

CS156 Session 7 - Principal component analysis (PCA)

≡ HCs/LOs	#cs156-unsupervisedlearning
↻ Pre-class	

Readings

(Option 1) Read sections 15.1-15.3 Barber, D. (2012). *Bayesian reasoning and machine learning*. Cambridge: Cambridge University Press.

This reading introduces PCA. It explains why we might need to use PCA, and gives the derivation for it. Focus on the high-level understanding, and be able to sketch out the derivation for PCA.

15.1 High-Dimensional Spaces – Low Dimensional Manifolds

- ML problem has high dimensional data: images, bag-of-words, gene expression
- the data does not occupy all spaces in high dimensions
 - linear dimension reductions → project high dimensional datapoint x onto a lower dimensional vector y

$$y = Fx + \text{const.}$$

$$F \rightarrow \dim(y) * \dim(x)$$

15.2 Principal Components Analysis

$$\mathbf{x}^n \approx \mathbf{c} + \sum_{j=1}^M y_j^n \mathbf{b}^j \equiv \tilde{\mathbf{x}}^n$$

- $b_j \rightarrow$ the basis vectors of the linear subspace (principal component coefficients or loadings)

15.2.1 Deriving the optimal linear reconstruction

- Find the best basis vector $B \rightarrow$ minimize the sum of squared difference

$$E(B, Y) = \sum_{n=1}^N \sum_{i=1}^D \left[x_i^n - \sum_{j=1}^M y_j^n b_i^j \right]^2 = \text{trace} \left((\mathbf{X} - \mathbf{BY})^\top (\mathbf{X} - \mathbf{BY}) \right)$$

- Optimal solution for B and Y is not unique \rightarrow need to constrain B to be an orthonormal matrix
- Since we wish to minimise $E(B)$, we therefore **define the basis using the eigenvectors with largest corresponding eigenvalues**

15.2.2 Maximum variance criterion

- To break the invariance of least squares projection with respect to rotations and rescaling, we need an additional criterion
- first searching for the single direction b such that the variance of the data projected onto this direction is maximal amongst all possible such projections

$$y^n = \sum_i b_i x_i^n$$

The projection of a datapoint onto a direction \mathbf{b} is $\mathbf{b}^\top \mathbf{x}^n$ for a unit length vector \mathbf{b} . Hence the sum of squared projections is

$$\sum_n (\mathbf{b}^\top \mathbf{x}^n)^2 = \mathbf{b}^\top \left[\sum_n \mathbf{x}^n (\mathbf{x}^n)^\top \right] \mathbf{b} = (N-1) \mathbf{b}^\top \mathbf{S} \mathbf{b} \quad (15.2.15)$$

- the optimal single b which **maximises the projection variance** is given by the eigenvector corresponding to the **largest eigenvalue of S**
- These maximal variance directions found by PCA are called the **principal directions**

15.2.4 PCA and nearest neighbours classification

- For high-dimensional data computing the squared Euclidean distance between vectors can be **expensive**, and also **sensitive to noise**.
→ project the data to a lower dimensional representation first.

$$\begin{aligned}(\mathbf{x}^a - \mathbf{x}^b)^T(\mathbf{x}^a - \mathbf{x}^b) &\approx (\mathbf{E}\mathbf{y}^a + \mathbf{m} - \mathbf{E}\mathbf{y}^b - \mathbf{m})^T(\mathbf{E}\mathbf{y}^a + \mathbf{m} - \mathbf{E}\mathbf{y}^b - \mathbf{m}) \\&= (\mathbf{y}^a - \mathbf{y}^b)^T \mathbf{E}^T \mathbf{E} (\mathbf{y}^a - \mathbf{y}^b) \\&= (\mathbf{y}^a - \mathbf{y}^b)^T (\mathbf{y}^a - \mathbf{y}^b)\end{aligned}$$

