



Linear Dimensionality Reduction and Linear Discriminant Analysis

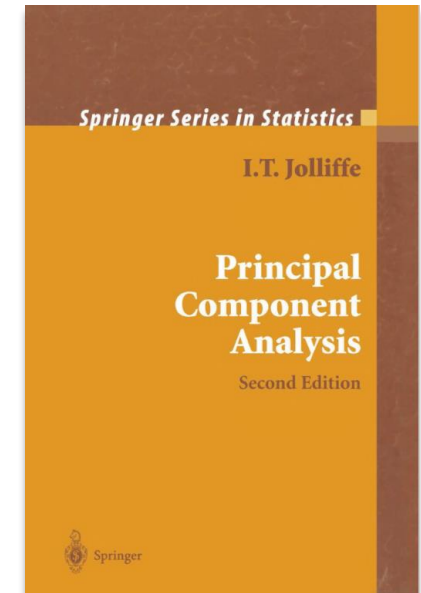
Assoc. Prof. Karl Ezra Pilario, Ph.D.

Process Systems Engineering Laboratory
Department of Chemical Engineering
University of the Philippines Diliman

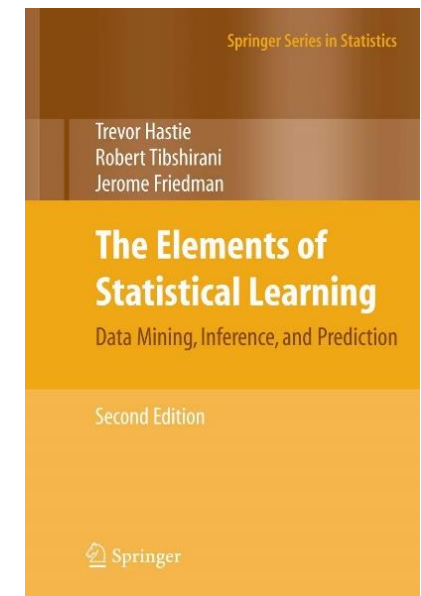
Outline

- Dimensionality Reduction
 - Curse of Dimensionality
 - Feature Selection vs. Feature Extraction
 - Principal Components Analysis (PCA)
 - Derivation
 - Non-negative Matrix Factorization (NMF)
 - Independent Components Analysis (ICA)
- Low-dimensional Classification
 - Linear Discriminant Analysis (LDA)

Jolliffe (2002)
Principal Component Analysis.
2nd Ed. Springer.

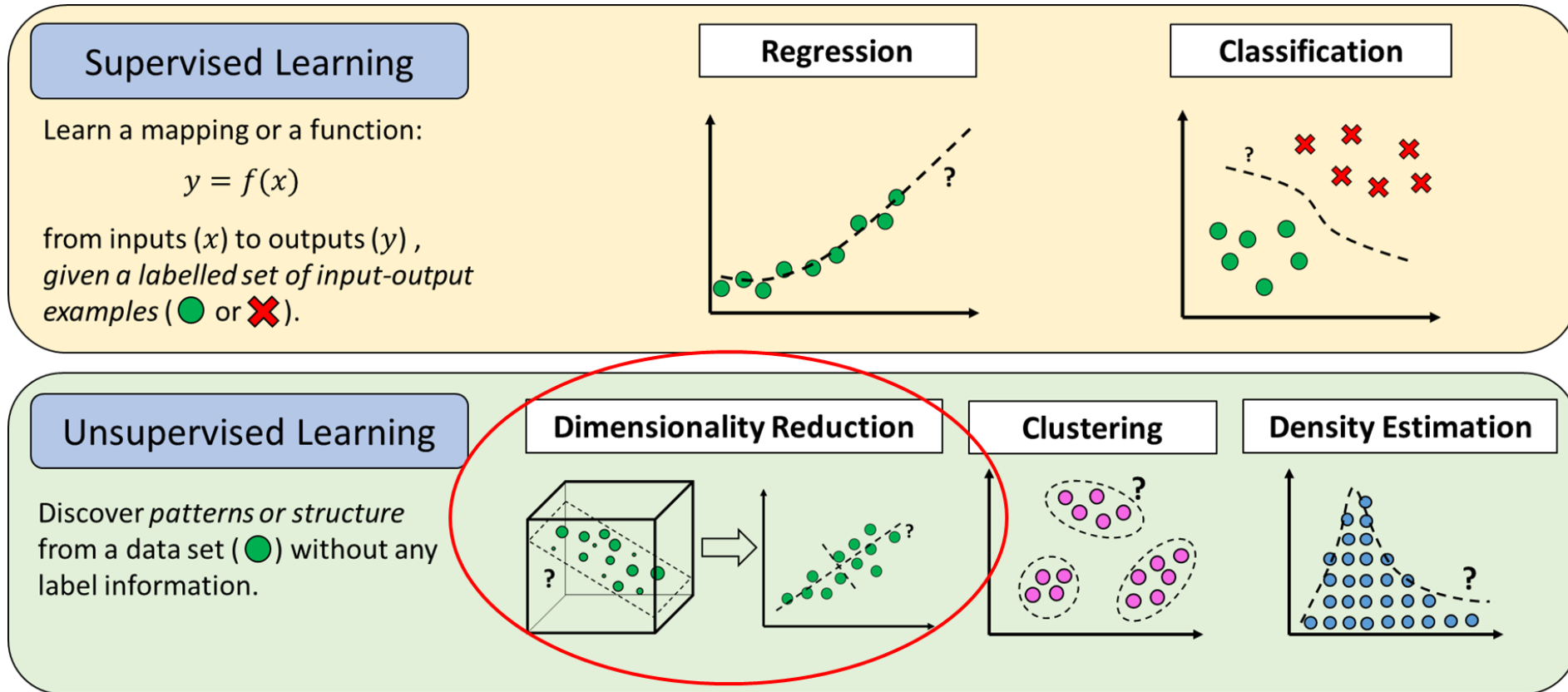


Hastie *et al.* (2008)
The Elements of Statistical Learning.
2nd Ed. Springer.



Overview

In the past few weeks, we mostly discussed supervised learning methods. Today, we'll move into unsupervised learning.



What is dimensionality?

Recall some data sets we discussed before:

	Fisher Iris Data Set	House Price Data Set	Handwritten Digits Images
Features	Sepal length Sepal width Petal length Petal width	Floor Area No. of Rooms Age	Grayscale values (x 64)
No. of Features	4	3	64

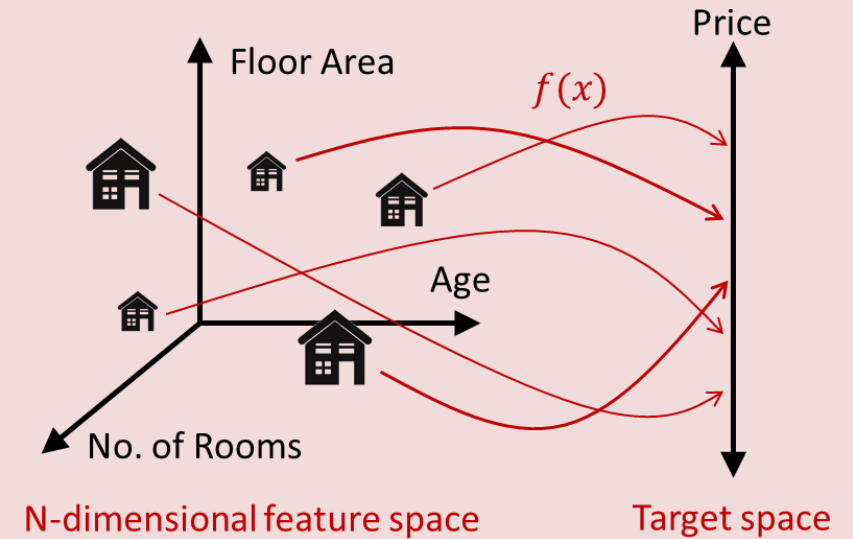
Most data sets have multiple features!

- Patient Data: Age, Height, Weight, BMI, Blood Type, ...
- Image Data: RGB values x No. of Pixels
- Weather Data: Temperature, Humidity, Wind speed, ...
- Car Data: Mileage, Horsepower, Weight, ...
- Etc.

Dimensionality

refers to the number of “attributes” or “features” in the data set.

If each sample is a point in the N-dimensional space defined by the features, then our **supervised learning model** $f(x)$ is nothing more than a *mapping* from feature space to target space.



The Curse of Dimensionality

More often, data sets are *high-dimensional*.
This is seen as a “curse” rather than a blessing.
(Bellman, 1961)



“The number of samples needed to estimate an arbitrary function with a given level of accuracy grows **exponentially** with respect to the number of input variables (i.e., dimensionality) of the function.”

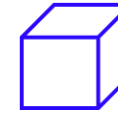
1-D
2 corners



2-D
4 corners



3-D
8 corners



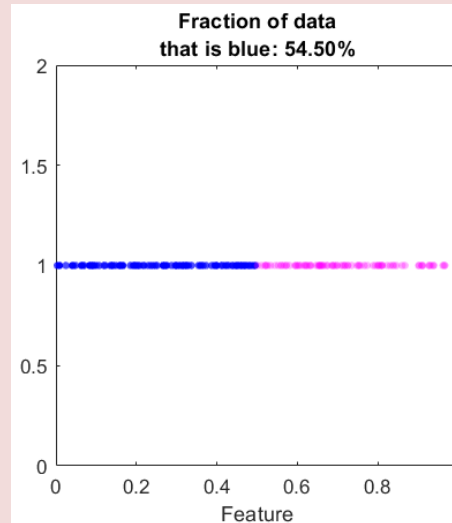
n-D
 2^n corners

N-dimensional
hypercube

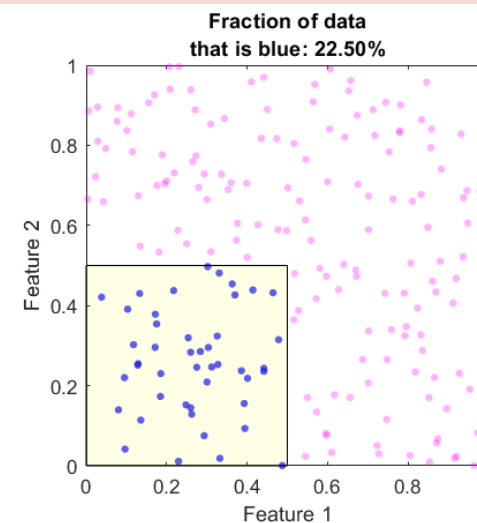
The more dimensions there are in your data, the more samples you need to cover the space.

Say, due to budget constraints, you only sampled *half* of the values in each dimension. As the number of dimensions increase, the fraction of the feature space that you covered becomes *a lot less than half*.

