



Data Management for Digital Health
Revision of Exercise IV

Borchert, Rasheed, Dr. Bayat, Dr. Schapranow

Data Management for Digital Health

Winter 2022

Exercise III

Topics

- Medical Use Case Infectious Diseases
- Medical Imaging
- Deep Learning
- Unsupervised Learning

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
2

Exercise III

Key Stats

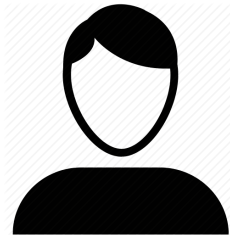
25 Questions
50 Points

40
Submissions
40 Passed

Average score
43.55 / 87.1%

Average time
57.6 min

<< 3h



Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022

3

Q3: What are effective measures to generally prevent spread of infectious diseases especially if you have no treatment at hand?

- ✗ Wearing a mask
- ✗ Use of high-dose antimycotics and antibiotics combinations to prevent additional infections
- ✓ Systematic contact tracing of infected persons
- ✓ Monitoring of persons suspected to be in close contact with infected persons for symptoms over the incubation time.

Frequently missed

Frequent incorrect answer

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022

4

What To Take Home?

- Infections may
 - Affect different body locations
 - Be triggered by numerous agents
 - Result in life-threatening events
 - Require intensive care
- If no therapy available: containment is the **only option** to fight a pandemic spread.



Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
5

Infectious Diseases: Definitions

- **Infection** := Invasion + multiplication of disease-causing agents + host reaction
- **Disease-causing agents** := Pathogenic microorganisms, e.g. bacteria, viruses, parasites or fungi
- Infectious diseases can be spread directly or indirectly from one person to another
- **Zoonotic diseases** := Infectious diseases transmitted from animals to humans or vice versa.

Evaluation Exercise IV

Q4: Please select all correct options for vaccinations as discussed in class.

- Inactivated vaccines are more efficient than mRNA-based vaccines.
- Vector-based vaccines make use of a non-pathogenic/disabled vector to invade body cells and release the disease-specific RNA there.
- Due to the high resistance of RNA molecules, mRNA-based vaccines can be stored at room temperature over a long period of time (months and year) before its use (long durable vaccines)
- A single vaccination dose is sufficient to obtain life-long immunity for most infectious diseases.

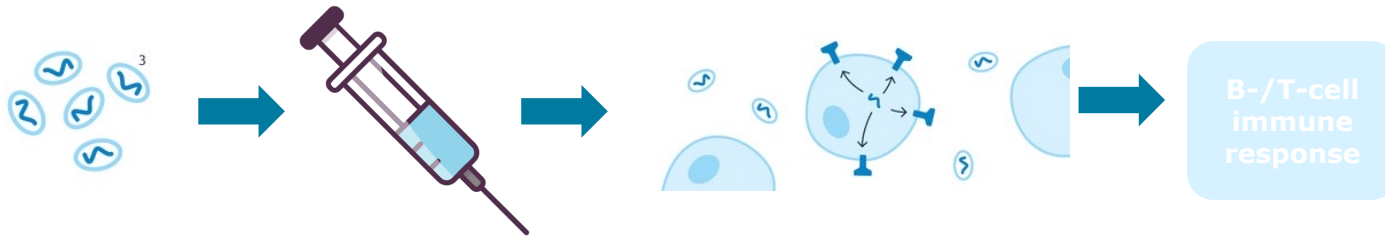
Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022

7

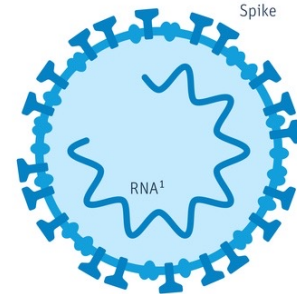
Types of Vaccination: mRNA Vaccine

- Requires synthesis and handling of very fragile RNA



- mRNA for spike protein covered by lipid hull

- Lipid hull allows to mRNA to enter body cells
- Ribosomes read RNA and assemble spike protein
- Proteins are released by cell



Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022

8

Q13: Which statements are true about unsupervised learning? Select all that apply

- Expert annotators are required to define baseline parameters.
- Clustering is a type of unsupervised learning.
- Explicit labels are required for unsupervised learning.
- Unsupervised learning always needs the number of clusters as input.

