

Data Science and Machine Learning in Python

Stephan Weyers

Part 1: Data Science

	Date	Topics covered
1	Apr 13 th	Course introduction Data Science motivation How to use Jupyter Notebook Python types and lists Loops, if/else, functions
2	Apr 20 th	Python tuples, lists, dictionaries Functions Numpy basics, operations Image processing
3	Apr 27 th	Pandas Series, DataFrame Pandas basic operations Import/export files
4	May 4 th	Principles of data visualization Data cleaning and preparation Join, combine and reshape data
5	May 11 th	Volkswahl Bund dataset Data visualization in Python How to write Data Science reports Data aggregation and grouping

Part 2: Machine Learning

	Date	Topics covered
6	Jun 1 st	Introduction to supervised learning Classification and regression scikit-learn k-Nearest Neighbors Linear regression (ridge and lasso)
7	Jun 8 th	Linear classification models Decision trees Random forests and gradient boosting
8	Jun 15 th	Kernel support vector machines Neural networks
9	Jun 22 nd	Introduction to unsupervised learning Preprocessing and scaling Dimensionality reduction Principal component analysis
10	Jun 29 th	k-means clustering Hierarchical clustering DBSCAN
11	Jul 6 th	Representing data Engineering features Model evaluation and improvement Text data analysis

Deadlines for Submission and Distribution of Grading

Student task	Deliverables	Deadline	Work	Share of grade
W01 Assignment	Code and results	Apr 26 th	Team A	5.0%
W02 Case Study	Code / presentation slides	May 22 nd	Team B	18.0%
W02 Case Study	Peer review*	May 31 st	Individual	2.0%
W03 Assignment	Code and results	May 29 th	Team B	5.0%
W04 Assignment	Code and results	Jun 12 th	Team C	10.0%
W05 Assignment	Code and results	Jun 26 th	Team D	7.0%
W06 Assignment	Code and results	Jul 8 th	Team D	13.0%
W07 Case Study	Code / presentation slides	Jul 17 th	Team D	22.0%
W07 Case Study	Peer review*	Jul 31 st	Individual	3.0%
DataCamp 1	Finish course	May 9 th	Individual	2.5%
DataCamp 2	Finish course	May 30 th	Individual	2.5%
DataCamp 3	Finish course	Jun 20 th	Individual	2.5%
DataCamp 4	Finish course	Jul 11 th	Individual	2.5%

* Peer review is mandatory. Quality of peer review itself is graded. Not providing peer review at all would result in high point deduction

Teams for assignment W05-W07

Team	Univ.	Name
D1	UV	Paula Piña
D1	UBA	Facundo Ignacio Zanalda
D1	UBA	Manuel Cabeza Galucci
D1	FHDO	Daniel Tobien
D2	UV	Adonis Nicola Cruz Navarrete
D2	UBA	Manuel Durán
D2	UBA	Lucas Trabanco
D2	FHDO	Bedirhan Abaz
D3	UV	Felipe Galdames
D3	UBA	Victoria Marquez
D3	FHDO	Minh Quan Dinh
D3	ESAN	Juan Jose A. Velasquez Leon
D4	UBA	Gian Franco Lancioni
D4	UBA	Kevin Michalewicz
D4	FHDO	Mohamed Elbaraka
D4	ESAN	Nayely Mayli Ore Ichpas
D5	UV	Nilari Berger Díaz
D5	UBA	Daniel Kundro
D5	UBA	Belen Ticona
D5	FHDO	Celine Cramer
D6	UBA	Francisco Rossi
D6	FHDO	René Frackmann
D6	FHDO	Jessica Heilig
D6	FHDO	Marius Meiners
D7	UV	Valentina Andrea Acuña Ponce
D7	UBA	Sofía Nieva
D7	FHDO	Fabian Herberholt
D7	FHDO	Arnold Urbano Olympio

Team	Univ.	Name
D8	UV	Manuel Orellana Hinojosa
D8	UV	Dian Arriagada
D8	UGTO	Abraham Morales Iturriaga
D8	ESAN	María Ximena Latorre Guzmán
D9	UV	Luis Martinez
D9	UV	Paula Riquelme
D9	UDEM	Jordana L.M. Apolinario Simon
D9	FHDO	Robin Drabon
D10	UV	Jaime Godoy
D10	UDEM	Mariana Gómez Gómez
D10	UBA	Francisco Alan Luna
D10	FHDO	Marco Vom Bovert
D11	UV	Joel Santana
D11	UV	Paula Toro
D11	UBA	Lucía Ailén Kasman
D11	UBA	Rocío Palacín Roitbarg
D11	FHDO	Intissar Boudi
D12	UV	Marcelo Leiton
D12	UV	Emmanuel Cuevas Parra
D12	UBA	Matías Nicolás Pereyra
D12	FHDO	Mamadama Cherif
D12	ESAN	Luiggy Johan Zea Guzman
D13	UV	Dietrich Ganz
D13	UV	Rodrigo Llano Orellana
D13	UBA	Juan Cruz Camacho
D13	FHDO	Justin Skupsch
D13	FHDO	Marco Kusnierek

Team	Univ.	Name
D14	UV	Jorge Rodriguez
D14	UV	Alejandra Valencia
D14	UTTEC	Hugo Isaac Vázquez Gutiérrez
D14	UDEM	Dilan Stiven Correa López
D14	UBA	Andrómeda P. Ovalles Castro
D15	UV	Diego Del Rio
D15	UV	Franco Garrido
D15	UTTEC	José Luís Godínez Vázquez
D15	FHDO	Tegar Fathir Muhammad
D15	ESAN	Jhossy J. Vargas Saldaña
D16	UV	Jose Ignacio Meneses Castillo
D16	UV	Benjamin Serra
D16	UV	Sofia Contreras Figueroa
D16	UGTO	Andrea Rodriguez Sotelo
D16	FHDO	Jakub Bogusz
D17	UV	Amaya Arroyo
D17	UV	Catalina Escobar
D17	UGTO	Frida Martinez Flores
D17	UBA	Victoria Cambriglia
D17	FHDO	Jannick Bröring
D18	UV	Maximiliano Arancibia Santana
D18	UV	Fernando Parada
D18	UGTO	Andrea Ortiz Alvarado
D18	UBA	Joaquin Ceppi
D18	ESAN	Angela Karin Paredes Solano

