



# Data Management for Digital Health

## Revision of Exercise III

Borchert, Dr. Schapranow  
Data Management for Digital Health  
Winter 2023

# Exercise III

## Topics

---

- Clinical Interpretation of Molecular Data (Guest Lecture)
- Natural Language Processing
- Neural Networks, Word Embeddings
- Foundation Models
- Medical Use Case Nephrology
- Clinical Prediction Models

### **Evaluation Exercise III**

Data Management for  
Digital Health, Winter  
2023  
2

# Exercise II Key Stats

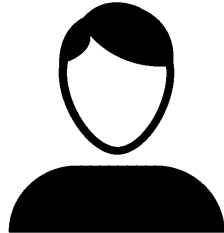
25 Questions  
50 Points

26 Students  
26 Passed

Average score  
44.2 / 88.4%

Average time  
92 min

<< 3h



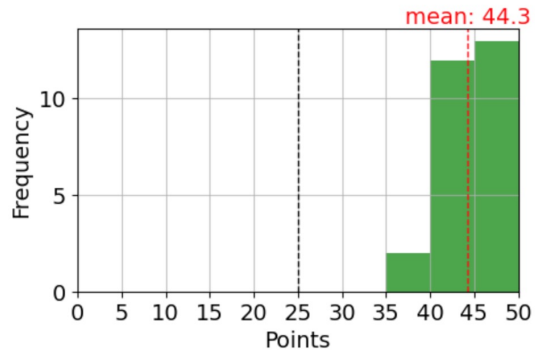
**Evaluation Exercise III**

Data Management for  
Digital Health, Winter  
2023

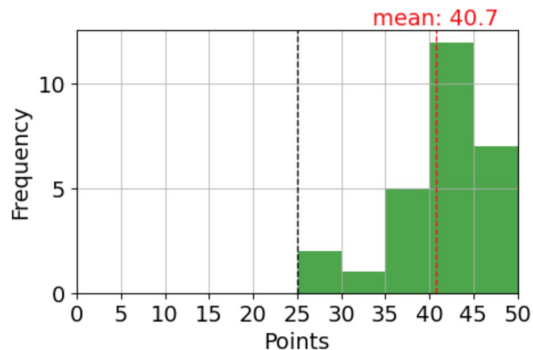
3

# Exercise III

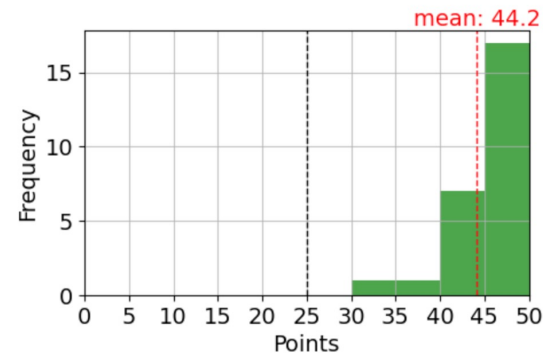
## Key Stats



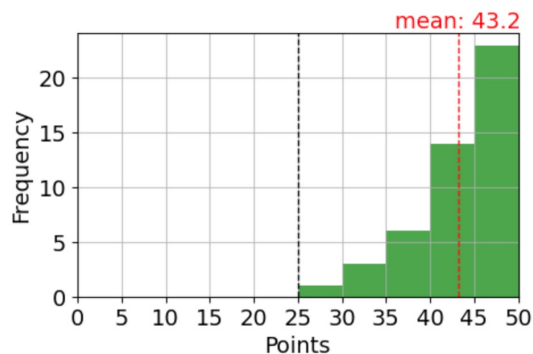
Exercise I (2023)



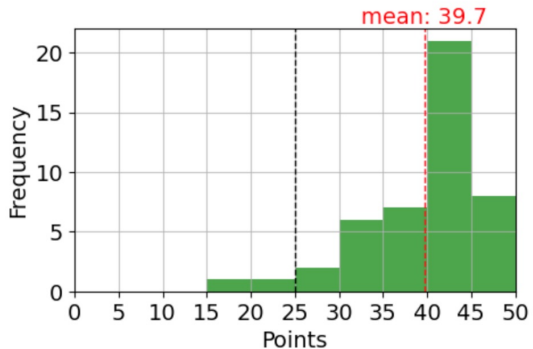
Exercise II (2023)



