

# INT104 ARTIFICIAL INTELLIGENCE

## REVIEW – WEEK 14

Sichen Liu

Sichen.Liu@xjtlu.edu.cn



1

### Data Cleaning

- Handling Missing Data
  - Get rid of the corresponding instance.
  - Get rid of the whole column.
  - Set the values to some value (zero, the mean, the median, etc.).
- Smooth Noisy Data
  - Identify or remove the outliers
  - Try to resolve the inconsistent (there is no one way to remove noise, or smooth out the noisiness in the data)



3

### Data Transformation

- Normalization
  - Min-max normalization.

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- Z-score normalization.

Normalizing every value in a dataset such that the mean of all of the values is 0 and the standard deviation is 1

$$x_{scaled} = \frac{x - mean}{sd}$$

- Normalization by decimal scaling.

$$x_{scaled} = \frac{x}{10^j}$$



5

### Data Collection



Lots of places that host/share data online, or you can collect them yourself.



Open data collections



Social media data



Multimodal data



2

### Practice: Data Cleaning

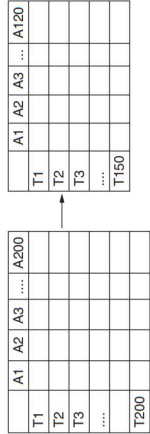
#	Country	Alcohol (L/person)	Deaths (Per 100k)	Heart (Per 100k)	Liver (Per 100k)	Free healthcare
1	Australia	2.5	785	211	15.3000019	Y
2	Austria	3.00000095	863	167	45.5999847	Y
3	Belg/Lux	2.90000095	883	131	20.7000076	N
4	Canada	2.40000095	793	NA	16.3999962	Y
5	Denmark	2.90000095	971	220	23.8999962	Y
6	Finland	0.80000012	970	297	19	N
7	France	9.10000381	751	11	37.9000153	N
8	Iceland	-0.80000012	743	211	11.1999981	Y
9	Ireland	0.69999988	1000	300	6.5	Y
10	Israel	0.60000024	-834	183	13.6999981	Y
11	Italy	27.90000095	775	107	42.2000076	Y
12	Japan	1.5	680	36	23.2000076	N
13	Netherlands	1.79999952	773	167	9.19999809	N
14	New Zealand	1.89999976	916	266	7.69999809	Y
15	Norway	0.080000012	806	227	12.1999981	N
16	Spain	6.5	724	NA	NA	Y
17	Sweden	1.60000024	743	207	11.1999981	N
18	Switzerland	5.80000191	693	115	20.2999924	N
19	UK	1.29999952	941	285	10.3000019	Y
20	US	1.20000048	926	199	22.1000038	N
21	West Germany	2.70000048	861	172	36.7000076	Y



4

### Data Reduction

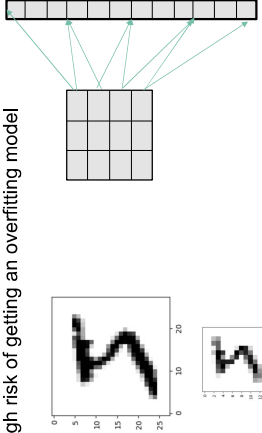
Data reduction is a key process in which a reduced representation of a dataset that produces the same or similar analytical results is obtained.



6

## Dimensionality Reduction

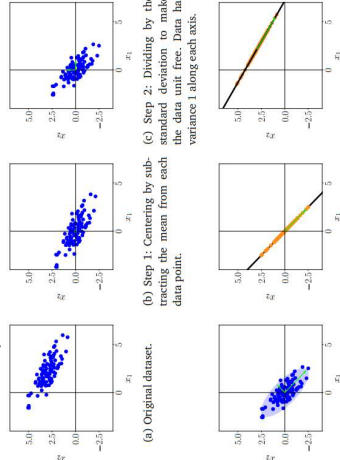
- Data with high dimensions:
- High computational complexity
  - May contain many irrelevant or redundant features
  - Difficulty in visualization
  - With high risk of getting an overfitting model



7

## PCA

Key steps of PCA in practice



(d) Step 3: Compute eigenvalues and eigenvectors (arrows) of the data covariance matrix to make the principal subspace.  
(e) Step 4: Project data onto the principal subspace.  
(f) Undo the standardization and move projected data back into the original data space from (e).



9

## API Information

[sklearn.cluster.KMeans](#)

class sklearn.cluster.KMeans(n\_clusters=8, random\_state=None): KMeans clustering.

### Parameters

**n\_clusters:** int, default=8

The number of clusters to form as well as the number of centroids to generate  
**random\_state:** int, RandomState instance or None, default=None  
Determines random number generation for centroid initialization. Use an int to make the randomness deterministic.

### Attributes

**labels\_:** ndarray of shape (n samples,). Labels of each point.

### Methods

**fit(X):** Compute k-means clustering.

Parameters

**X:** array-like, sparse matrix of shape (n samples, n features)

Training instances to cluster.

Returns

**self:** object

Fitted estimator.



11

## Principal Component Analysis (PCA)

Preserving the Variance:

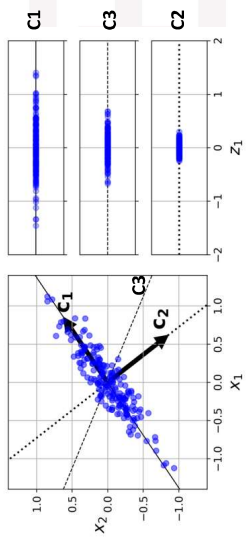


Figure 8-7. Selecting the subspace to project on

PCA identifies the axis that accounts for the largest amount of variance in the training set.



8

## API Information

[numpy.linalg.norm](#)

linalg.norm(x, axis=None): Matrix or vector norm.

### Parameters

**x:** array like

Input array. If axis is None, x must be 1-D or 2-D.

### Returns

**n:** float or ndarray

Norm of the matrix or vector(s).



10

## Indentation matters!

- Code is grouped by its indentation
- Indentation is the number of whitespace or tab characters before the code.
- If you put code in the wrong block, then you will get unexpected behavior

```
Line 1 x = True
Line 2 if x:
Line 3     print("Executing if")
Line 4 else:
Line 5     print("Executing else")
Line 6     print("Prints regardless of the if-else block")

Executing if
Prints regardless of the if-else block
```



12