Agenda for online lecture 9

Session	Topic	Mode	Materials used	Minutes	End
14:30-16:00	Organizational questions	Q&A		10	14:40
	Scaling of Data	Lecture / Q&A	Lecture slides	10	14:50
	Curse of Dimensionality	Lecture / Q&A	Lecture slides	5	14:55
	Projections and PCA	Lecture / Q&A	Lecture slides	20	15:15
	Questions for discussion	Team work in break-out rooms	Lecture slides	25	15:40
	Olivetti Faces	Lecture / Q&A	Lecture slides	10	15:50
16:10-17:40	Olivetti Faces	Lecture / Q&A	Lecture 09a notebook	15	16:25
	Manifold learning, t-SNE	Lecture / Q&A	Lecture slides	5	16:30
	MNIST with PCA / t-SNE	Joint coding in main room	Lecture 09b notebook	60	17:30
17:50-19:20	Questions for discussion	Team work in break-out rooms	Lecture slides	30	18:20
	OCEAN Personality Traits	Lecture / Q&A	Lecture slides	15	18:35
	OCEAN Personality Traits	Lecture / Q&A	Lecture 09c notebook	15	18:50
	Organizational questions	Q&A		10	19:00

Supervised Approaches

- Labeled data
- Target values known

Classification

- Predict category

Regression

- Predict numeric value

Unsupervised Approaches

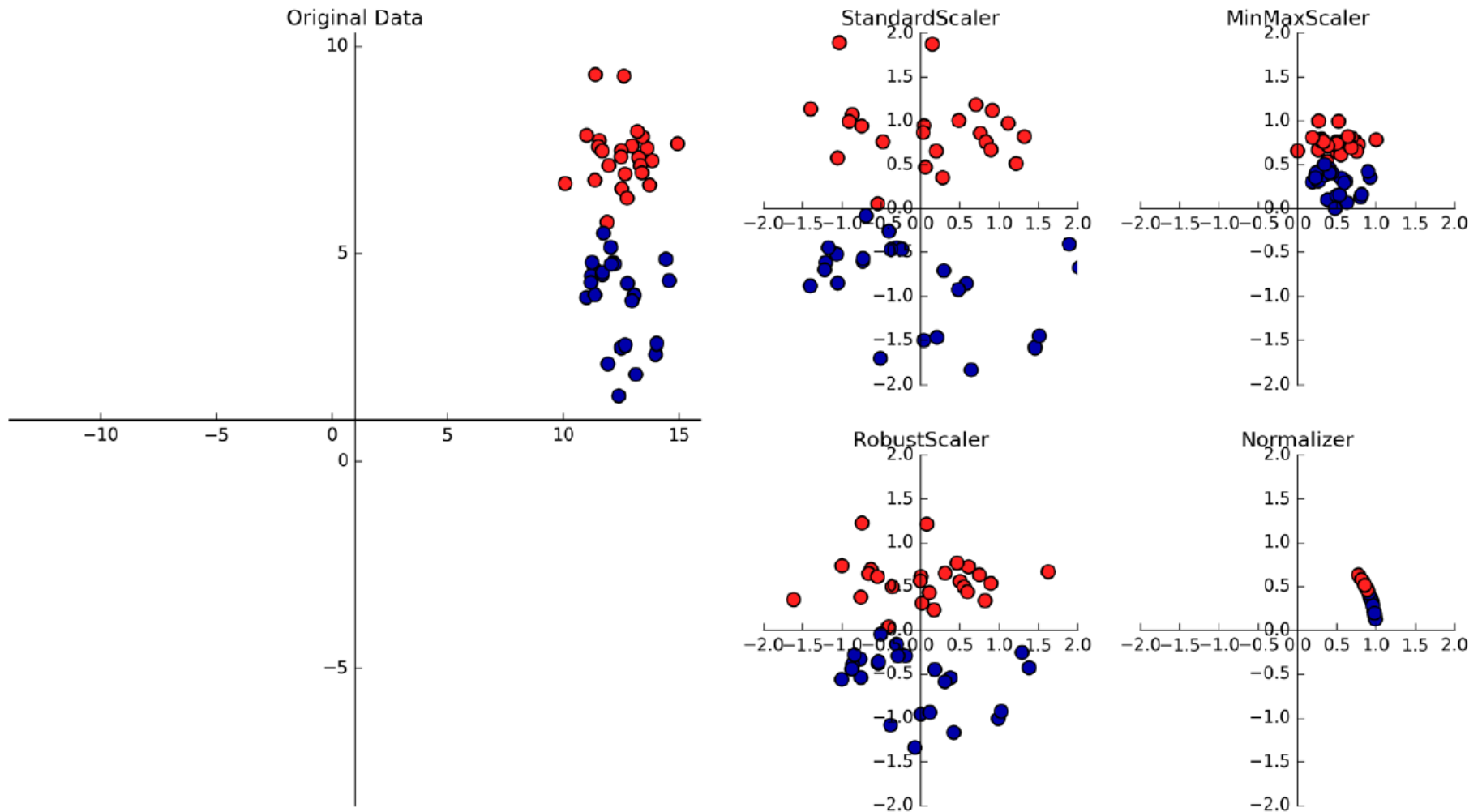
- Unlabeled data
- No target value provided