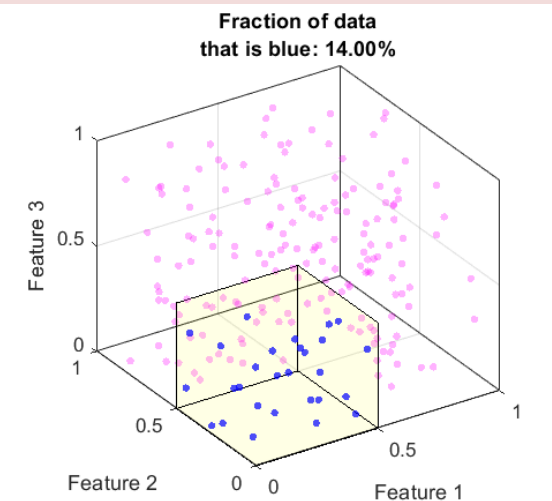
1-D



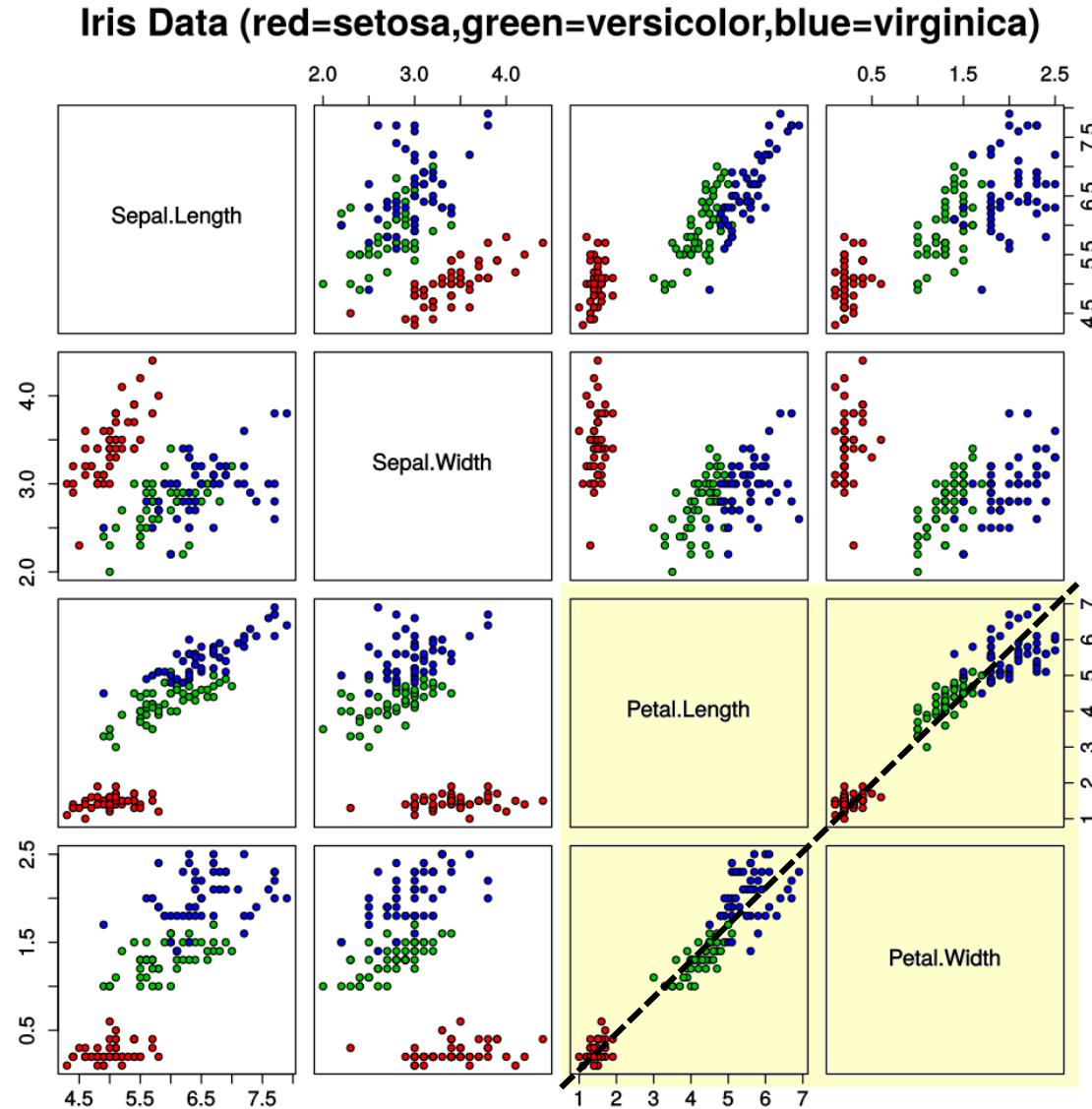
2-D



3-D



The Curse of Dimensionality



High-dimensionality comes with other problems such as:

- multi-collinear or redundant features, and
- unimportant features.

Multi-collinear or Redundant Features

- Features that are highly correlated with each other.
- *Positively correlated*: When one increases, the other also tends to increase. When one decreases, the other also tends to decrease.
- If 2 features are correlated, they could contain the same information content. Hence, one of them is *redundant*.

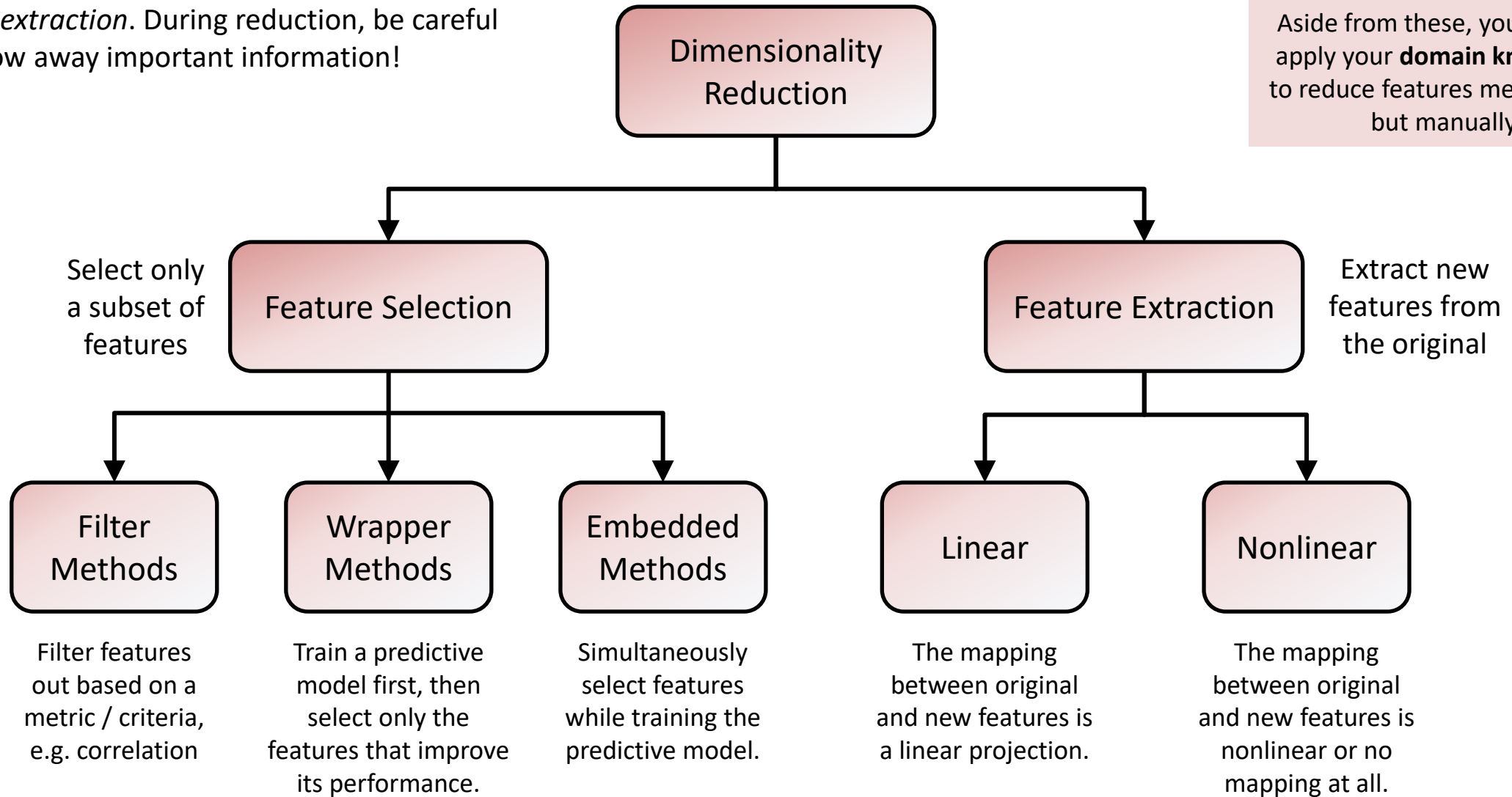
Unimportant Features

- Features that are *irrelevant* to the prediction task.
- These features can fool the learning algorithm, thinking that they contain valuable information, but they don't.

Dimensionality Reduction

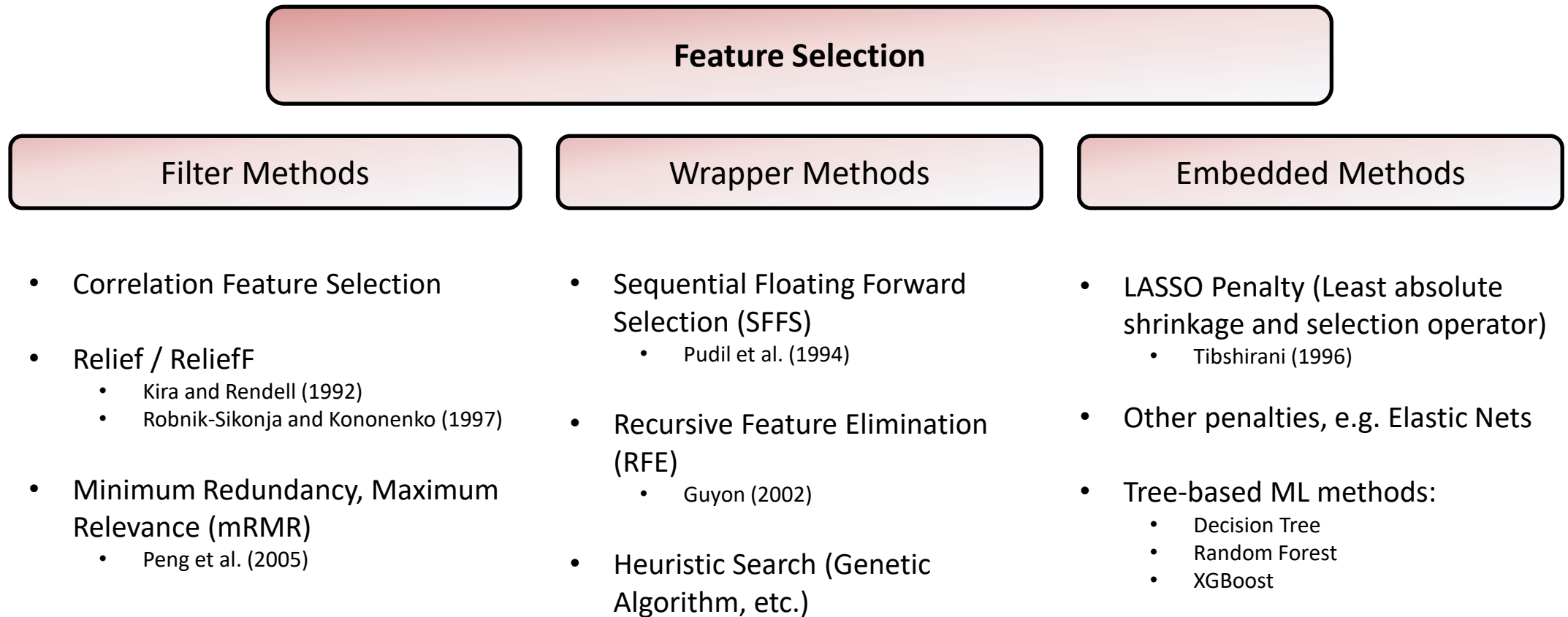
We can reduce features either by *feature selection* or *feature extraction*. During reduction, be careful not to throw away important information!