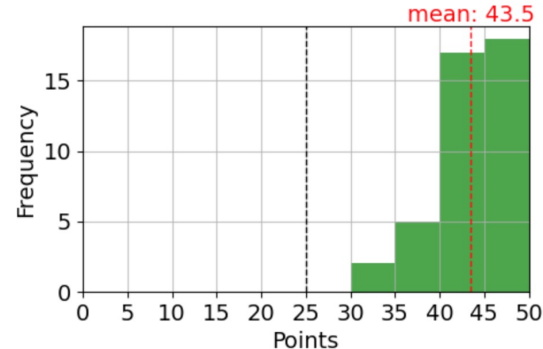
Exercise III (2023)



Exercise I (2022)



Exercise II (2022)



Exercise III (2022)

## Q1 Which of the following statements are true for the use of molecular genetic data in course of precision oncology [...]

- ✓ Clinical relevance of genetic variants is assessed through knowledge bases like CIViC or OnkoKB.
- ✓ Today, whole genome sequencing is applied only for a limited number of severe ill patients, e.g., when enrolled in the DKTK MASTER program.
- ✗ A biological rationale constitutes the highest level of evidence, e.g., using molecular pathways only.
- ✗ Clinical interpretation of molecular data is already fully automated through general-purpose large language models.

### Evaluation Exercise III

Data Management for  
Digital Health, Winter  
2023  
5

# Q1 Which of the following statements are true for the use of molecular genetic data in course of precision oncology [...]?

## Evidence-based Assessment of Clinical Cancer Therapies

- The higher the GoR and LoE, the more evidence about a success of the desired treatment is available

GoR	LoE	Type of Study
A	1a	Systematic review of (homogeneous) randomized controlled trials
A	1b	Individual randomized controlled trials (with narrow confidence intervals)
B	2a	Systematic review of (homogeneous) cohort studies of "exposed" and "unexposed" subjects
B	2b	Individual cohort study / low-quality randomized control studies
B	3a	Systematic review of (homogeneous) case-control studies
B	3b	Individual case-control studies
C	4	Case series, low-quality cohort or case-control studies
D	5	Expert opinions based on non-systematic reviews of results or mechanistic studies

[https://guides.library.stonybrook.edu/evidence-based-medicine/levels\\_of\\_evidence](https://guides.library.stonybrook.edu/evidence-based-medicine/levels_of_evidence)

### Medical Use Case Oncology

Data Management for Digital Health, Winter 2023  
48

### Evaluation Exercise III

Data Management for Digital Health, Winter 2023  
6

#### NGS (accredited/validated)

- Oncomine Focus/Precision DNA Assay
- Oncomine Focus/Precision RNA Assay
- ColonLung Panel V2
- Cancer Hotspot Panel
- Myeloid Panel (Custom)
- (B-cell) Lymphoma Panel
- Oncomine cfDNA (Liquid Biopsies)
- Breast cfDNA Panel (Liquid Biopsies)
- BRCA1/2 Panel
- Tumor Mutational Burden (1.7 Mbases)
- Molecular Health 600+ Panel (3 Mbases); NextSeq
- Oncomine Comprehensive Assay V4 (500+) Panel
- TSO500 (DNA/RNA) Panel
- Ig/TCR Clonality Panel
- Archer RNA Panel

#### IHC/FISH

- nTRK screening
- TMB
- Other Targets (e.g. HER2, Fusion Gene validation)

#### Other

- e.g. EPIC (Methylom)

#### DKTK MASTER

- WES/WGS
- RNASeq

#### ExLiquid

- ctDNA

#### Functional Analyses

#### Single cell analyses

# Q1 Which of the following statements are true for the use of molecular genetic data in course of precision oncology [...]?