Cluster Analysis

- Organize similar cases into segments

Dimensionality reduction

- Reduce number of features



```
In [7]: ▶ matrix = [[0,4,-1],[10,6,-1],[20,2,2],[10,4,2]]
df = pd.DataFrame(matrix, columns=["X1","X2","X3"])
df
```

Out[7]:

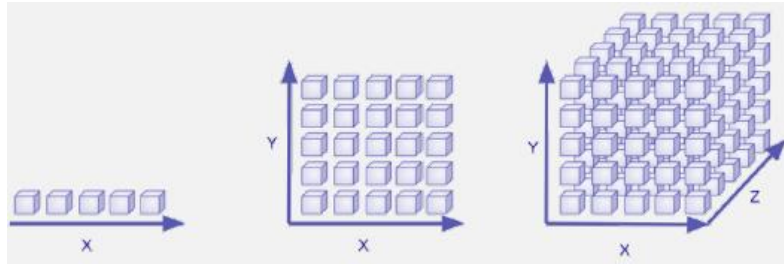
	X1	X2	X3
0	0	4	-1
1	10	6	-1
2	20	2	2
3	10	4	2

```
In [8]: ▶ from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(df)
df_scaled = scaler.transform(df)
df_scaled = pd.DataFrame(df_scaled, columns=df.columns)
df_scaled
```

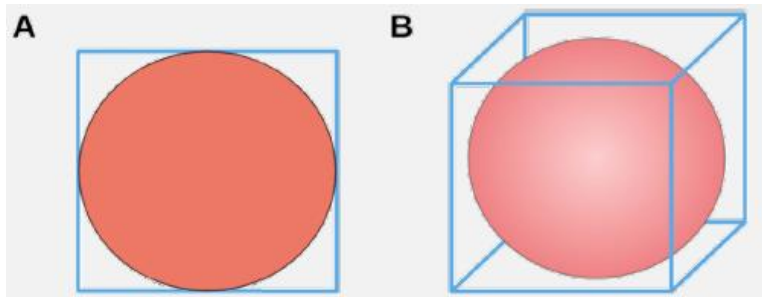
Out[8]:

	X1	X2	X3
0	-1.414214	0.000000	-1.0
1	0.000000	1.414214	-1.0
2	1.414214	-1.414214	1.0
3	0.000000	0.000000	1.0

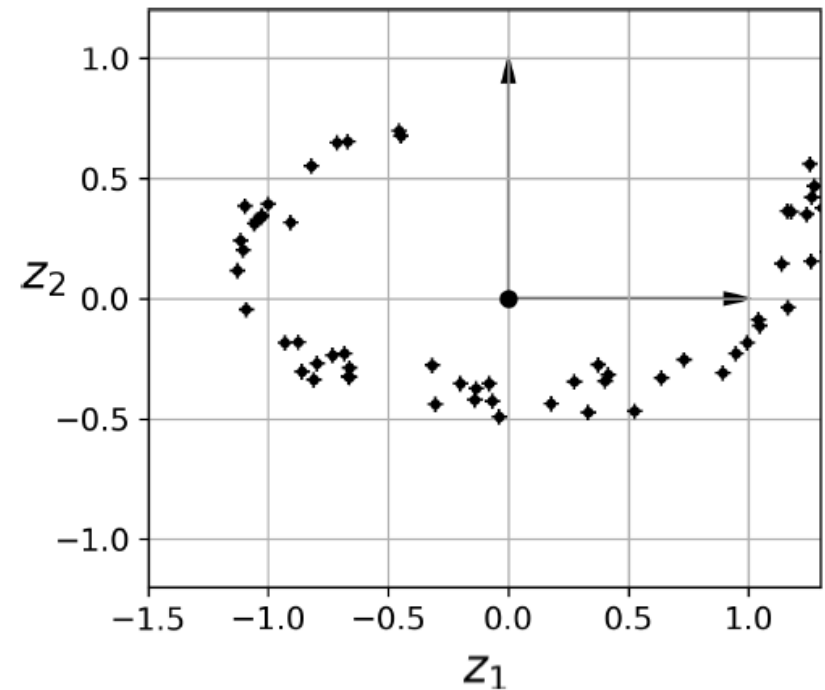
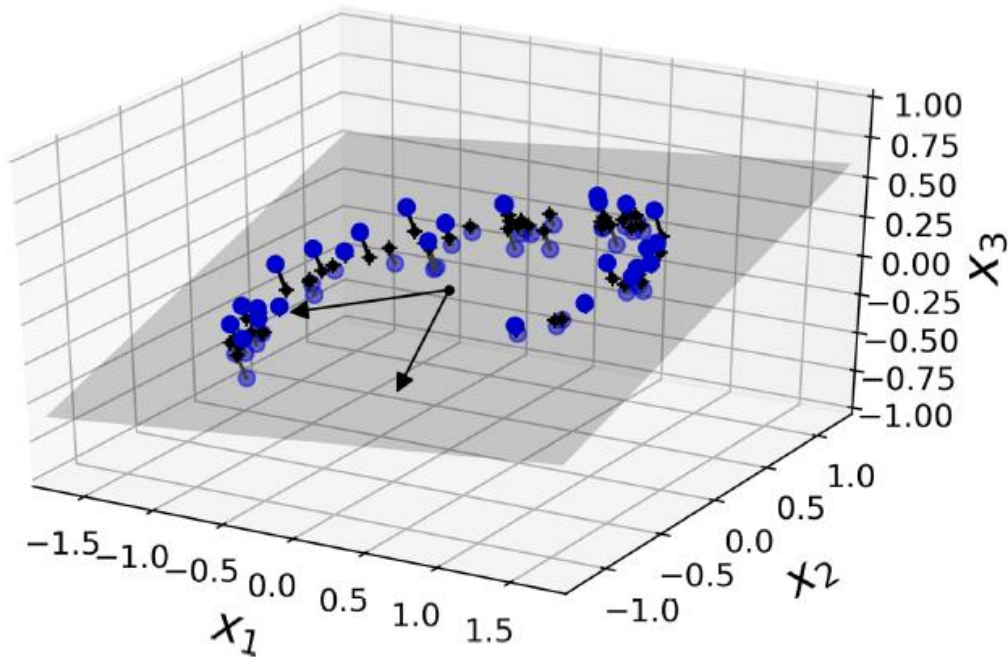
Curse of Dimensionality – Too Many Features