Aside from these, you can also apply your **domain knowledge** to reduce features meaningfully but manually.



Dimensionality Reduction

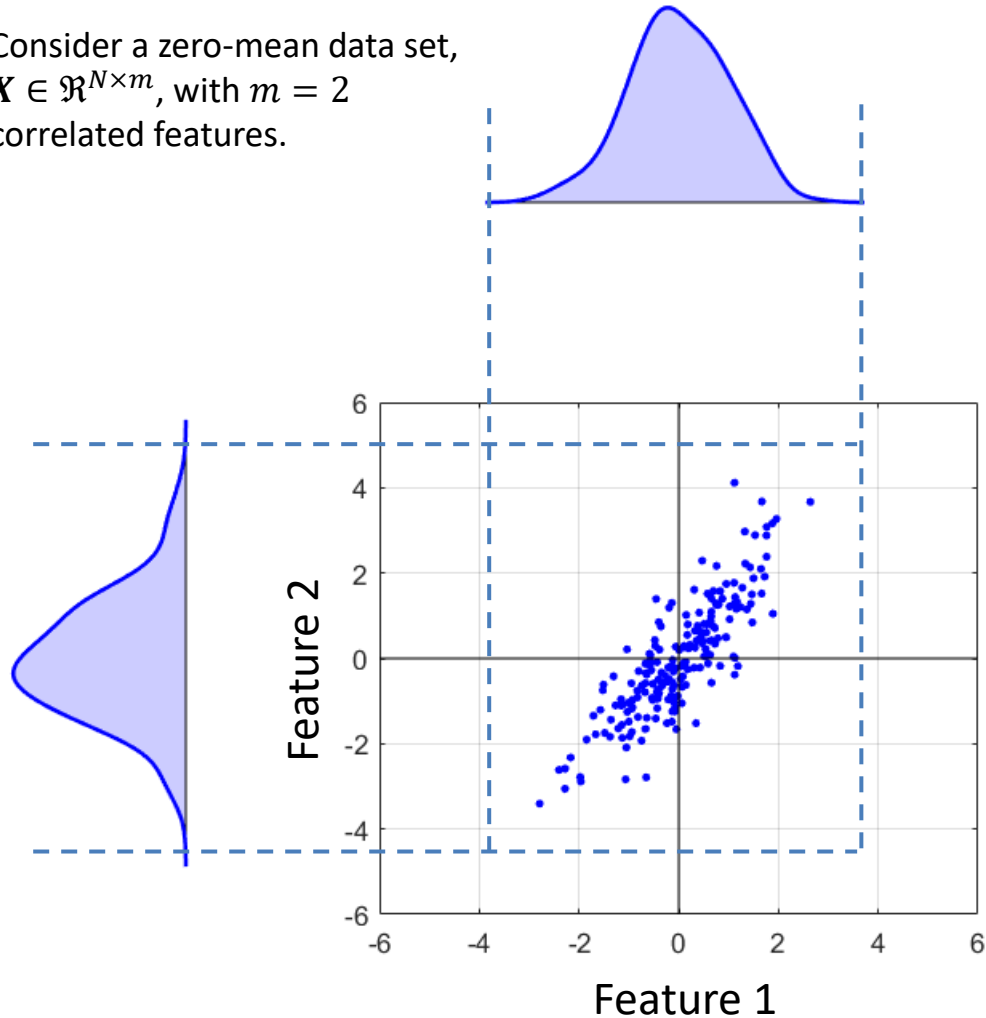
Here are some of the typical algorithms used in feature selection.



Principal Components Analysis

A popular linear feature extraction-based dimensionality reduction method due to Pearson (1901).

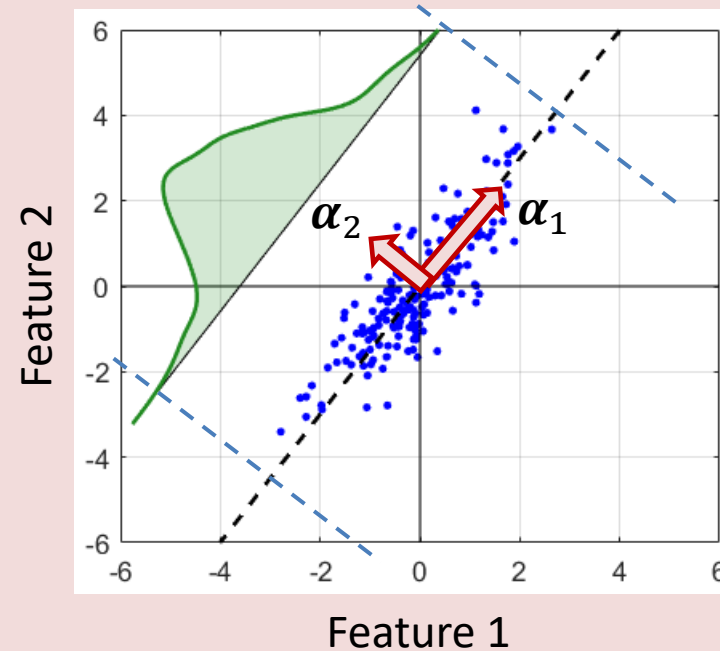
Consider a zero-mean data set, $\mathbf{X} \in \mathbb{R}^{N \times m}$, with $m = 2$ correlated features.



The goal of PCA is to find a projection matrix, $\mathbf{P} \in \mathbb{R}^{m \times m}$, such that \mathbf{P} is **orthonormal** and the *variance* of the projected data, $\mathbf{Y} \in \mathbb{R}^{N \times m}$, is **maximized**:

$$\mathbf{Y} = \mathbf{XP}$$

where: $\mathbf{Y} = [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \dots \quad \mathbf{y}_m]$ \mathbf{y} 's are called **scores**.
 $\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_m]$
 $\mathbf{P} = [\boldsymbol{\alpha}_1 \quad \boldsymbol{\alpha}_2 \quad \dots \quad \boldsymbol{\alpha}_m]$ $\boldsymbol{\alpha}$'s are called **loadings / coefficients**.



$(\mathbf{y}_1, \boldsymbol{\alpha}_1)$ is the **1st** principal component.
 $(\mathbf{y}_2, \boldsymbol{\alpha}_2)$ is the **2nd** principal component.
...and so on...

Why does PCA accomplish dimensionality reduction?

After PCA, we can take only the *first few* scores as the new extracted features, then discard the rest.

$\begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_m \end{matrix} \Rightarrow \begin{matrix} \mathbf{y}_1 \text{ (PC1)} \\ \mathbf{y}_2 \text{ (PC2)} \end{matrix}$

Principal Components Analysis

In the next few slides, we will derive how to obtain the PCA projection matrix, P .

$$Y = XP$$



$$y_1 = X\alpha_1$$

$$y_2 = X\alpha_2$$

$$y_3 = X\alpha_3$$

$$y_4 = X\alpha_4$$

...and so on...

Step 1

What is the variance that we are trying to maximize?

The **variance** of y_1 is given as:

$$\begin{aligned}\text{var}[y_1] &= (X\alpha_1)^T(X\alpha_1) \\ &= (\alpha_1^T X^T)(X\alpha_1) \\ &= (N-1)\alpha_1^T \Sigma \alpha_1\end{aligned}$$

where:

$$\Sigma = \frac{1}{N-1} X^T X$$

is the **sample covariance** of X .

Step 2

How is the optimization problem in PCA formulated?

Ans: It is actually a series of optimization problems!

Maximize var $[y_1]$

$$\max_{\alpha_1} \alpha_1^T \Sigma \alpha_1$$

Subject to:

$$\alpha_1^T \alpha_1 = 1$$

(α_1 should be a normal vector)

Maximize var $[y_2]$

$$\max_{\alpha_2} \alpha_2^T \Sigma \alpha_2$$

Subject to:

$$\alpha_1^T \alpha_2 = 0$$

$$\alpha_2^T \alpha_2 = 1$$

(α_2 should be a normal vector and orthogonal to α_1)

Maximize var $[y_3]$...

$$\max_{\alpha_3} \alpha_3^T \Sigma \alpha_3$$

Subject to:

$$\alpha_1^T \alpha_3 = 0$$

$$\alpha_2^T \alpha_3 = 0$$

$$\alpha_3^T \alpha_3 = 1$$

(α_3 should be a normal vector and orthogonal to α_1, α_2)

...and so on...

Reference:

[1] Jolliffe and Cadima (2016). Principal component analysis: a review and recent developments. <http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2015.0202>

[2] Jolliffe (2002). Principal Component Analysis. 2nd Ed. Springer Verlag-Berlin. <http://onlinelibrary.wiley.com/doi/10.1002/0470013192.bsa501/full>

Principal Components Analysis

In the next few slides, we will derive how to obtain the PCA projection matrix, P .

Step 3

How can we solve each optimization problem?

The standard approach in constrained optimization is to use **Lagrange multipliers**, λ :

1. Form the Lagrange equation, \mathcal{L} , as the objective function minus all the equality constraints multiplied to λ .
2. Optimize the Lagrange equation by differentiating it with respect to the variables, then equate them to 0.

Maximize var [y_1]

$$\max_{\alpha_1} \alpha_1^T \Sigma \alpha_1$$

Subject to:

$$\alpha_1^T \alpha_1 = 1$$

(α_1 should be a normal vector)

Form the Lagrange equation:

$$\mathcal{L}(\alpha_1, \lambda) = \alpha_1^T \Sigma \alpha_1 - \lambda(\alpha_1^T \alpha_1 - 1)$$

Differentiate and equate to 0:

$$\frac{\partial \mathcal{L}(\alpha_1, \lambda)}{\partial \alpha_1} = 2\Sigma \alpha_1 - 2\lambda \alpha_1 = 0$$

Rearrange the equation:

$$\Sigma \alpha_1 - \lambda \alpha_1 = 0$$

$$(\Sigma - \lambda I) \alpha_1 = 0$$

Insight: α_1 must be an eigenvector of Σ .

Substituting λ into the objective:

$$\max_{\alpha_1} \alpha_1^T \Sigma \alpha_1 = \alpha_1^T \lambda \alpha_1 = \lambda \overbrace{\alpha_1^T \alpha_1}^1 = \lambda$$

Insight: α_1 must be an eigenvector of Σ , corresponding to the largest eigenvalue λ .

Reference:

[1] Jolliffe and Cadima (2016). Principal component analysis: a review and recent developments. <http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2015.0202>

[2] Jolliffe (2002). Principal Component Analysis. 2nd Ed. Springer Verlag-Berlin. <http://onlinelibrary.wiley.com/doi/10.1002/0470013192.bsa501/full>

Principal Components Analysis

In the next few slides, we will derive how to obtain the PCA projection matrix, P .