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
9

Most common problem settings in Machine Learning

Supervised Learning (Labels available for training)

Classification

Categorical output

e.g. $x \in \text{Fruits}$, $y \in \{\text{"apple"}, \text{"orange"}\}$

$f(\text{apple}) = \text{"apple"}$

$f(\text{orange}) = \text{"orange"}$

Regression

Continuous output

e.g.: $x \in \text{Fruits}$, $y \in \mathbb{R}_+ \triangleq \text{t until ripe}$

$f(\text{apple}) = 12 \text{ days}$

Structured Prediction

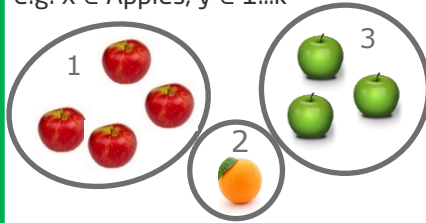
e.g. $x \in \mathbb{R}^{w \times h \times d}$, $y \in \mathbb{R}^{w \times h} \triangleq \text{pixels}$

$f(\text{apple image}) = \text{apple mask}$

Unsupervised Learning (No labels during training)

Clustering

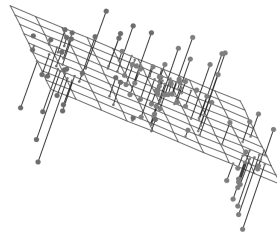
e.g. $x \in \text{Apples}$, $y \in 1 \dots k$



Dimensionality reduction

$x \in \mathbb{R}^d$, $x' \in \mathbb{R}^p$, $p < d$

e.g., projecting all features of a fruit to 2 dimensions for visualization

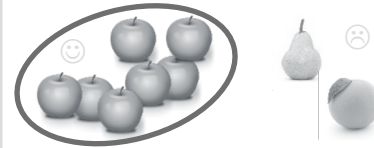


Semi-Supervised Learning (Some labels for training)

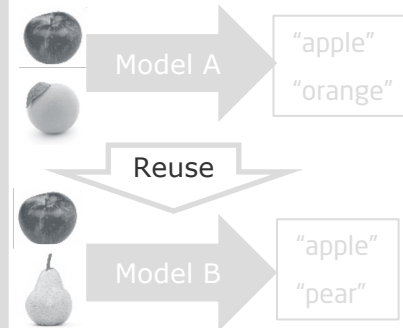
Anomaly / novelty detection

trained only on "normal" samples

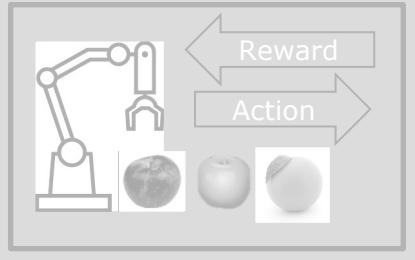
e.g. $x \in \text{Apples}$, $y \in \{\text{😊}, \text{😞}\}$



Transfer Learning



Reinforcement Learning



<https://en.wikipedia.org/wiki/Apple>
<https://cdn4.vectorstock.com/i/1000x1000/16/58/robot-arm-line-icon-sign-on-vector-17841658.jpg>

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
10

- Pros:
 - Do not need number of cluster
 - Could identify noise in data
- Cons
 - Only assign one data to one cluster point. No hierarchy
 - Slow to run on large amounts of data
- Other density-based clustering algorithms
 - HDBscan: Improved version of Dbscan uses **hierarchy** of clusters
 - Optics: Uses **K-nearest neighbour search** to improve algorithm