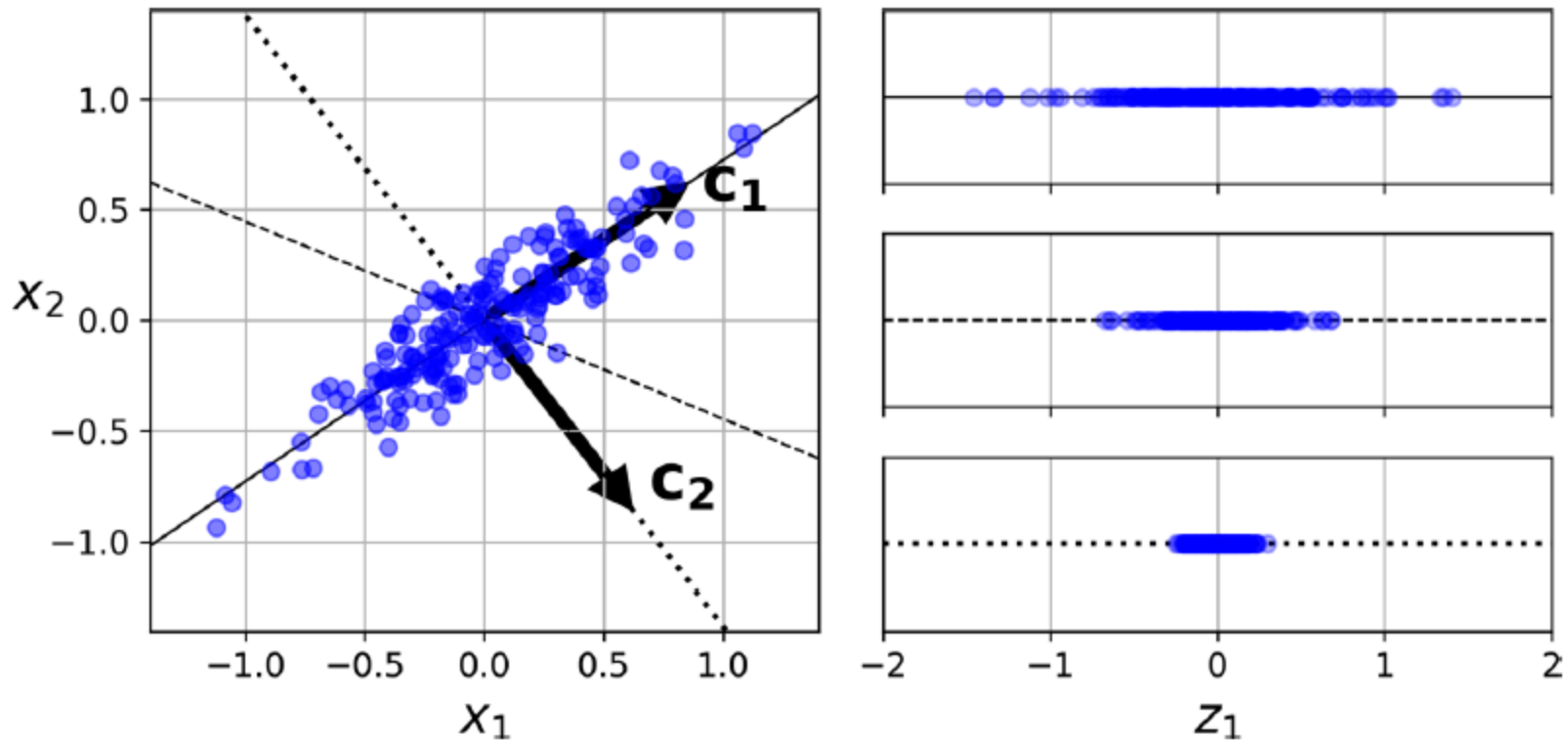


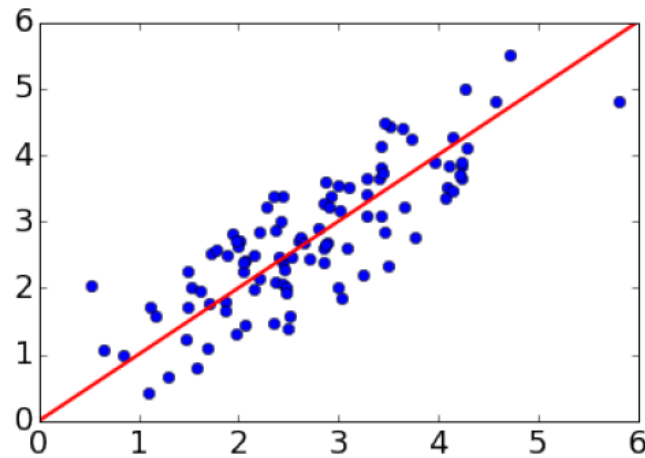
p	# combinations for 5 groups in p-dim
1	5
2	25
3	125
5	3,125
10	9,765,625
15	30,517,578,125