15.2.5 Comments on PCA

- The ‘intrinsic’ dimension of data
 - the reconstruction error is proportional to the sum of the discarded eigenvalues
→ the number of large eigenvalues is the indication of the number of degrees of freedom in the data
 - small eigenvalues are noise
- Non-linear dimension reduction

15.3 High Dimensional Data

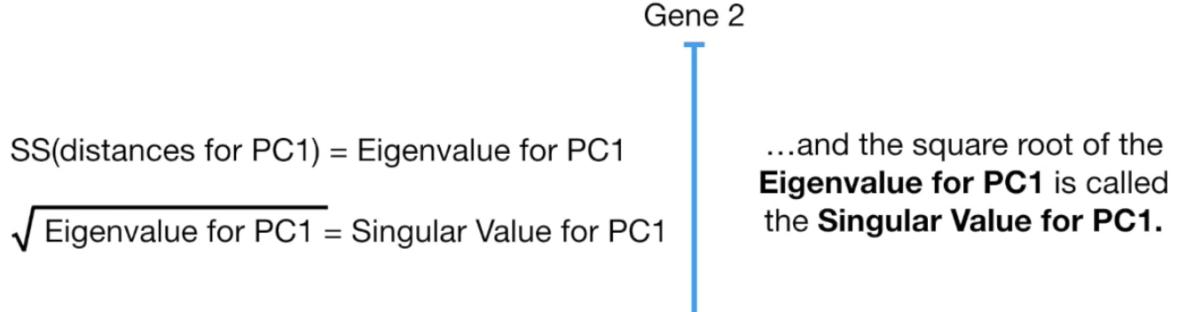
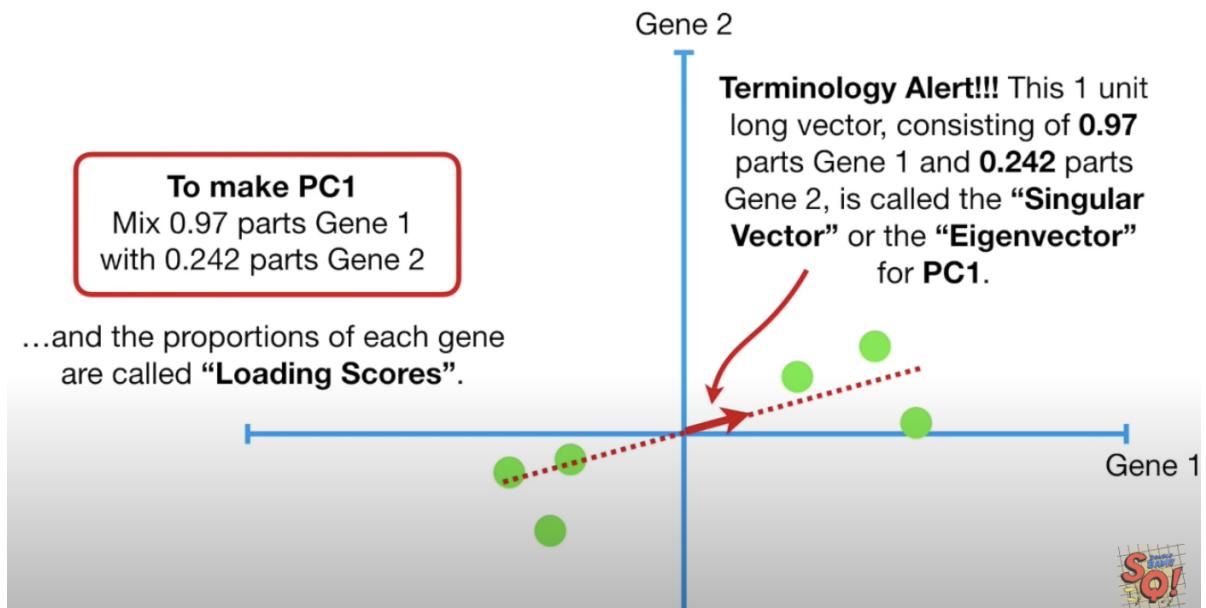
15.3.2 PCA via Singular value decomposition

Starmer, J. (2018). StatQuest: Principal Component Analysis (PCA), Step-by-Step Retrieved January 31, 2018, from

<https://www.youtube.com/watch?v=FgakZw6K1QQ> ##### Retrieved from <https://www.youtube.com/watch?v=FgakZw6K1QQ>

This video helps to understand PCA visually. Watch this first if you have no prior background.