Gleiche Tumorentität	m1A	In der <b>gleichen Tumorentität</b> wurde der prädiktive Wert des Biomarkers oder die klinische Wirksamkeit in einer <b>Biomarker-stratifizierten Kohorte</b> einer adäquat gepowerten <b>prospektiven Studie</b> oder <b>Metaanalyse</b> gezeigt.
	m1B	In der <b>gleichen Tumorentität</b> wurde der prädiktive Wert des Biomarkers oder die klinische Wirksamkeit in einer <b>retrospektiven Kohorte</b> oder <b>Fall-Kontroll-Studie</b> gezeigt.
	m1C	Ein oder mehrere <b>Fallberichte</b> in der <b>gleichen Tumorentität</b> .
Andere Tumorentität	m2A	In einer <b>anderen Tumorentität</b> wurde der prädiktive Wert des Biomarkers oder die klinische Wirksamkeit in einer <b>Biomarker-stratifizierten Kohorte</b> einer adäquat gepowerten <b>prospektiven Studie</b> oder <b>Metaanalyse</b> gezeigt.
	m2B	In einer <b>anderen Tumorentität</b> wurde der prädiktive Wert des Biomarkers oder die klinische Wirksamkeit in einer <b>retrospektiven Kohorte</b> oder <b>Fall-Kontroll-Studie</b> gezeigt.
	m2C	Unabhängig von der Tumorentität wurde beim Vorliegen des Biomarkers eine <b>klinische Wirksamkeit</b> in einem oder mehreren <b>Fallberichten</b> gezeigt.
In vitro oder Tiermodell	m3	<b>Präklinische Daten</b> ( <i>in vitro/in vivo</i> -Modelle, funktionelle Untersuchungen) zeigen eine Assoziation des Biomarkers mit der Wirksamkeit der Medikation, welche durch eine wissenschaftliche Rationale gestützt wird.
Biologische Rationale	m4	Eine <b>wissenschaftliche, biologische Rationale</b> legt eine Assoziation des Biomarkers mit der Wirksamkeit der Medikation nahe, welche bisher <b>nicht durch (prä)klinische Daten</b> gestützt wird.

## Evaluation Exercise III

Data Management for  
Digital Health, Winter  
2023

Q4 [...] search for “FOLFIRI antineoplastic chemotherapy regimen”. [...] What can you infer from the SNOMED CT ontology?

- ✓ The concept is connected to its included substances through the “Direct Substance” attribute.
- ✗ The concept only has a single parent, as SNOMED CT is a classification system.
- ✗ The concept is connected to its included substances through the “Method” attribute.
- ✗ “FOLIFIRINOX” is a preferred synonym for the concept.

### Evaluation Exercise III

Data Management for  
Digital Health, Winter  
2023

8



# Q4 [...] search for "FOLFIRI antineoplastic chemotherapy regimen". [...] What can you infer from the SNOMED CT ontology?

Summary Details Diagram Expression Refsets Members History References

Type at least 3 characters ✓ Example: shou fra

FOLFIRI antineoplastic chemotherapy regimen

2 matches found in 0.363 seconds.

FOLFIRI antineoplastic chemotherapy regimen	Folinic acid, fluorouracil and irinotecan antineoplastic chemotherapy regimen (regime/therapy)
FOLFIRINOX antineoplastic chemotherapy regimen	Fluorouracil, leucovorin, irinotecan and oxaliplatin antineoplastic chemotherapy regimen (regime/therapy)

All results are displayed

**Parents**

- Administration of vitamin (procedure)
- Antimetabolite therapy (procedure)
- Antineoplastic chemotherapy regimen (regime/therapy)

**Folinic acid, fluorouracil and irinotecan antineoplastic chemotherapy regimen (regime/therapy)**

SCTID: 870249004

870249004 | Folinic acid, fluorouracil and irinotecan antineoplastic chemotherapy regimen (regime/therapy) |

- en Folinic acid, fluorouracil and irinotecan antineoplastic chemotherapy regimen (regime/therapy)
- en FOLFIRI antineoplastic chemotherapy regimen
- en FOLFIRI (folinic acid, fluorouracil, irinotecan) antineoplastic chemotherapy regimen
- en Folinic acid, fluorouracil and irinotecan antineoplastic chemotherapy regimen

Method → Administration - action  
Direct substance → Irinotecan

Method → Administration - action  
Direct substance → Fluorouracil

Method → Administration - action  
Direct substance → Folinic acid

Has intent → Therapeutic intent

**Children (1)**

- Fluorouracil, leucovorin, irinotecan and oxaliplatin antineoplastic chemotherapy regimen (regime/therapy)

## Evaluation Exercise III

Data Management for  
Digital Health, Winter  
2023  
9

Q14 Which statements are true about few-shot prompting in the context of large language models (LLMs)? Select all that apply:

- ✓  Few-shot prompting is a type of in-context learning, where provided examples are used as context for the LLM predictions.
- ✗  Few-shot prompting requires a fixed set of examples in the prompt to be used for each execution on a new input sample.
- ✗  As opposed to zero-shot prompting, few-shot prompting requires parameter updates to the underlying model.
- ✗  Few-shot prompting prevents LLM hallucinations.

### Evaluation Exercise III

Data Management for  
Digital Health, Winter  
2023  
10

Your task is to determine the smoking status of the person described in a clinical note.