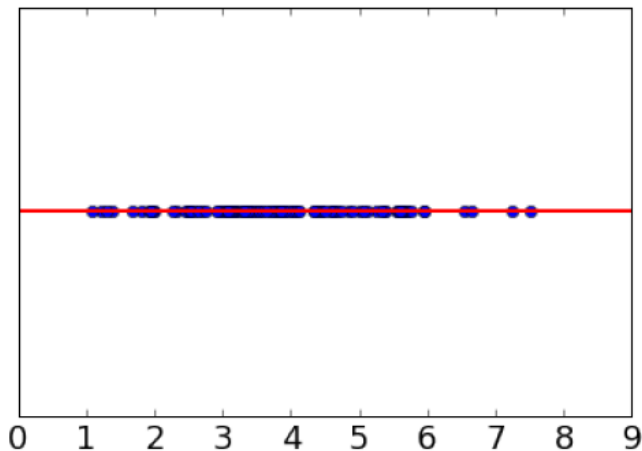
p	p-dim volume of p-ball with diameter 1
1	100.000%
2	78.540%
3	52.360%
4	30.843%
5	16.449%
10	0.249%
15	0.001%



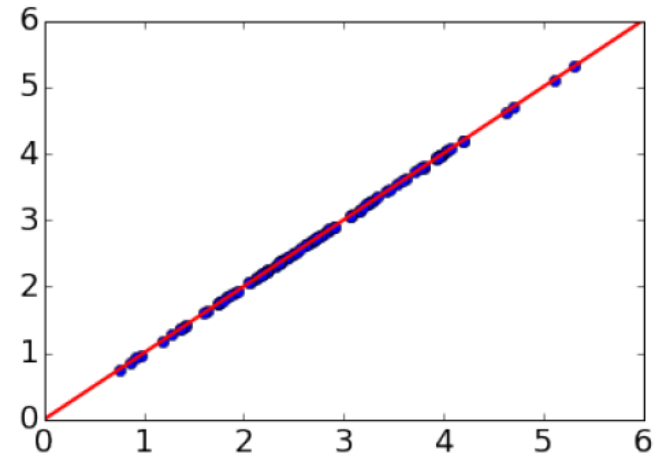




Projection onto \mathbb{R} :



Reconstruction in \mathbb{R}^2 :



Original data

$$\begin{pmatrix} 2 & 4 & -1 \\ -3 & 2 & 0 \\ 8 & -2 & 6 \\ 1 & -3 & 4 \\ 4 & 4 & 2 \end{pmatrix}$$

Adjusted by mean

$$\begin{pmatrix} -0.4 & 3.0 & -3.2 \\ -5.4 & 1.0 & -2.2 \\ 5.6 & -3.0 & 3.8 \\ -1.4 & -4.0 & 1.8 \\ 1.6 & 3.0 & -0.2 \end{pmatrix}$$

Projection

$$\begin{pmatrix} -3.45 & -2.57 & 0.96 \\ -5.41 & 2.37 & -0.35 \\ 7.37 & -0.73 & 0.11 \\ 1.91 & 4.19 & 0.13 \\ -0.42 & -3.27 & -0.87 \end{pmatrix}$$

Mean (2.4 1 2.2)

(0 0 0)

(0 0 0)

Var (13.0 8.8 6.7)

(13.0 8.8 6.7)

(19.8 8.2 0.4)

Principal components in rows

$$\begin{pmatrix} 0.69 & -0.47 & 0.55 \\ -0.66 & -0.72 & 0.21 \\ 0.30 & -0.50 & -0.81 \end{pmatrix}$$

Original data

Projection

Reconstructed data

$$\begin{pmatrix} 2 & 4 & -1 \\ -3 & 2 & 0 \\ 8 & -2 & 6 \\ 1 & -3 & 4 \\ 4 & 4 & 2 \end{pmatrix}$$

$$\begin{pmatrix} -3.45 & -2.57 \\ -5.41 & 2.37 \\ 7.37 & -0.73 \\ 1.91 & 4.19 \\ -0.42 & -3.27 \end{pmatrix}$$

$$\begin{pmatrix} 1.71 & 4.48 & -0.22 \\ -2.90 & 1.83 & -0.28 \\ 7.97 & -1.94 & 6.09 \\ 0.96 & -2.93 & 4.11 \\ 4.26 & 3.56 & 1.30 \end{pmatrix}$$

Mean

$$(2.4 \quad 1 \quad 2.2)$$

Principal components in rows

$$\begin{pmatrix} 0.69 & -0.47 & 0.55 \\ -0.66 & -0.72 & 0.21 \end{pmatrix}$$

data

	X1	X2	X3
0	2	4	-1
1	-3	2	0
2	8	-2	6
3	1	-3	4
4	4	4	2

```
from sklearn.decomposition import PCA  
pca = PCA(n_components=3)  
pca.fit(data)
```

```
PCA(n_components=3)
```

```
pca.components_
```

```
array([[ 0.69, -0.47,  0.55],  
       [-0.66, -0.72,  0.21],  
       [ 0.3 , -0.5 , -0.81]])
```

```
pca.transform(data)
```

```
array([[ -3.45, -2.57,  0.96],  
       [-5.41,  2.37, -0.35],  
       [ 7.37, -0.73,  0.11],  
       [ 1.91,  4.19,  0.13],  
       [-0.42, -3.27, -0.87]])
```

```
pca.explained_variance_ratio_
```

```
array([0.7 , 0.29, 0.01])
```

data

	X1	X2	X3
0	2	4	-1
1	-3	2	0
2	8	-2	6
3	1	-3	4
4	4	4	2

```
pca = PCA(n_components=2).fit(data)
print(pca.components_)
print("\n")
print(pca.transform(data))
print("\n")
print(pca.inverse_transform(pca.transform(data)))
```

```
[[ 0.69 -0.47  0.55]
 [-0.66 -0.72  0.21]]
```