Step 3

How can we solve each optimization problem?

The standard approach in constrained optimization is to use **Lagrange multipliers**, λ :

1. Form the Lagrange equation, \mathcal{L} , as the objective function minus all the equality constraints multiplied to λ .
2. Optimize the Lagrange equation by differentiating it with respect to the variables, then equate them to 0.

Maximize var [y_2]

$$\max_{\alpha_2} \alpha_2^T \Sigma \alpha_2$$

Subject to:

$$\alpha_1^T \alpha_2 = 0$$

$$\alpha_2^T \alpha_2 = 1$$

(α_2 should be a normal vector and orthogonal to α_1)

Form the Lagrange equation:

Differentiate and equate to 0:

Pre-multiply α_1^T :

Now that $\lambda_1 = 0$, we can proceed:

$$\mathcal{L}(\alpha_2, \lambda) = \alpha_2^T \Sigma \alpha_2 - \lambda_1 \alpha_1^T \alpha_2 - \lambda_2 (\alpha_2^T \alpha_2 - 1)$$

$$\frac{\partial \mathcal{L}(\alpha_2, \lambda)}{\partial \alpha_2} = 2\Sigma \alpha_2 - 2\lambda_1 \alpha_1 - 2\lambda_2 \alpha_2 = 0$$

$$\Sigma \alpha_2 - \lambda_1 \alpha_1 - \lambda_2 \alpha_2 = 0$$

$$\underbrace{\alpha_1^T \Sigma \alpha_2}_0 - \lambda_1 \underbrace{\alpha_1^T \alpha_1}_1 - \lambda_2 \underbrace{\alpha_1^T \alpha_2}_0 = 0$$

Hence, $\lambda_1 = 0$.

$$\Sigma \alpha_2 - \lambda_2 \alpha_2 = 0$$

$$(\Sigma - \lambda_2 I) \alpha_2 = 0$$

Insight: α_2 must be an eigenvector of Σ , corresponding to the 2nd largest eigenvalue, λ_2 .

Reference:

- [1] Jolliffe and Cadima (2016). Principal component analysis: a review and recent developments. <http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2015.0202>
[2] Jolliffe (2002). Principal Component Analysis. 2nd Ed. Springer Verlag-Berlin. <http://onlinelibrary.wiley.com/doi/10.1002/0470013192.bsa501/full>

Principal Components Analysis

In the next few slides, we will derive how to obtain the PCA projection matrix, P .

Step 4

Finally, we realize that the series of optimization problems in PCA can be solved by simply taking the eigenvalue decomposition of the covariance matrix, Σ . Each eigenvalue-eigenvector pair corresponds to a principal component!

Recall: Eigenvalue Decomposition

$$\Sigma = V \times \Lambda \times V^T$$

$$\Sigma = \begin{bmatrix} a_1 & b_1 & c_1 & \dots \\ a_2 & b_2 & c_2 & \dots \\ a_3 & b_3 & c_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} a_1 & a_2 & a_3 & \dots \\ b_1 & b_2 & b_3 & \dots \\ c_1 & c_2 & c_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

V is our desired projection matrix, $V = P$.

Note: For symmetric matrices such as Σ , the matrix V satisfies $V^{-1} = V^T$.

Reference:

[1] Jolliffe and Cadima (2016). Principal component analysis: a review and recent developments. <http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2015.0202>

[2] Jolliffe (2002). Principal Component Analysis. 2nd Ed. Springer Verlag-Berlin. <http://onlinelibrary.wiley.com/doi/10.1002/0470013192.bsa501/full>

Maximize var [y_1]

Maximize var [y_2]

Maximize var [y_3]...

$$\max_{\alpha_1} \alpha_1^T \Sigma \alpha_1$$



$$\max_{\alpha_2} \alpha_2^T \Sigma \alpha_2$$



$$\max_{\alpha_3} \alpha_3^T \Sigma \alpha_3$$

Subject to:

$$\alpha_1^T \alpha_1 = 1$$

(α_1 should be a normal vector)

Subject to:

$$\alpha_1^T \alpha_2 = 0$$

$$\alpha_2^T \alpha_2 = 1$$

(α_2 should be a normal vector and orthogonal to α_1)

Subject to:

$$\alpha_1^T \alpha_3 = 0$$

$$\alpha_2^T \alpha_3 = 0$$

$$\alpha_3^T \alpha_3 = 1$$

(α_3 should be a normal vector and orthogonal to α_1, α_2)

...and so on...



Solutions:

1. α_1 must be an eigenvector of Σ , corresponding to the **largest** eigenvalue λ_1 .
2. α_2 must be an eigenvector of Σ , corresponding to the **2nd largest** eigenvalue λ_2 .
3. α_3 must be an eigenvector of Σ , corresponding to the **3rd largest** eigenvalue λ_3 .
4. And so on...

Principal Components Analysis

Now, we present the main PCA algorithm.

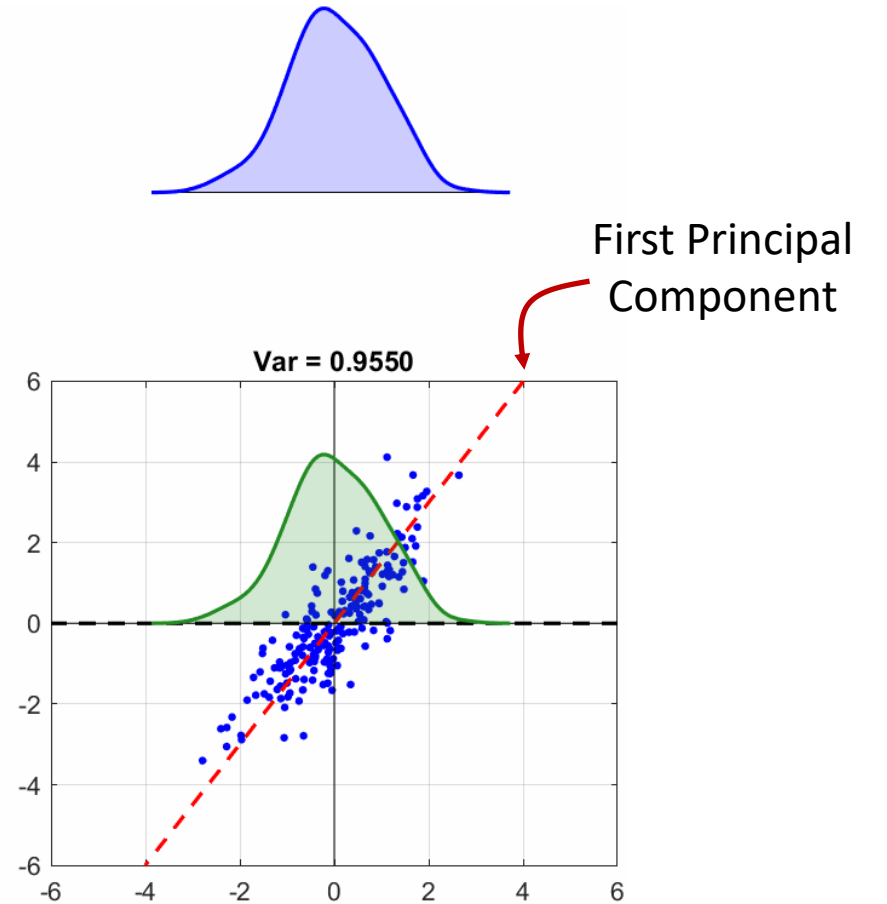
PCA Algorithm

1. Standardize the Data (zero-mean, unit-variance)

2. Compute the covariance of \mathbf{X} : $\Sigma = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$

3. Compute the eigenvalue decomposition of Σ :
 $\Sigma = \mathbf{V} \Lambda \mathbf{V}^T$

4. Choose only n principal components, then get \mathbf{Y} :
 $\mathbf{P} = \mathbf{V}_n$
 $\mathbf{Y} = \mathbf{X} \mathbf{P}$



Reference:

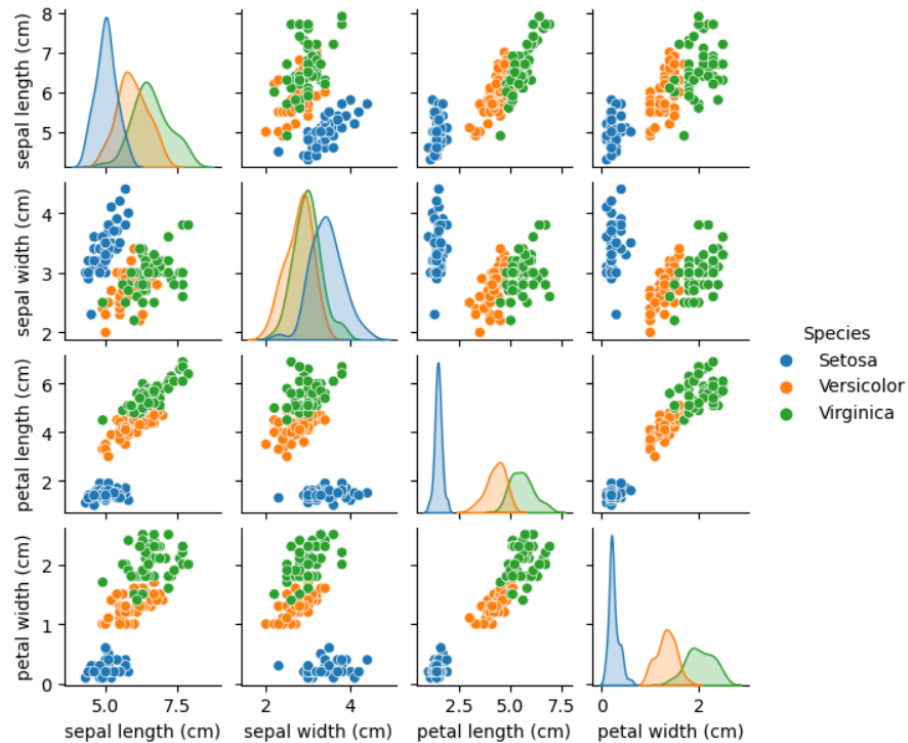
[1] Jolliffe and Cadima (2016). Principal component analysis: a review and recent developments. <http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2015.0202>

[2] Jolliffe (2002). Principal Component Analysis. 2nd Ed. Springer Verlag-Berlin. <http://onlinelibrary.wiley.com/doi/10.1002/0470013192.bsa501/full>