▼ Notes



$$\frac{\text{SS}(\text{distances for PC1})}{n - 1} = \text{Variation for PC1}$$

$$\frac{\text{SS}(\text{distances for PC2})}{n - 1} = \text{Variation for PC2}$$

(Option 2) Read sections 12.2 and 12.3 of Murphy, K. (2012). Machine Learning: A Probabilistic Perspective. Cambridge: The MIT Press.

Read the following: 1.3.2 Discovering latent factors (p.11) - an introduction to why we need dimensionality reduction 12.2 Principal components analysis (PCA) (p.387) -ignore the first paragraph that references previous material, the sections after 12.2.1 are optional (Proof, Singular value decomposition, Probabilistic PCA, EM algorithm for PCA) 12.2.1 Classical PCA: - statement of the theorem (p.387) 12.3.2 Model selection for PCA (p. 399) - This section introduces metrics to assess how many latent dimensions (components) to choose.

1.3.2 Discovering latent factors

- **latent factors:** the data may appear high dimensional, there may only be a small number of degrees of variability
- **PCA:**
 - an unsupervised version of (multi-output) linear regression
 - observe the high-dimensional response y , but not the low-dimensional “cause” z .
 - Thus the model has the form $z \rightarrow y$; we have to “invert the arrow”
 - **infer the latent low-dimensional z from the observed high-dimensional y**

12.2.1 Classical PCA: statement of the theorem

Theorem 12.2.1. Suppose we want to find an orthogonal set of L linear basis vectors $\mathbf{w}_j \in \mathbb{R}^D$, and the corresponding scores $\mathbf{z}_i \in \mathbb{R}^L$, such that we minimize the average reconstruction error

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 \quad (12.26)$$

where $\hat{\mathbf{x}}_i = \mathbf{W}\mathbf{z}_i$, subject to the constraint that \mathbf{W} is orthonormal. Equivalently, we can write this objective as follows:

$$J(\mathbf{W}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{WZ}^T\|_F^2 \quad (12.27)$$

where \mathbf{Z} is an $N \times L$ matrix with the \mathbf{z}_i in its rows, and $\|\mathbf{A}\|_F$ is the **Frobenius norm** of matrix \mathbf{A} , defined by

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2} = \sqrt{\text{tr}(\mathbf{A}^T \mathbf{A})} = \|\mathbf{A}(:,\cdot)\|_2 \quad (12.28)$$

The optimal solution is obtained by setting $\hat{\mathbf{W}} = \mathbf{V}_L$, where \mathbf{V}_L contains the L eigenvectors with largest eigenvalues of the empirical covariance matrix, $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$. (We assume the \mathbf{x}_i have zero mean, for notational simplicity.) Furthermore, the optimal low-dimensional encoding of the data is given by $\hat{\mathbf{z}}_i = \mathbf{W}^T \mathbf{x}_i$, which is an orthogonal projection of the data onto the column space spanned by the eigenvectors.

Study Guide

Please consult the study guide given in the course repo:

<https://github.com/minerva-schools/cs156>

(Don't forget to double check that you are doing the correct session; both the session number and the session title should be the same on Forum as in the repo.)

Pre-class Work

Please consult the pre-class work given in the course repo. Ensure you submit your work to your personal repository in an appropriately-named folder before class starts.