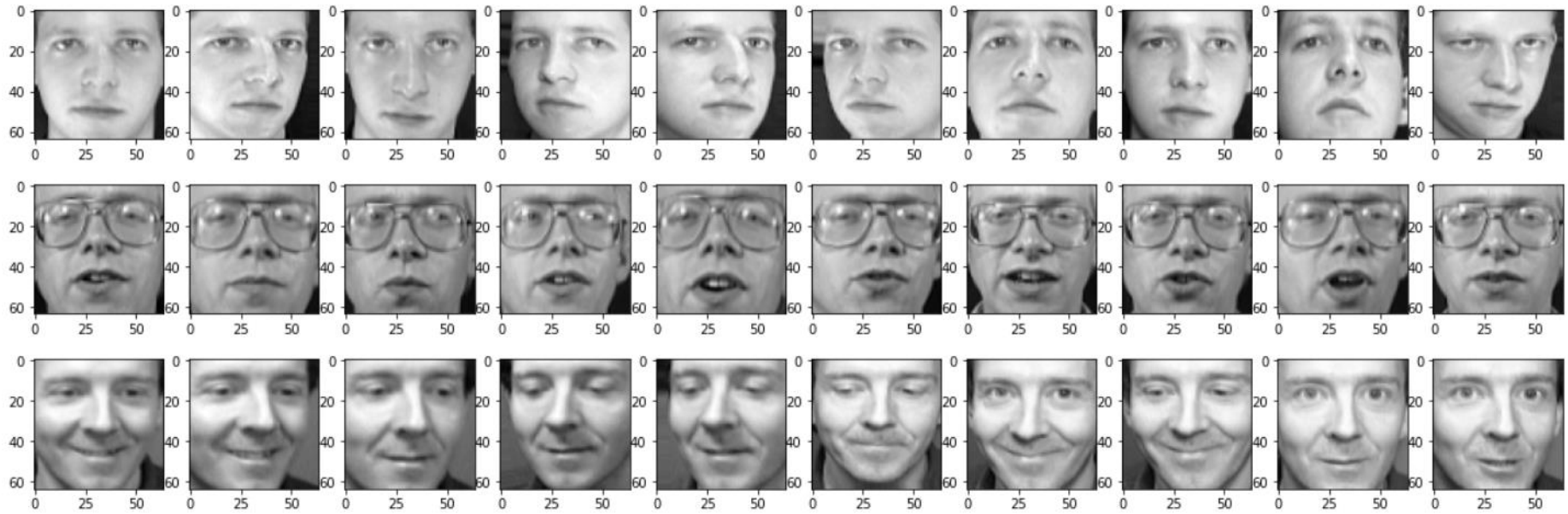
```
[[ -3.45 -2.57]
 [-5.41  2.37]
 [ 7.37 -0.73]
 [ 1.91  4.19]
 [-0.42 -3.27]]
```

```
[[ 1.71  4.48 -0.22]
 [-2.9   1.83 -0.28]
 [ 7.97 -1.94  6.09]
 [ 0.96 -2.93  4.11]
 [ 4.26  3.56  1.3  ]]
```

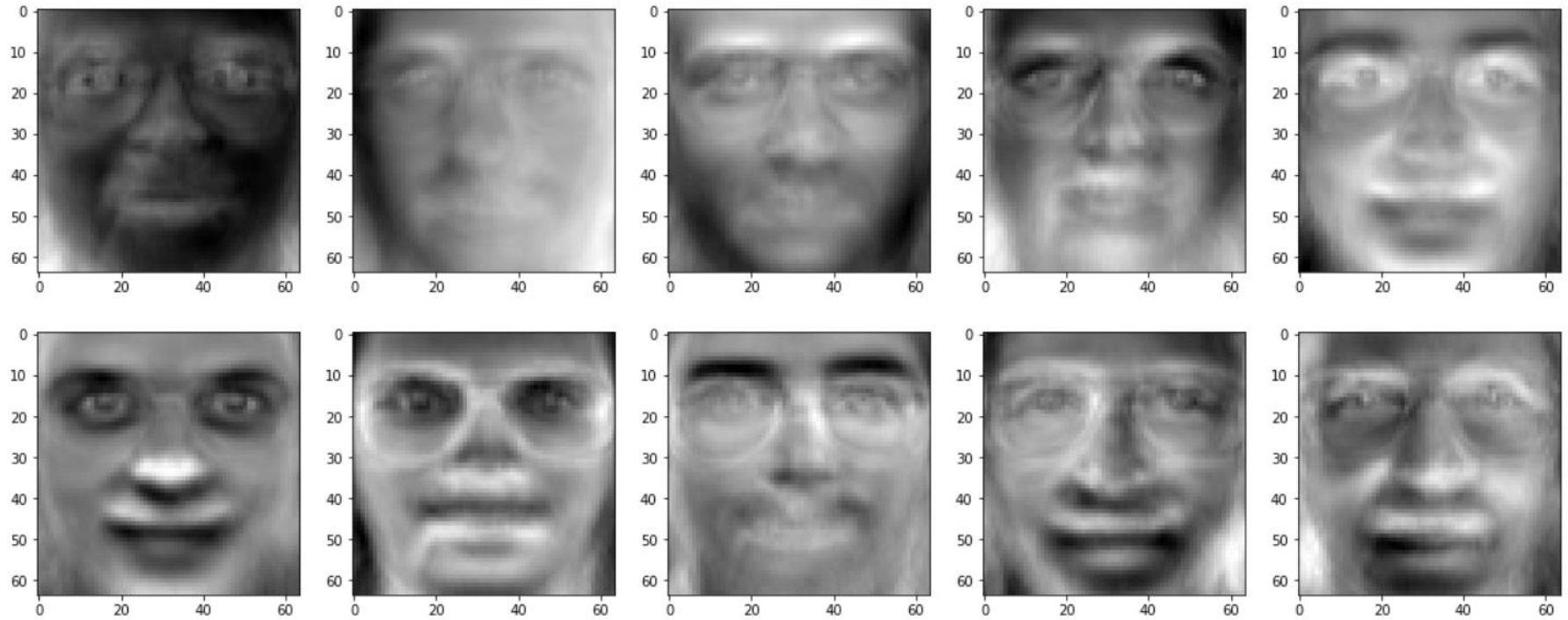

Questions for discussion

1. What are the main motivations for reducing a dataset's dimensionality? What are the main drawbacks?
2. What is the curse of dimensionality?
3. Once a dataset's dimensionality has been reduced, is it possible to reverse the operation? If so, how? If not, why?