Here are some examples:

Input: "Smoker until 1999"

Output: ex-smoker

Input: "SOCIAL HISTORY: Widowed since 1972, no tobacco, no alcohol, lives alone."

Output: non-smoker

Input: "He is a heavy smoker and drinks 2-3 shots per day at times."

Output: current smoker

Input:

"Social History: No alcohol use and quit tobacco greater than 25 years ago with a 10-pack year smoking history."

Output:

## Prompt Template

## Few-Shot Examples

## Input

- In-context learning (ICL) refers to the model's ability to adapt and respond based on the immediate context provided within a prompt.
- Few-shot prompting is a type of ICL where the model is given a few examples to illustrate the task
- Examples can be fixed or based on the input

## Foundation Models

Data Management for  
Digital Health, Winter  
2023

Q17 Please select all appropriate statements about the Glomerular Filtration Rate (GFR) as discussed in class.

- ✓ GFR can be used as a good indicator for the kidney function and its cleaning quality.
- ✓ Amongst others, personal serum creatine, sex, and age can be used to estimate GFR.
- ✗ Amongst others, renal blood flow (RBF) and renal plasma flow (RPF) can be used to estimate GFR.
- ✗ GFR is defined as the quotient of renal plasma flow (RPF) and urine flow rate (UFR).

Frequently missed

Frequent incorrect answer

**Evaluation Exercise III**

Data Management for  
Digital Health, Winter  
2023  
12

# Kidney Diseases: Glomerular Filtration Rate

- Measuring GFR (mGFR) is complex → Estimated Glomerular Filtration Rate (eGFR) is good approximation for the kidney's function and cleaning quality
- eGFR depends on multiple factors, e.g. age, sex, ethnicity, etc.
- $eGFR = 141 \times \min(\text{Scr}/\kappa, 1)^\alpha \times \max(\text{Scr}/\kappa, 1)^{-1.209} \times 0.993^{\text{Age}} \times 1.018$  [if female]  
\_  $1.159$  [if African American] [1]
  - with Scr = serum creatinine,
  - $\kappa$  is 0.7 for females, 0.9 for males,
  - $\alpha$  is -0.329 for females, -0.411 for males,
  - the minimum of Scr/ $\kappa$  or 1, and
  - the maximum of Scr/ $\kappa$  or 1.

## Medical Use Case Nephrology

Data Management for  
Digital Health, Winter  
2023

13

Q23 During training of your CPM, you figure out that the model performance lacks behind your expectations. Which of the following aspect can be helpful to address this?

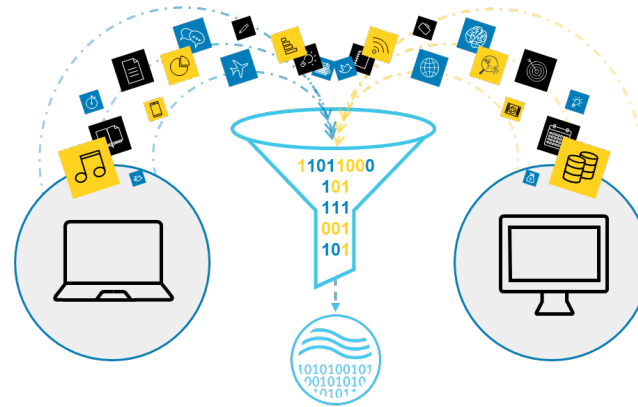
- ✓  Conduct imputation of missing values to improve prediction quality, e.g., use mean values for missing eGFR measurements.
- ✓  Consider feature engineering to improve model accuracy, e.g., select most important features for the clinical endpoint of interest and remove less important features for model training.
- ✗  Ask for additional training data from the same hospitals from 30 years ago.
- ✗  Perform min-max and variance scaling for all features of the complete training dataset.

### Evaluation Exercise III

Data Management for  
Digital Health, Winter  
2023  
14

# 3. Data Preparation: Data Transformation

- **Scaling:** Data may contain attributes with a mixtures of scales, but some ML methods require attributes to have the same scale
- **Decomposition:** Features may represent a complex concept that may be more useful to a ML method when split into its parts, e.g. data, zip code, etc.
- **Aggregation:** Features that might be aggregated into a single feature



<https://blog.dellemc.com/en-us/digital-transformation-just-got-easier-with-analytic-insights/>

## 3 Data Preparation

- Exploration
- Quality assessment
- Cleansing
- Feature engineering
- Labeling

## Clinical Predictive Modeling

Data Management for  
Digital Health, Winter  
2023  
15