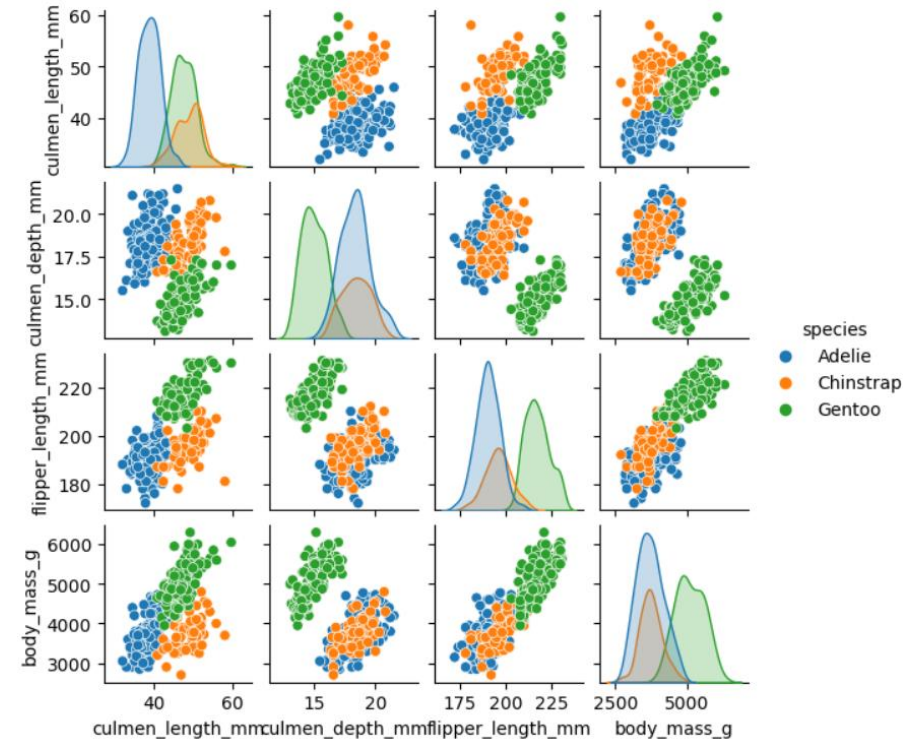
Principal Components Analysis

Example 1: Fisher Iris Data Set and Palmer Penguins Data Set

The following data set has 4 features of **150 Iris flowers**: *sepal length*, *sepal width*, *petal length*, *petal width*. Use PCA to extract 2 principal components then visualize the projected data.



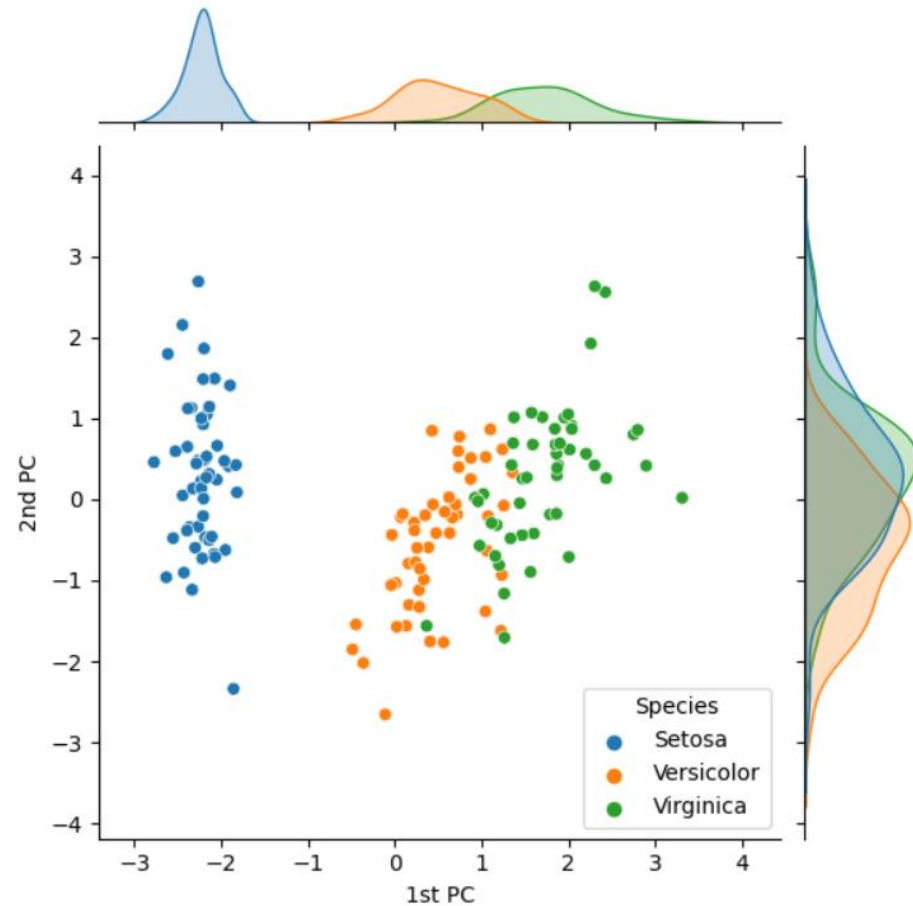
The following data set has 4 features of **345 penguins** from the Palmer Archipelago, Antarctica: *culmen length*, *culmen depth*, *flipper length*, *body mass index*. Use PCA to extract 2 principal components then visualize the projected data.



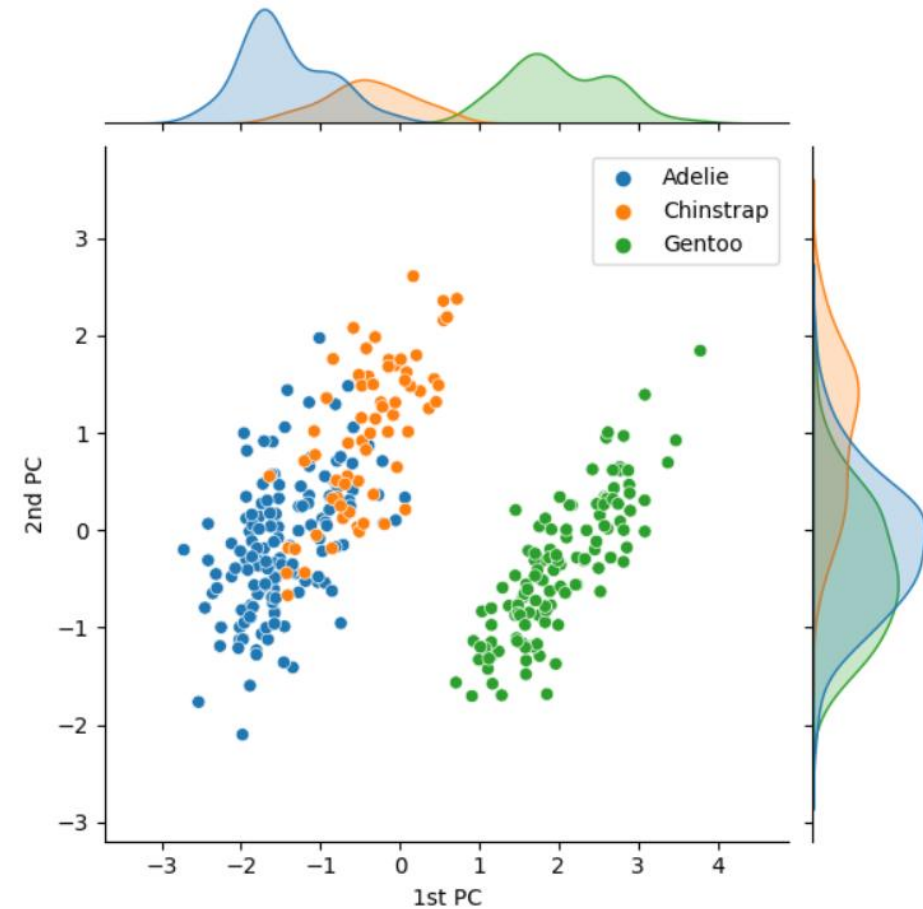
Principal Components Analysis

Example 1: Fisher Iris Data Set and Palmer Penguins Data Set

PCA Result on Fisher Iris Data



PCA Result on Palmer Penguins Data



Principal Components Analysis

How to choose the number of new features to retain (no. of principal components, n)?

It depends on your purpose:

- If you wish to **visualize** data in 2-D or 3-D, choose $n = 2$ or $n = 3$.
- If you have an idea about the **intrinsic dimensionality** of the data (from prior knowledge), then use that as the value of n .
- You can use criteria such as the **CPV** (cumulative percent variance) or the **Scree Plot**.

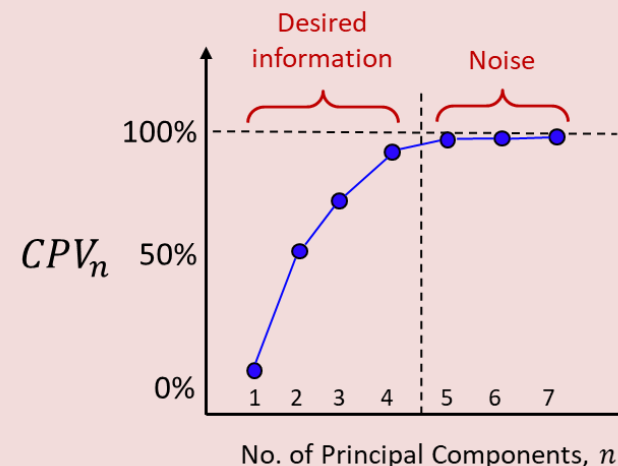
Cumulative Percent Variance (CPV) and Scree Plots

- The sum of all eigenvalues, $\sum \lambda_i$, explains the **total variance** in the data.
- Hence, each cumulative eigenvalue explains the **cumulative variance**.

$$\Sigma = \begin{bmatrix} a_1 & b_1 & c_1 & \cdots \\ a_2 & b_2 & c_2 & \cdots \\ a_3 & b_3 & c_3 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots \\ 0 & \lambda_2 & 0 & \cdots \\ 0 & 0 & \lambda_3 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} a_1 & a_2 & a_3 & \cdots \\ b_1 & b_2 & b_3 & \cdots \\ c_1 & c_2 & c_3 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \quad \boxed{CPV_n = \frac{\sum_{i=1}^n \lambda_i}{\sum \lambda_i}}$$

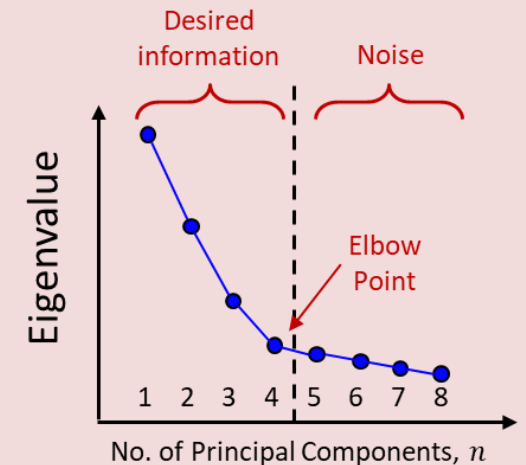
CPV Plot

Choose the first n PC's that cover, say, 95% CPV.