Olivetti Faces – Dataset



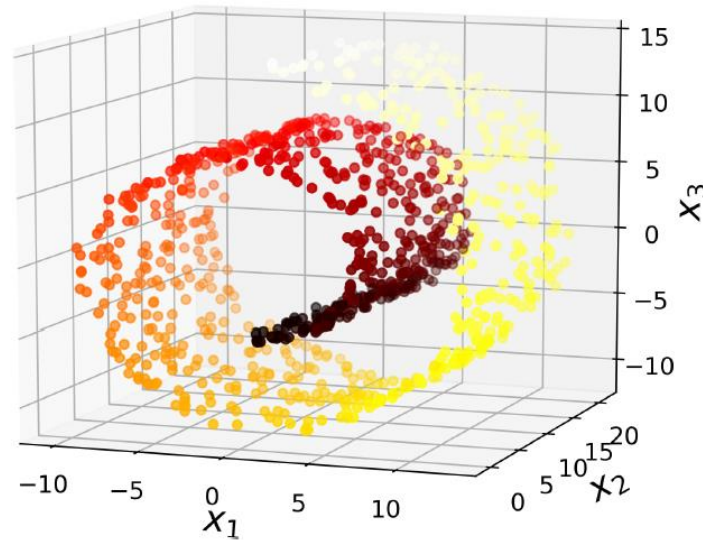
Olivetti Faces – EigenFaces



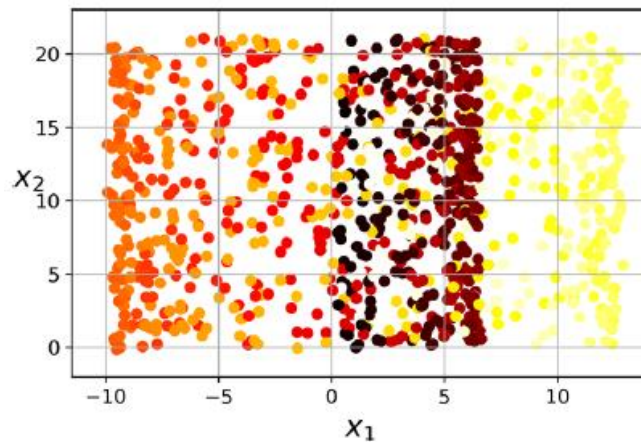
Olivetti Faces – Reconstructed



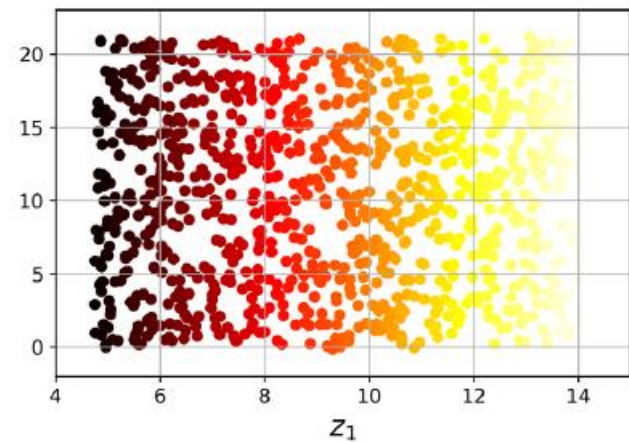
Original data



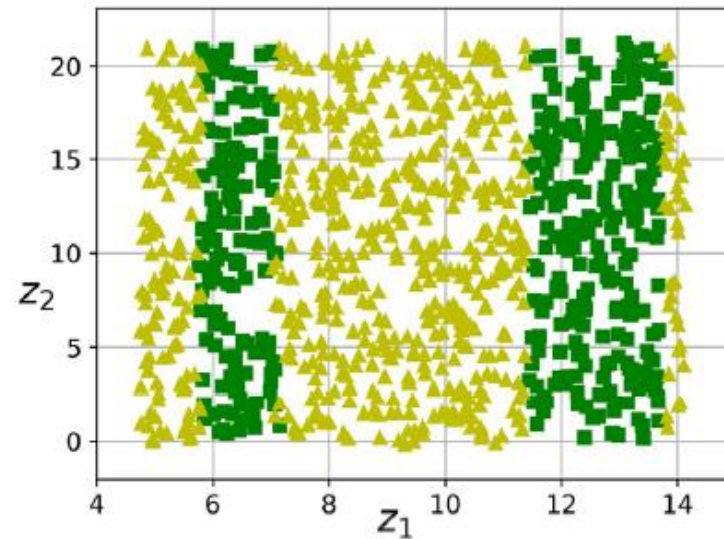
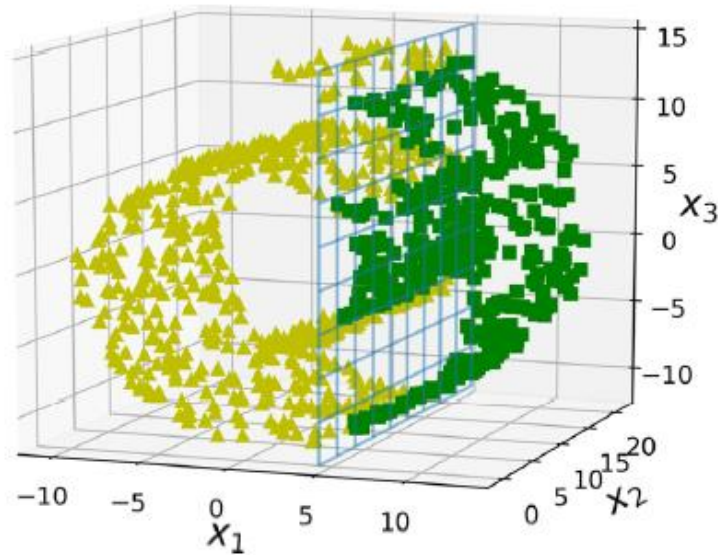
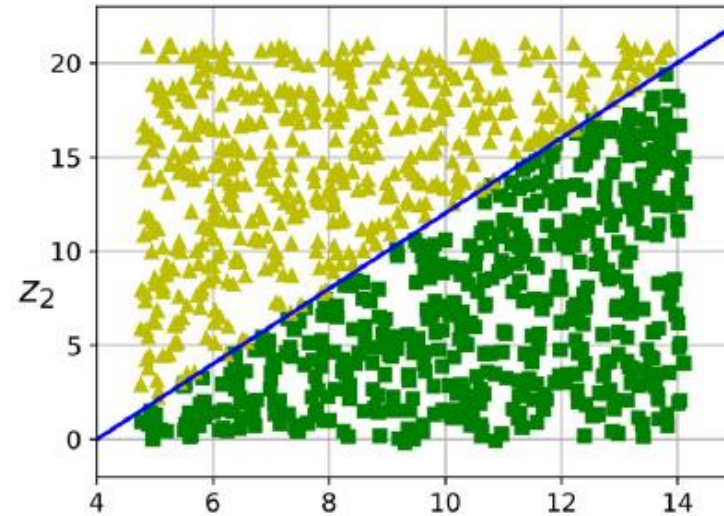
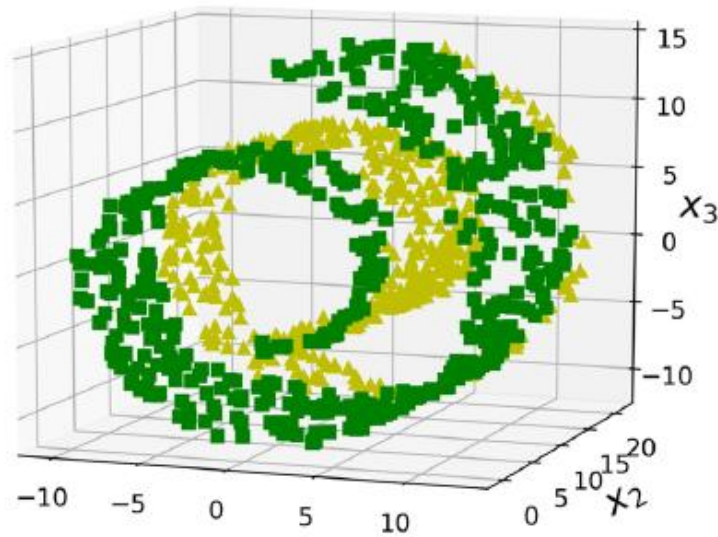
Projection



Unrolling manifold



The decision boundary may not always be simpler with lower dimensions



Questions for discussion

1. Can PCA be used to reduce the dimensionality of a highly nonlinear dataset?
2. Suppose you perform PCA on a 1,000-dimensional dataset, setting the explained variance ratio to 95%. How many dimensions will the resulting dataset have?
3. How can you evaluate the performance of a dimensionality reduction algorithm on your dataset?
4. Does it make any sense to chain two different dimensionality reduction algorithms?

- Lexical hypothesis: most important personality characteristics have become encoded in natural language.
- Allport and Odbert (1936): identified 4500 words describing personality traits.
- Group these words into (approximate) synonyms, by manual clustering
- Data collection: Ask persons whether these words describe them.

Spirit
Talkativeness
Sociability
Spontaneity
Boisterousness
Adventure
Energy
Conceit
Vanity