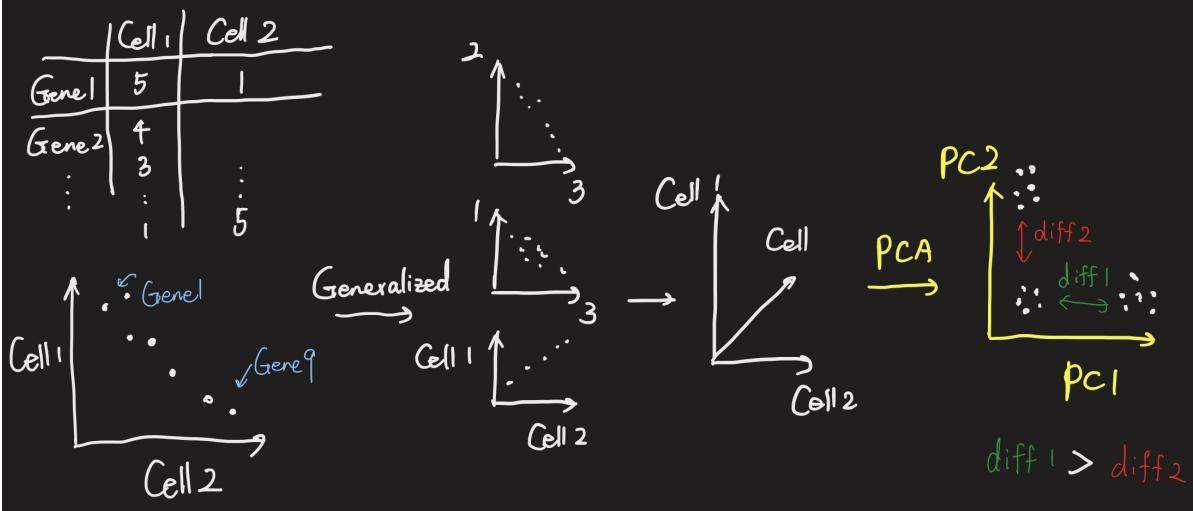
▼ PCA algor concepts

PCA

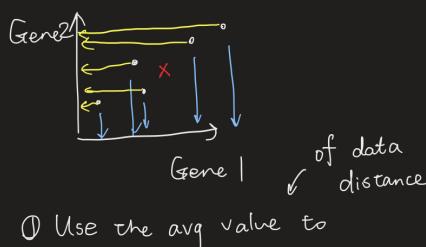
Saturday, January 29, 2022 2:09 PM

Convert correlations of data into 2D graphs → Dimension reduction

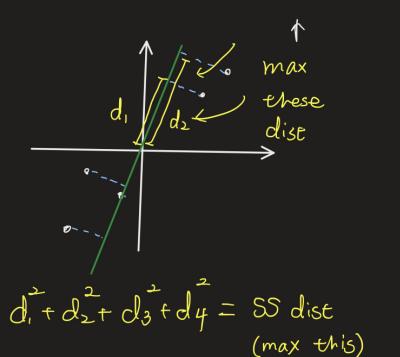
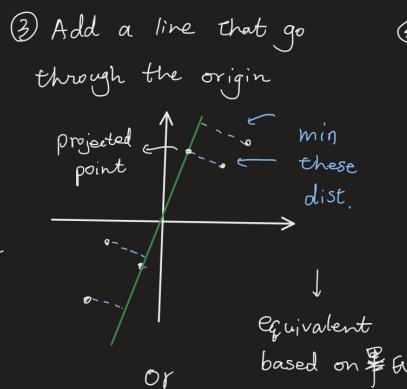
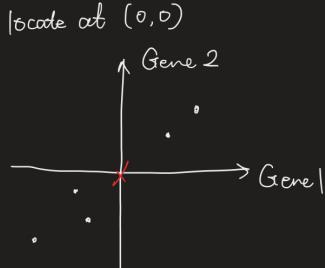
Example:



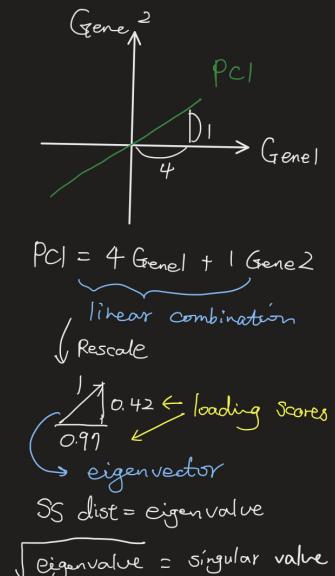
Steps for PCA

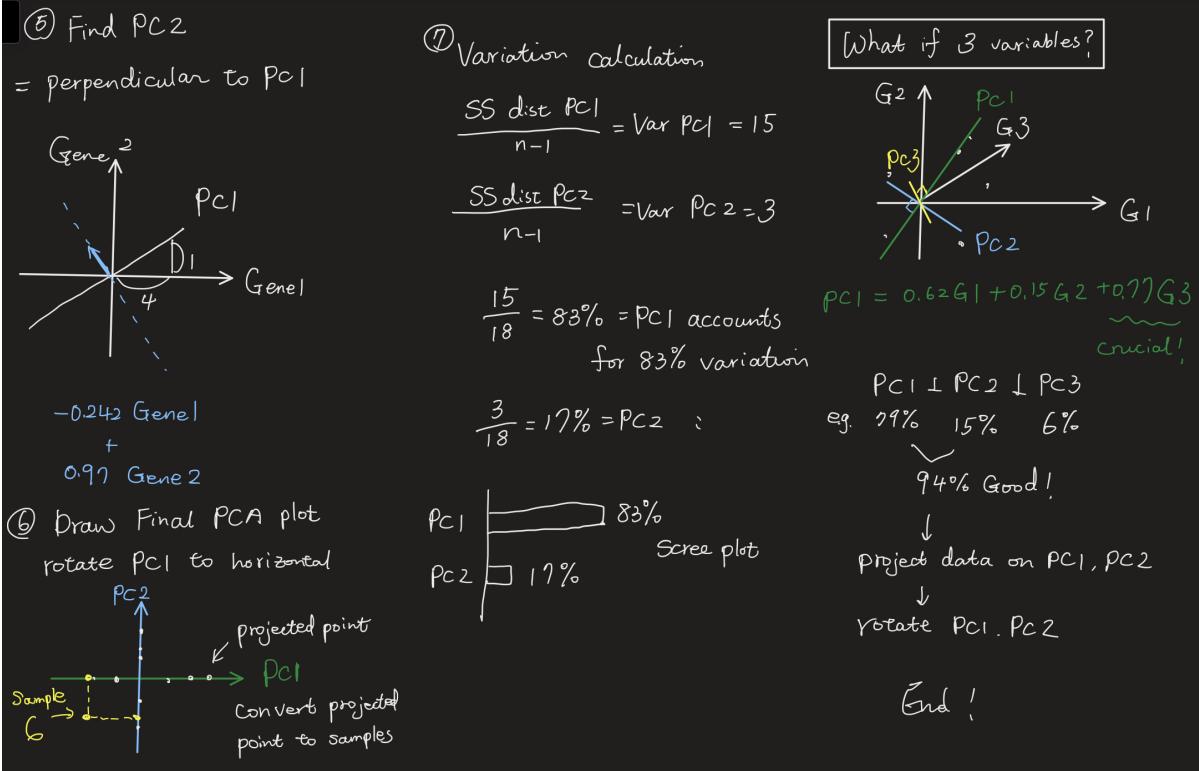


- ① Use the avg value to find the center \bar{x}
- ② Shift the data \rightarrow make \bar{x}



- ③ Add a line that go through the origin
 - ④ Find PC1
- = the line that max SS dist