Scree Plot

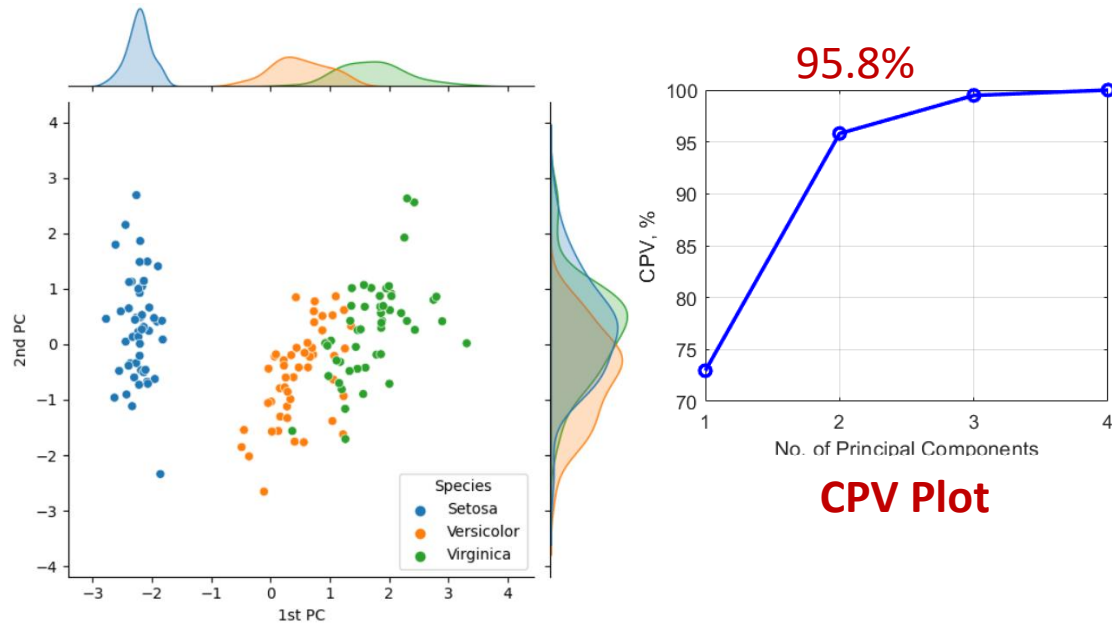
Choose the first n PC's until the elbow point occurs.



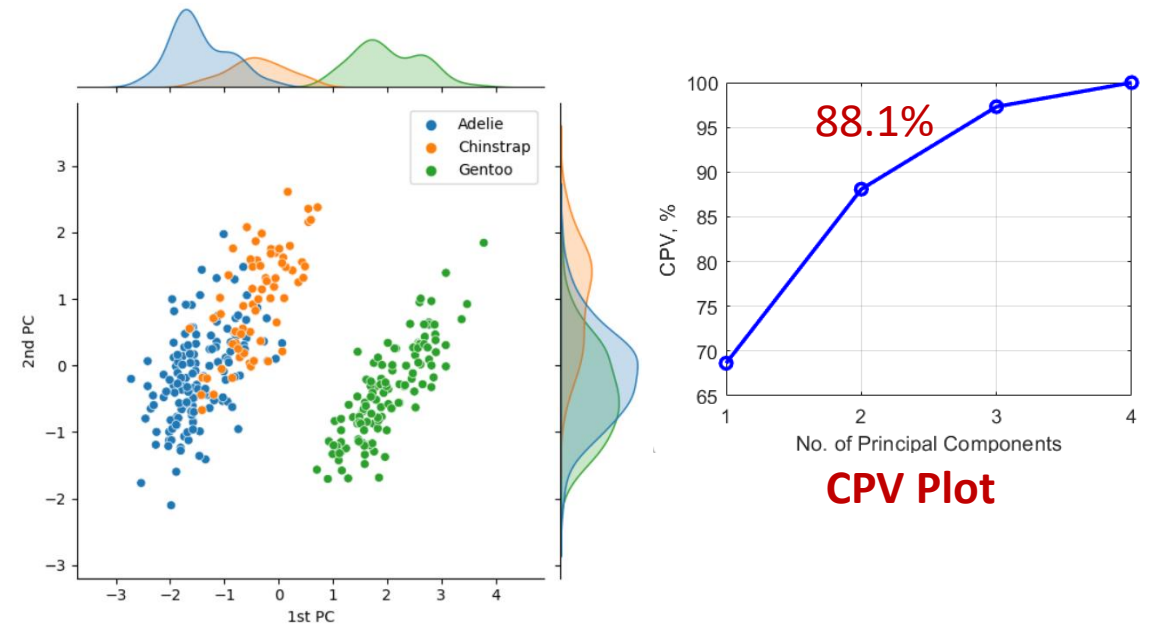
Principal Components Analysis

Example 1: Fisher Iris Data Set and Palmer Penguins Data Set

PCA Result on Fisher Iris Data



PCA Result on Palmer Penguins Data



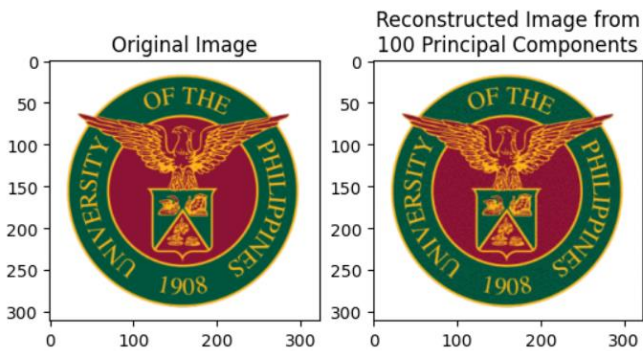
For the same number of principal components, the Fisher Iris data have a higher explained variance than the Palmer Penguins data set. This means that the measurements from the Iris flowers are a lot more correlated. The first 2 principal components already captured 95% of the variation in all 4 features, compared to only 88% in the first 2 principal components of the Palmer Penguins data.

Principal Components Analysis

Other applications of PCA:

Image Reconstruction using Less Information

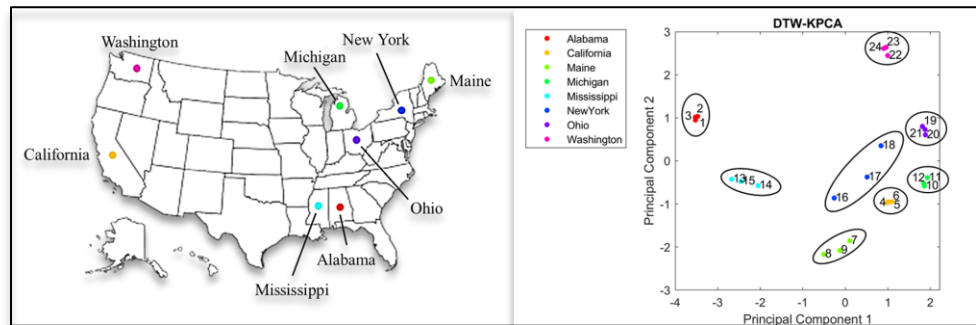
Image compression and reconstruction can be done using eigenfaces.



Discrimination of Substances via PCA on Chemometric Data

Bee substance (propolis) collected from different origins can be traced by applying PCA on their chromatographic data. (Pilario et al., 2022)

<https://www.sciencedirect.com/science/article/pii/S0957417421012926>



Finding Personality Traits using Factor Analysis

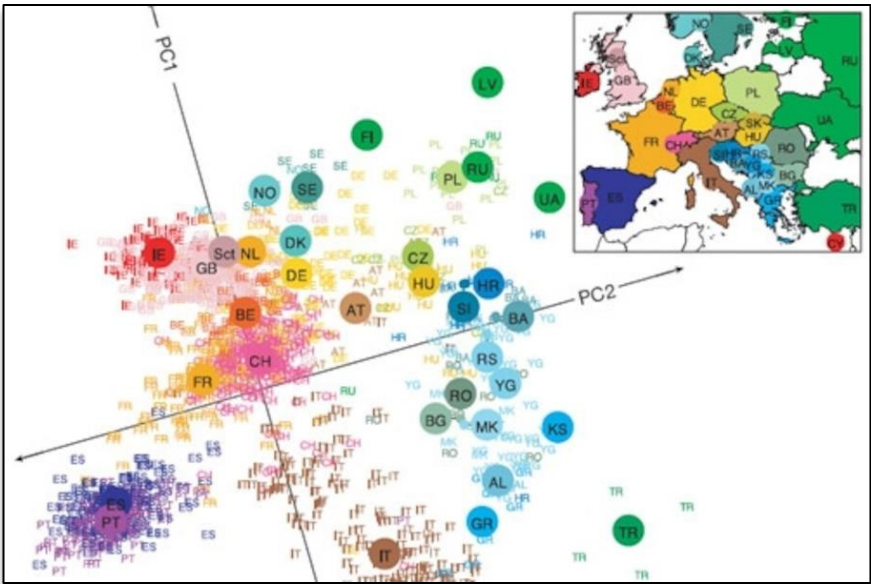
Factor Analysis is a variant of PCA used in Psychology to study personalities.

Trait Dimension	Endpoints of the Dimension	
Conscientiousness	Organized	Disorganized
	Careful	Careless
	Disciplined	Impulsive
Agreeableness	Soft-hearted	Ruthless
	Trusting	Suspicious
	Helpful	Uncooperative
Neuroticism (emotional stability vs. instability)	Calm	Anxious
	Secure	Insecure
	Self-satisfied	Self-pitying
Openness	Imaginative	Practical
	Preference for variety	Preference for routine
	Independent	Conforming
Extraversion	Sociable	Retiring
	Fun-loving	Sober
	Affectionate	Reserved

PCA on Single Nucleotide Polymorphisms (SNPs)

European genes mirror European geography (2008)

<https://www.nationalgeographic.com/science/article/european-genes-mirror-european-geography>



PCA, NMF, ICA

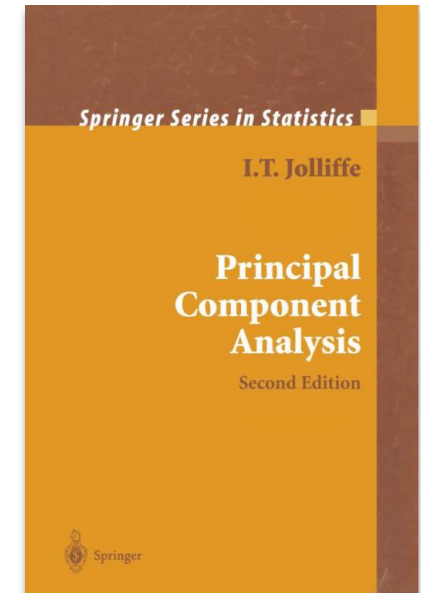
These algorithms are related but they differ in only a few aspects. They each have their own applications.