Indiscretion
Sensuality

Jolly, merry, witty, lively, peppy
Talkative, articulate, verbose, gossipy
Companionable, social, outgoing
Impulsive, carefree, playful, zany
Mischievous, rowdy, loud, prankish
Brave, venturesome, fearless, reckless
Active, assertive, dominant, energetic
Boastful, conceited, egotistical
Affected, vain, chic, dapper, jaunty
Nosy, snoop, indiscreet, meddling
Sexy, passionate, sensual, flirtatious

	shy	merry	tense	boastful	forgiving	quiet
Person 1	4	1	1	2	5	5
Person 2	1	4	4	5	2	1
Person 3	2	4	5	4	2	2
		⋮				

Correlation of first five principle components with personality traits

Extraversion

- : quiet (-.83), reserved (-.80), shy (-.75), silent (-.71)
- + : talkative (.85), assertive (.83), active (.82), energetic (.82)

Agreeableness

- : fault-finding (-.52), cold (-.48), unfriendly (-.45), quarrelsome (-.45)
- + : sympathetic (.87), kind (.85), appreciative (.85), affectionate (.84)

Conscientiousness

- : careless (-.58), disorderly (-.53), frivolous (-.50), irresponsible (-.49)
- + : organized (.80), thorough (.80), efficient (.78), responsible (.73)

Neuroticism

- : stable (-.39), calm (-.35), contented (-.21)
- + : tense (.73), anxious (.72), nervous (.72), moody (.71)

Openness

- : commonplace (-.74), narrow (-.73), simple (-.67), shallow (-.55)
- + : imaginative (.76), intelligent (.72), original (.73), insightful (.68)