▼ Preclass: prove eigenvectors are orthogonal to each other

Eigenvectors orthogonal

$$y_2 = X \cdot a_2$$

$$y_1 = X \cdot a_1$$

projection

Goal: $\max(y_1) = a_1^T \Sigma a_1$

Constraints: $a_1^T a_1 = 1$

Max $a_1^T \Sigma a_1 - \lambda(a_1^T a_1 - 1)$

$$\Sigma a_1 - \lambda a_1 = 0$$

$$(\Sigma - \lambda I_p) a_1 = 0$$

y_1 will be max for biggest λ_1 , (if?

$$\lambda_1 a_1 = \Sigma a_1$$

$$\text{Cov}(y_1, y_2) = \text{Cov}(a_1^T X, a_2^T X) = a_1^T \Sigma a_2 = a_2^T \Sigma a_1 = a_2^T \lambda_1 a_1 = \lambda_1 a_2^T a_1 = 0 \text{ if and only if } a_2^T a_1 = 0$$

Lagrangian

Example

$$R(x, y) = x^2 e^y y$$

$$B(x, y) = x^2 + y^2 = b$$

$$\nabla R = \lambda \nabla B$$

$$\text{constraint}$$

$$L(x, y, \lambda) = R(x, y) - \lambda(B(x, y) - b)$$

$$\nabla L = 0$$

$$\begin{bmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \\ \frac{\partial L}{\partial \lambda} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\frac{\partial L}{\partial x} = \frac{\partial R}{\partial x} - \lambda \frac{\partial B}{\partial x} = 0 \leftarrow$$

$$\frac{\partial L}{\partial y} = \frac{\partial R}{\partial y} - \lambda \frac{\partial B}{\partial y} = 0 \leftarrow$$

$$\frac{\partial L}{\partial \lambda} = -(B(x, y) - b) = 0 \quad B(x, y) = b$$

<https://stats.stackexchange.com/questions/266652/why-are-pca-eigenvectors-orthogonal-and-what-is-the-relation-to-the-pca-scores-b>

- why the biggest lambda will max the y?
- why we want cov matrix = 0 ?
 - bc cov matrix will be symmetric $x^T y = 0$
- who is the eigen vector here? (y or a?)
- Why is the goal looks like that?

▼ Math steps to find PCA dimensions(aka eigenvector)

PCA math steps

- ① Calculate covariance matrix of the dataset
- ② Find eigenvectors & eigenvalues
- ③ Sort eigenvalues & vec

eigenvalue 1 & \vec{v}_1 , eigenvector
✓
✓ eigenvalue 2 & \vec{v}_2
✓ eigenvalue 3 & \vec{v}_3

$$W = \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \\ \vec{v}_3 \end{bmatrix}$$

- ④ Use eigenvector matrix W to transform the dataset

▼ PCA math foundations (Questions)

PCA math and question list

Vector projection

$$\vec{a} \text{ on } \vec{b} = \frac{\vec{a} \cdot \vec{b}}{\|\vec{b}\|} = \text{proj}_{\vec{b}} \vec{a}$$

$$\vec{a} \cdot \vec{b} = \|\vec{a}\| \|\vec{b}\| \cos \theta$$

Question.

① Why is C a correlation matrix?

isn't it dot product of two vectors?

$$\begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \end{bmatrix} \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \end{bmatrix}^T = \begin{bmatrix} (1, -2) \\ (1, -2) \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} (1, -2) \\ (1, -2) \\ \vdots \end{bmatrix}$$

$$= \begin{bmatrix} 5 \\ \vdots \end{bmatrix}$$

② Differentiable a matrix? a^T

③ Is correlational matrix diagonal?

Connect to proof of eigenvectors are orthogonal to each other.