	Principal Components Analysis (PCA)	Non-negative Matrix Factorization (NMF)	Independent Components Analysis (ICA)
Transformation of Data, X	$Y = \textcolor{red}{X}P$ or $\textcolor{red}{X} = YP^T$	$\textcolor{red}{X} = WH$	$S = W\textcolor{red}{X}$
Requirements	<ul style="list-style-type: none">• Variances of Y are <i>maximized</i>.• P is <i>orthonormal</i>.	<ul style="list-style-type: none">• W and H elements must all be <i>non-negative</i>: $W \geq 0, H \geq 0$	<ul style="list-style-type: none">• S contains <i>maximally independent</i> components.
Cost Function / Objective Func.	<ul style="list-style-type: none">• $\max_P \text{var}(Y)$	<ul style="list-style-type: none">• $\min_{H,W} \ X - WH\ ^2$ (Frobenius norm)• $\min_{H,W} D(X WH)$ (Kullback-Leibler divergence)	<ul style="list-style-type: none">• <i>Maximum Independence</i> can mean:<ul style="list-style-type: none">• Min: Mutual Information• Max: Non-Gaussianity<ul style="list-style-type: none">• Max: Negentropy• Max: Kurtosis
Solver/s	<ul style="list-style-type: none">• Eigenvalue Decomposition (EVD), or,• Singular Value Decomposition (SVD)	<ul style="list-style-type: none">• Fixed-point iteration (Lee and Seung, 2000)<ul style="list-style-type: none">• <i>Multiplicative Update Rules</i>• Any optimization solver	<ul style="list-style-type: none">• EVD or SVD• Fixed-point iteration (Hyvarinen, 2000)• Any optimization solver
Other Related Algorithms / Applications	<ul style="list-style-type: none">• Factor Analysis• Correspondence Analysis• Kernel PCA, Sparse PCA, Robust PCA• Directional Component Analysis• Canonical Correlation Analysis• PARAFAC	<ul style="list-style-type: none">• K-means clustering ($HH^T = I$)• Probabilistic Latent Semantic Analysis (PLSA)• Collaborative Filtering (Recommendation systems)	<ul style="list-style-type: none">• Blind Source Separation• Projection Pursuit• Infomax• FastICA

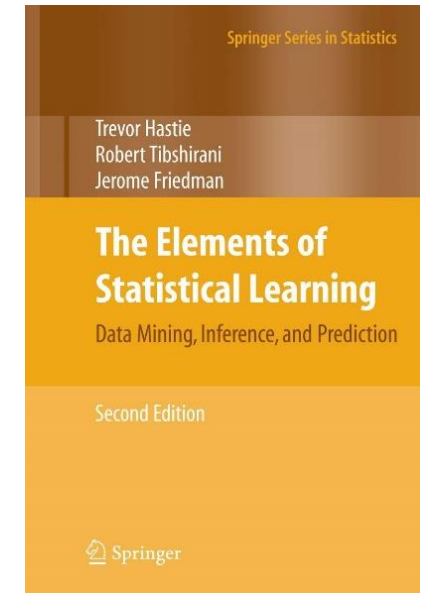
Outline

- Dimensionality Reduction
 - Curse of Dimensionality
 - Feature Selection vs. Feature Extraction
 - Principal Components Analysis (PCA)
 - Derivation
 - Non-negative Matrix Factorization (NMF)
 - Independent Components Analysis (ICA)
- **Low-dimensional Classification**
 - **Linear Discriminant Analysis (LDA)**

Jolliffe (2002)
Principal Component Analysis.
2nd Ed. Springer.



Hastie *et al.* (2008)
The Elements of Statistical Learning.
2nd Ed. Springer.

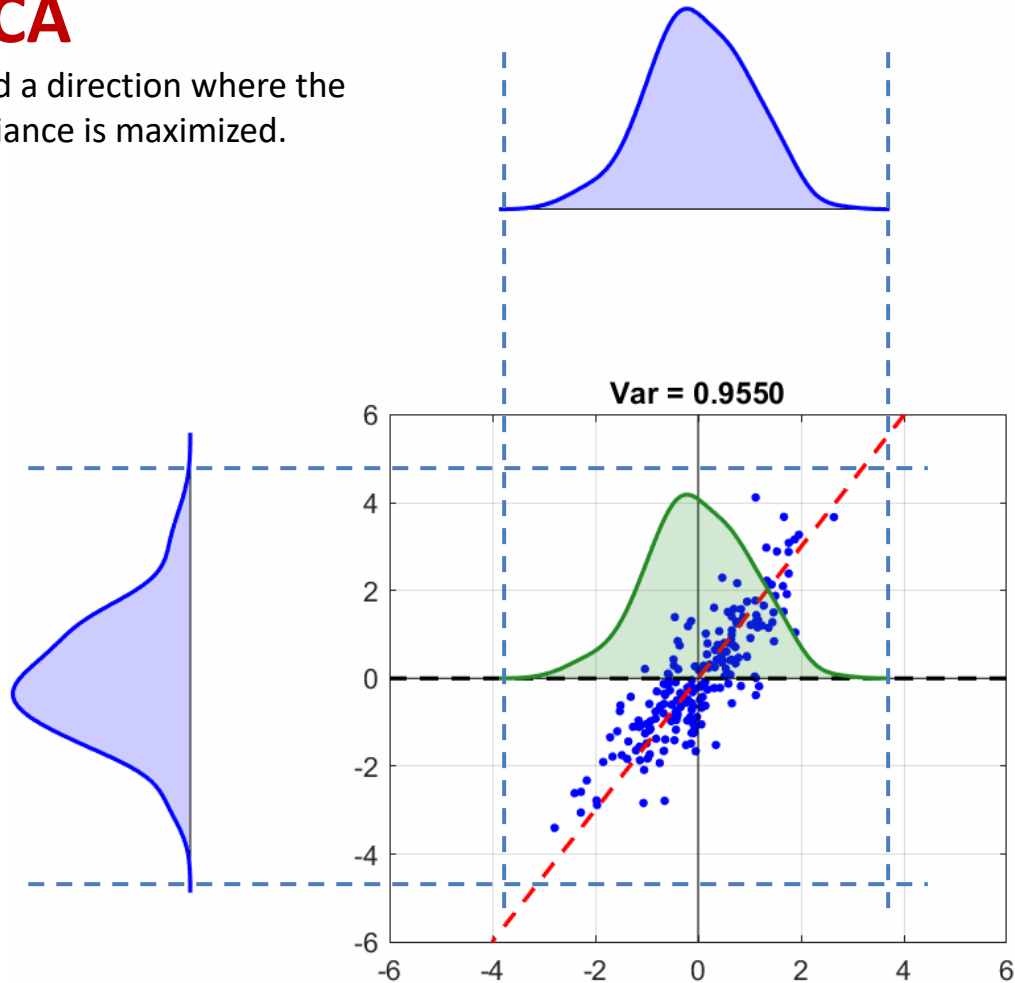


Linear Discriminant Analysis

LDA performs dimensionality reduction, but its main purpose is classification. Since it requires knowledge of class labels, it is a **supervised** learning method.

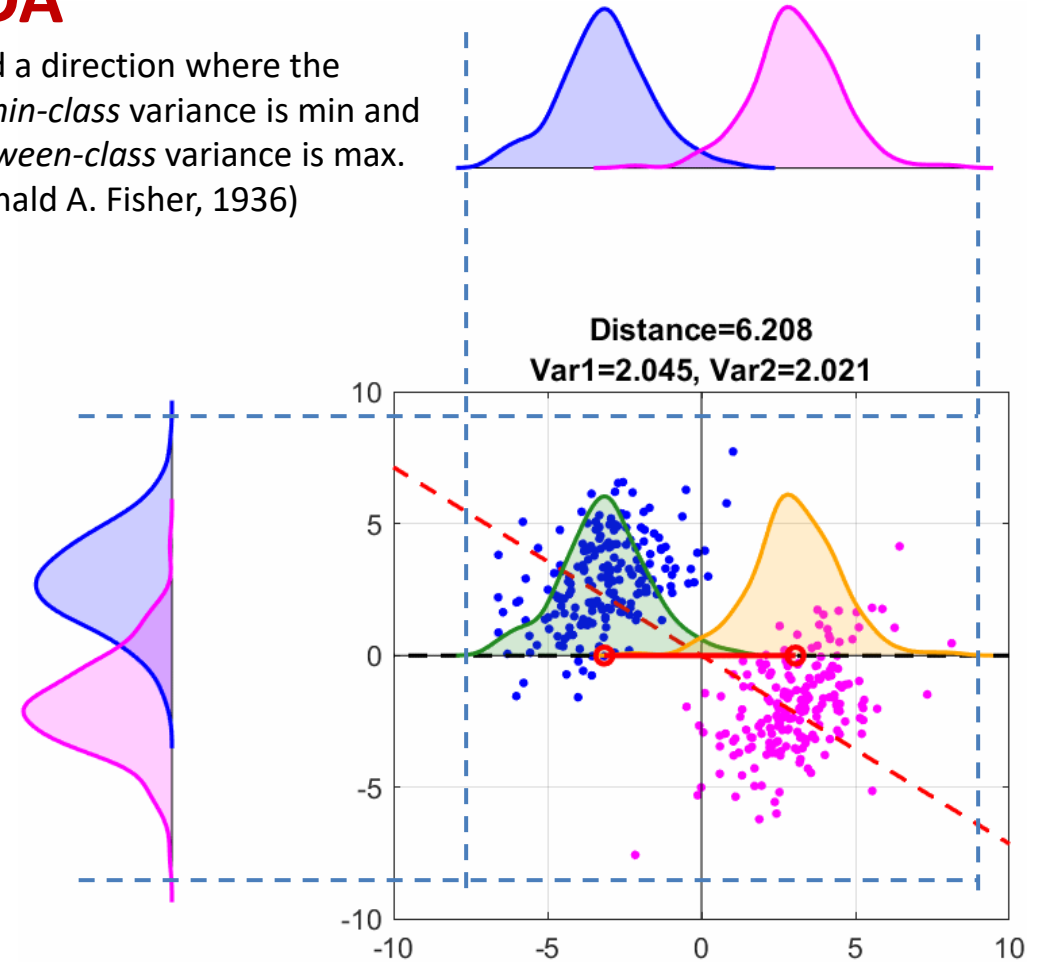
PCA

Find a direction where the variance is maximized.



LDA

Find a direction where the *within-class* variance is min and *between-class* variance is max.
(Ronald A. Fisher, 1936)



Linear Discriminant Analysis

Similar to PCA, the solution to LDA also involves an eigenvalue decomposition.

