E}[X_i]$ and finite variance $\sigma^2 = \text{Var}(X_i)$. Define the sample mean as:

Then, for any $\varepsilon > 0$:



Weak Law of Large Numbers $E[X_n] = \mu$ $Vou(\bar{x}_n) = \underline{\Gamma}^2$ Var (Xn) Proof? $\mathbb{P}\left(\left|\overline{X}_{n}-\mu\right|>\varepsilon\right)\leq$ P=0.5 Ö 0: 5





Proof? What is the probability that you get all heads when you toss a fair coin?



I estimate the sum of n random real numbers by rounding each to the nearest integer, and adding the resulting integers. What is the probability that the total error is at most $\pm \sqrt{n}$?



Let X_1, X_2, \ldots, X_n be independent and identically distributed (i.i.d.) random variables with expected value $\mathbb{E}[X_i] = \mu < \infty$ and variance $0 < \operatorname{Var}(X_i) = \sigma^2 < \infty$. Then, as $n \to \infty$



Let X_1, X_2, \ldots, X_n be independent and identically distributed (i.i.d.) random variables with expected value $\mathbb{E}[X_i] = \mu < \infty$ and variance $0 < \operatorname{Var}(X_i) = \sigma^2 < \infty$. Then, as $n \to \infty$

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$



Let X_1, X_2, \ldots, X_n be independent and identically distributed (i.i.d.) random variables with expected value $\mathbb{E}[X_i] = \mu < \infty$ and variance $0 < \operatorname{Var}(X_i) = \sigma^2 < \infty$. Then, as $n \to \infty$

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Viz

Central Limit Theorem



I estimate the sum of *n* random real numbers by rounding each to the nearest integer, and adding the resulting integers. What is the probability that the total error is at most $\pm \sqrt{n}$?

$$P\left(\Xi X_{n} \leq I R\right) \in \mathbb{R}$$

1-200