Covariance matrix & Correlation matrix

Covariance matrix

→ measure how much x, y changes together

$$\text{Cov}(x, y) = \sum \frac{(x_i - \bar{x})(y_i - \bar{y})}{N} = \sum \frac{x_i y_i}{N} - \frac{1}{N} (x^T y)$$

X	N	Y	$x - \bar{x}$	$y - \bar{y}$
1	2	-2	-4	-2
2	8	-1	2	-1
3	6	0	0	0
4	4	1	-2	1
5	10	2	4	2

Correlation Matrix C

→ find coeffs for a set of vars (x, y)

$$C = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}$$

$$= \frac{\overline{z^T z} - \overline{z} \overline{z^T}}{\overline{z^T z} - \overline{z^T z}}$$

diagonal

Cov matrix Z

$$Z = \begin{bmatrix} -2 & -4 \\ -1 & 2 \\ 0 & 0 \\ 1 & -2 \\ 2 & 4 \end{bmatrix} \quad \frac{1}{N-1} Z^T Z = \frac{1}{4} \begin{bmatrix} -2 & -1 & 0 & 1 & 2 \\ -4 & 2 & 0 & -2 & 4 \end{bmatrix} \begin{bmatrix} -2 & -4 \\ -1 & 2 \\ 0 & 0 \\ 1 & -2 \\ 2 & 4 \end{bmatrix}$$

$$= \frac{1}{4} \begin{bmatrix} 10 & 12 \\ 12 & 40 \end{bmatrix} = \begin{bmatrix} 2.5 & 3 \\ 3 & 10 \end{bmatrix}$$

covariance = $\begin{bmatrix} S_{xx} & S_{xy} \\ S_{xy} & S_{yy} \end{bmatrix}$

Variance of y

Eigenvector review

Eigenvector

- special vectors that remains unchanged after a linear transformation
- find a vector \vec{u}_i such that $T(\vec{u}_i) = \lambda_i \vec{u}_i$

λ_i : eigenvalues
 \vec{u}_i : eigenvectors

Example $T(x, y) = (-x, y)$

$$T(\vec{u}) = \lambda \vec{u}$$

$$T(u_1, u_2) = (-u_1, u_2)$$

$$\lambda(u_1, u_2) = (-u_1, u_2)$$

$$\begin{cases} \lambda u_1 = -u_1 \\ \lambda u_2 = u_2 \end{cases} \xrightarrow{\text{solve}} \lambda = \begin{cases} 1 & \Rightarrow u_2 = 0 \quad (1, 0) \\ -1 & \Rightarrow u_1 = 0 \quad (0, 1) \end{cases}$$

if we have $\vec{v} = (2, 1)$
 $\vec{v} = 2(1, 0) + (0, 1)$
 $T(\vec{v}) = 2T(1, 0) + T(0, 1) = (-2, 1)$

Math about coef & cov & eigen

The Mathematics Behind Principal Component Analysis

The central idea of principal component analysis (PCA) is to reduce the dimensionality of a data set consisting of a large number of interrelated variables while retaining as much as possible of the

<https://towardsdatascience.com/the-mathematics-behind-principal-component-analysis-fff2d7f4b643>

coef matrix

How to express a correlation matrix in terms of a covariance matrix?

I'd like to know, if I have the following Normal multivariate structure $\left[\begin{array}{c c} Y_{r_1} & Y_{r_2} \end{array} \right] \sim \mathcal{N}_{r} \left(\begin{array}{a} \dots \end{array} \right)$

<https://stats.stackexchange.com/questions/413033/how-to-express-a-correlation-matrix-in-terms-of-a-covariance-matrix>

$$E((u \cdot x)^2)$$

https://www.math.uci.edu/icamp/courses/math77b/lecture_12w/pdfs/PCA.pdf

Finn's note: [here](#)

Pre-class

1. Barber Exercise

Exercise 15.2. Consider a dataset in two dimensions where the data lies on the circumference of a circle of unit radius. What would be the effect of using PCA on this dataset, in which we attempt to reduce the dimensionality to 1? Suggest an alternative one dimensional representation of the data.

If the data is evenly distributed, as long as the PCA dimension crosses the origin, the effect of PCA will all be similar. Another method for dimensional representation is to use radio basis function to transform circular data into a straight line.