Given $\mathbf{X} \in \mathbb{R}^{N \times m}$ and class labels $\mathbf{y} \in \mathbb{R}^N$, let:

\mathcal{X}_j = set of samples that belong to class j

p = no. of classes

m = no. of features

N = total no. of samples

$$\mathbf{S}_w = \sum_{j=1}^p \mathbf{S}_j$$

(Within-class
Scatter Matrix)

where:

$$\mathbf{S}_j = \sum_{\mathbf{x}_i \in \mathcal{X}_j} (\mathbf{x}_i - \text{mean}(\mathcal{X}_j)) (\mathbf{x}_i - \text{mean}(\mathcal{X}_j))^T$$

(Within-scatter Matrix for class j)

$$\mathbf{S}_b = \sum_{i=1}^p n_i (\text{mean}(\mathcal{X}_i) - \text{mean}(\mathbf{X})) (\text{mean}(\mathcal{X}_i) - \text{mean}(\mathbf{X}))^T$$

(Between-class
Scatter Matrix)

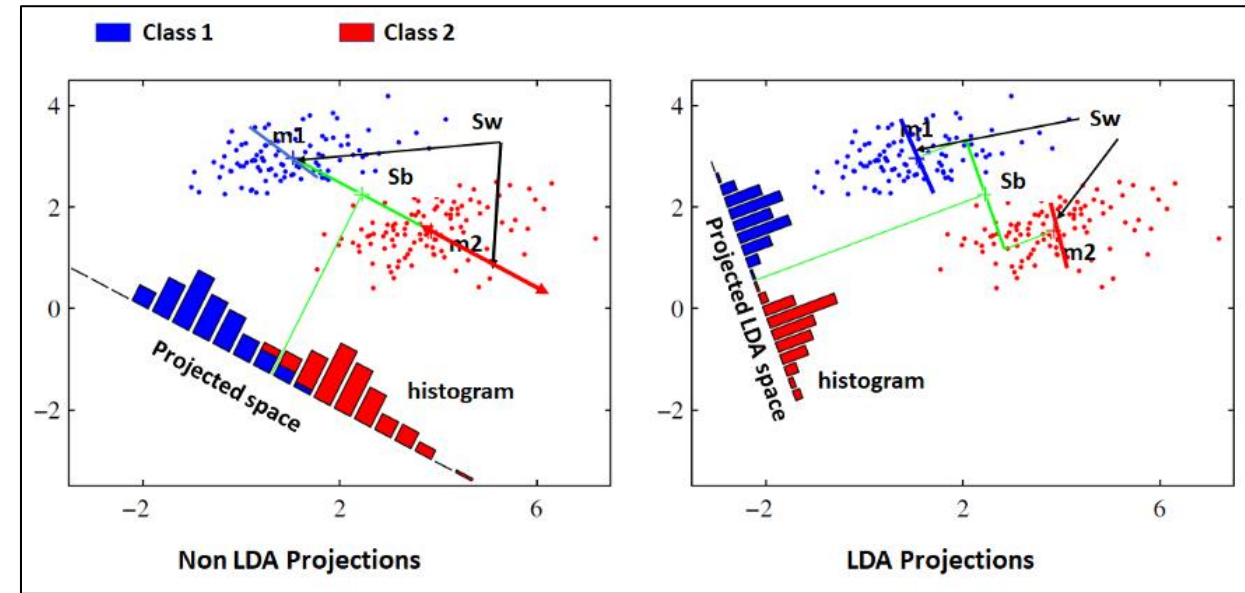
where: n_i = no. of samples in class i

$$\mathbf{S}_t = \sum_{i=1}^N (\mathbf{x}_i - \text{mean}(\mathbf{X})) (\mathbf{x}_i - \text{mean}(\mathbf{X}))^T$$

(Total Scatter
Matrix)

a.k.a. covariance matrix

Note: $\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w$



Source: <https://blog.devgenius.io/part-3-linear-discriminant-analysis-b311fbef7369>

The goal of LDA is to find a projection matrix $\mathbf{W}_r = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_r] \in \mathbb{R}^{m \times r}$, such that the following **Fisher criterion** is maximized:

$$\max_{\mathbf{w} \neq 0} \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

which is equivalent to solving the following EVD:

$$\mathbf{S}_b \mathbf{w}_i = \lambda_i \mathbf{S}_w \mathbf{w}_i$$

The new features extracted from LDA are computed as: $\mathbf{z}_i = \mathbf{W}_r^T \mathbf{x}_i$

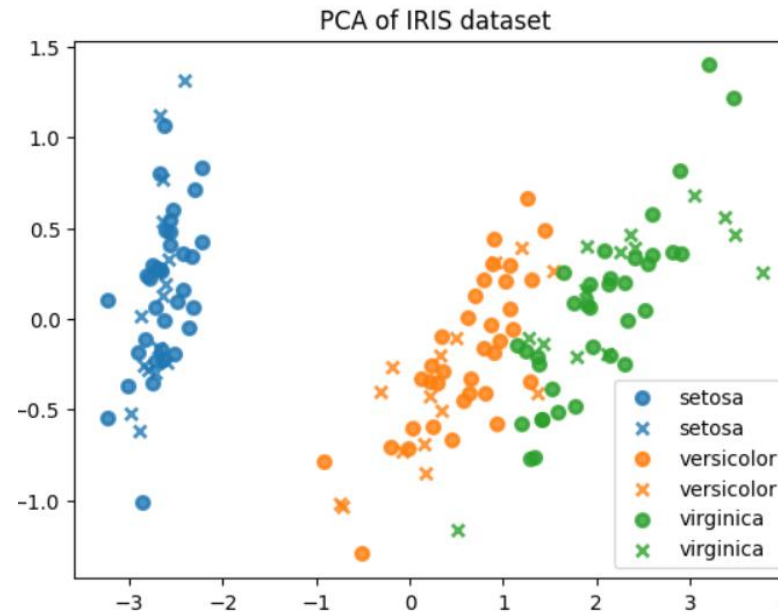
Dimensionality reduction occurs since the data is now projected onto r -dimensional space. However, since the rank of \mathbf{S}_b is less than p (no. of classes), we can only choose r to be at most $p - 1$.

Linear Discriminant Analysis

Example: PCA+SVM vs. LDA on Fisher Iris Data Set

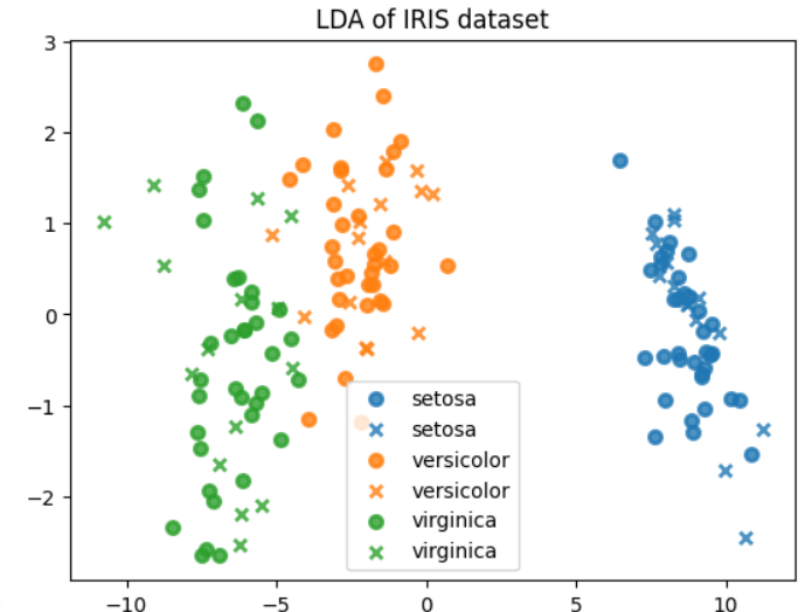
In the Iris Data Set, split the data into 70%-30% train-test with stratification. Project the training data using PCA and LDA, then transform the test data using those projections. Compare the PCA vs. LDA results in 2D.

In addition, compare the result of PCA + SVM vs. LDA in terms of training and testing accuracy.



PCA + SVM Classification

Training Accuracy: 98.10%
Testing Accuracy: 93.33%



LDA Classification

Training Accuracy: 99.05%
Testing Accuracy: 95.56%

Note: PCA and LDA may fail if

- Noise in data are multi-modal or non-Gaussian.
- The underlying manifold is highly nonlinear.

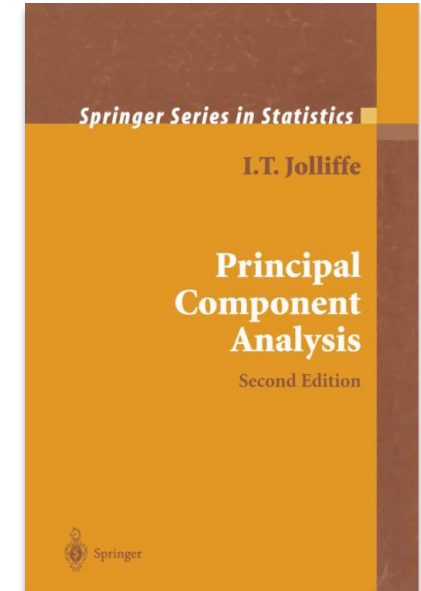
Outline

- Dimensionality Reduction
 - Curse of Dimensionality
 - Feature Selection vs. Feature Extraction
 - Principal Components Analysis (PCA)
 - Derivation
 - Non-negative Matrix Factorization (NMF)
 - Independent Components Analysis (ICA)
- Low-dimensional Classification
 - Linear Discriminant Analysis (LDA)

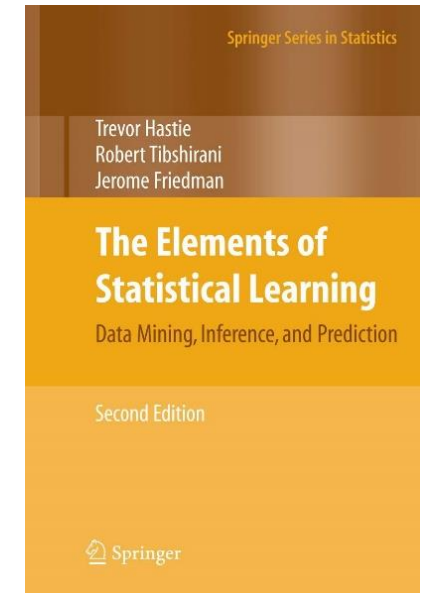
Related methods:

- *PLS (Partial Least Squares)*
 - Low-dimensional Regression
- *NCA (Neighborhood Components Analysis)*
 - Another low-dimensional classifier like LDA.
- *CCA (Canonical Correlation Analysis)*
 - Maximize correlation instead of covariance as in PCA.

Jolliffe (2002)
Principal Component Analysis.
2nd Ed. Springer.



Hastie et al. (2008)
The Elements of Statistical Learning.
2nd Ed. Springer.



Further Reading

- “Curse of Dimensionality.” Bellman R.E. Adaptive Control Processes. Princeton University Press, Princeton, NJ, 1961.
- Chandrashekar and Sahin (2014). A Survey on Feature Selection Methods. *Computers and Electrical Engineering*. Vol. 40, No. 1, 16-28.
<https://www.sciencedirect.com/science/article/pii/S0045790613003066>
- <https://www.izen.ai/blog-posts/feature-selection-filter-method-wrapper-method-and-embedded-method/>
- Ron Kohavi and George H. John (1997). "Wrappers for feature subset selection", *Artificial Intelligence*. Vol. 97, 273-324. doi: 10.1016/S0004-3702(97)00043-X
- Jolliffe and Cadima (2016). Principal component analysis: a review and recent developments.
<http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2015.0202>
- Lee and Seung (2000). Algorithms for Non-negative Matrix Factorization. *Advances in Neural Information Processing Systems*, 13.
<https://papers.nips.cc/paper/2000/file/f9d1152547c0bde01830b7e8bd60024c-Paper.pdf>
- Hyvarinen and Oja (2000). ICA: Algorithms and Applications. *Neural Networks*, 13(4-5): 411-430.
<https://www.cs.helsinki.fi/u/ahyvarin/papers/NN00new.pdf>
- <https://www.kaggle.com/code/tirendazacademy/penguin-dataset-data-visualization-with-seaborn>
- https://scikit-learn.org/stable/auto_examples/decomposition/plot_pca_vs_lda.html
- https://scikit-learn.org/stable/auto_examples/decomposition/plot_pca_iris.html
- https://scikit-learn.org/stable/auto_examples/preprocessing/plot_scaling_importance.html
- Van der Maaten et al. (2009). Dimensionality Reduction: A Comparative Review, *Journal of Machine Learning Research*.
- <https://medium.com/logicai/non-negative-matrix-factorization-for-recommendation-systems-985ca8d5c16c>