Q15: What are evaluation metrics for clustering algorithms?
Select all that apply:

- Euclidean distance
- Pearson's correlation coefficient
- Rand index
- Manhattan distance

Frequently missed

Frequent incorrect answer

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
12

- Evaluated based on the data that was clustered itself
- **Ground truth** not needed
- Assigns the best score to clusters with **high similarity** within a cluster and **low similarity** between clusters.
- Silhouette Coefficient
- Dunn Index
- Davies-Bouldin index





- Uses **ground truth** to perform the cluster evaluation
- Evaluates the ability of clustering algorithms ability to separate class compared to ground truth
- **Rand Index**: Measures the similarity of the two assignments.
- Mutual information based: Uses ideas from **Shannon's entropy**
 - Mutual Information Gain
 - Homogeneity, completeness, and V-measure

$$H(U) = - \sum_{i=1}^{|U|} P(i) \log(P(i))$$

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022

Q21: Select all correct answers about Convolutional Neural Networks (CNN)

-  The values of filters are fixed during training of a CNN, as it is not a learnable parameter.
-  Pooling is a type of down sampling method used in CNN.
-  When defining a CNN architecture, one should specify the number of filters along with the filter size and stride of each filter.
-  Padding in a convolution layer averages all the value in a window.

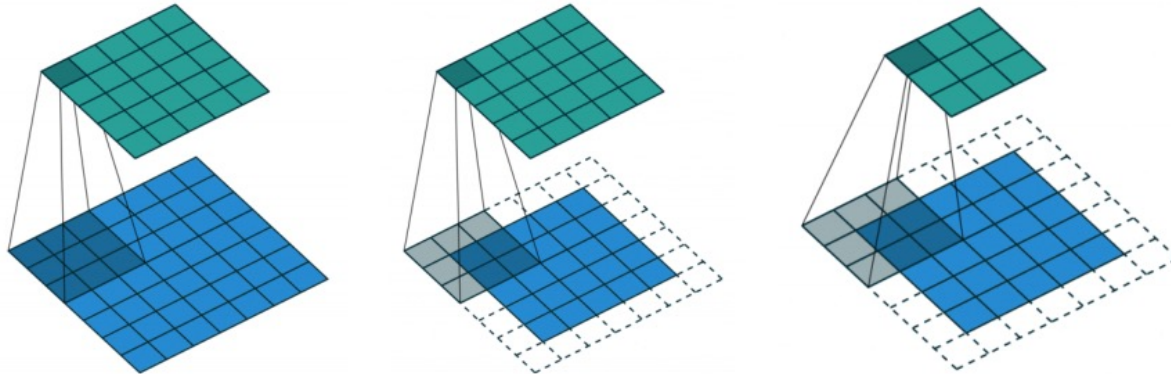
Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
15



CNN

Convolution Layer

- In practice, a CNN learns the values of these **filters** on its own during the **training process**
- Although we still need to specify parameters such as **number of filters**, **filter size**, **padding**, and **stride** before the training process



Convolution Layer Padding

- Adding **extra pixels** of filler around the boundary of input image → Increasing the effective size of the image
- Typically, set values of the extra pixels to zero
- **Valid convolutions**  no padding
- **Same convolution**  pad so that output size is the same as the input size

0	0	0	0	0	0	0	0	0
0	3	3	4	4	7	0	0	0
0	9	7	6	5	8	2	0	0
0	6	5	5	6	9	2	0	0
0	7	1	3	2	7	8	0	0
0	0	3	7	1	8	3	0	0
0	4	0	4	3	2	2	0	0
0	0	0	0	0	0	0	0	0

$6 \times 6 \rightarrow 8 \times 8$

*

1	0	-1
1	0	-1
1	0	-1

3×3

=

-10	-13	1			
-9	3	0			

6×6

<http://datahacker.is/what-is-padding-cnn/>

Evaluation Exercise IV

Data Management for
Digital Health, Winter
2022
17