

---

# Appendix for Whose Health Matters in Healthcare Models? Unmasking Data Bias for Data-Aware Modeling

---

Anonymous Author(s)

Affiliation

Address

email

## A Datasheet for MIMIC Database

In the advancement of healthcare machine learning (HML) models, the quest for unbiased predictive outcomes is not just technical—it’s fundamentally ethical (Rajkomar et al. (2018)). Despite substantial strides facilitated by HML, the specter of unfairness looms (Buolamwini and Gebru (2018), Corbett-Davies and Goel (2018)), propelled by the dual engines of data and algorithmic biases (Mehrabi et al. (2021)). The Impossibility Theorem in Fairness asserts that models with perfect performance and optimal fairness are simultaneously unattainable except special cases. Nevertheless, we can endeavor to find an optimal balance—a perfect trade-off tailored to specific domains. In the realm of healthcare machine learning (HML), this balance is not just about mitigating algorithmic bias; it’s also about understanding and rectifying data biases. Researchers often face the colossal task of sifting through extensive data documentation to unearth task-specific anomalies. To streamline this process, we present the Datasheet for MIMIC IV v2.0. This resource empowers researchers to discern and address data inconsistencies, guides the selection of sensitive attributes for fairness assessments, and facilitates the creation of robust, just, and data-conscious HML models. More than a mere inventory, the datasheet provides comprehensive insights into the database’s structure, data collection methodologies, management practices, and potential biases, ensuring researchers are well-informed and vigilant.

The MIMIC IV v2.0 datasheet was created based on the template provided by Gebru et al. (2021), with necessary modifications to accommodate the complexities of the database structure in clinical research databases (CRDs). This adaptation was crucial as the original template in Gebru et al. (2021) was designed for datasets rather than databases. To ensure comprehensive coverage of CRD-related information, several changes were made. For example, the question “Can/How the dataset be/are created from the MIMIC database?” was included in the datasheet (indicated by ‘+’), and the question “What is the composition of the dataset?” was replaced with “What is the composition of the database?”. The complete MIMIC IV v2.0 datasheet is provided below.

### A.1 Motivation

*For what purpose was the database created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.*

The creation of the MIMIC-IV CRD aimed to improve patient care through knowledge discovery and algorithm development using a historically collected medical dataset. It was developed with an approach that allows permissive access, enabling extensive utilization of the MIMIC-IV database. Consequently, the database has been widely utilized in various healthcare applications, including assessing treatment effectiveness in specific patient groups and predicting critical outcomes such as mortality, readmission and length of stay (Johnson et al. (2023)).

35 **Who created the database (e.g., which team, research group) and on behalf of which entity (e.g.,**  
36 **company, institution, organization)?**

37 The MIMIC-IV database Johnson et al. (2023), developed by Alistair Johnson, Lucas Bulgarelli,  
38 Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark from the Massachusetts Institute of  
39 Technology at the MIT Laboratory for Computational Physiology, is a collaborative effort involving  
40 various research groups.

41 **Who funded the creation of the database?**

42 The work was supported by grants from the National Institute of Biomedical Imaging and Bioengi-  
43 neering (NIBIB) of the National Institutes of Health (NIH) under award numbers R01-EB001659  
44 (2003-2013) and R01-EB017205 (2014-2018)<sup>1</sup>.

45 **Any other comments?**

46 MIMIC is a large and freely available database that contains deidentified health-related data from  
47 patients admitted to the critical care units of the Beth Israel Deaconess Medical Center. There are  
48 multiple versions of MIMIC that have been released:

- 49 1. **MIMIC-IV** encompasses data collected from 2008 to 2019, obtained from Metavision  
50 bedside monitors Johnson et al. (2023).
- 51 2. **MIMIC-III** comprises data collected from 2001 to 2012, obtained from both Metavision  
52 and CareVue bedside monitors Johnson et al. (2023).
- 53 3. **MIMIC-II** includes data collected from 2001 to 2008, obtained exclusively from CareVue  
54 bedside monitors. While MIMIC-II is no longer publicly available, its data can still be  
55 obtained from MIMIC-III by selectively including the data from the CareVue monitors  
56 Johnson et al. (2023).

57 Throughout the surveyed timeline, MIMIC III and MIMIC IV have been extensively utilized for  
58 healthcare machine learning (HML) prediction models. The datasheet provided is specifically for  
59 MIMIC IV v2.0, the latest available version of the database.

60 **License** - The licensing for the MIMIC files can be found in the *PhysioNet Credentialed Health*  
61 *Data License 1.5.0 (MIT-LCP)*<sup>1</sup>.

## 62 **A.2 Composition**

63 **What is the composition of the database?**

64 MIMIC IV database is grouped into four modules: MIMIC IV (hosp, and icu) , MIMIC IV-ED (ed),  
65 MIMIC IV-Note (note) and MIMIC-CXR (cxr)<sup>1</sup>

- 66 1. **MIMIC IV**Johnson et al. (2023)
  - 67 (a) The **Hosp** module grants access to diverse data extracted from the hospital’s electronic  
68 health record system, including patient and admission details, laboratory measurements,  
69 microbiology information, medication administration records, and billed diagnoses.  
70 These data are organized in the form of tables, including patient and admission-related  
71 tables (patients, admissions, transfers), laboratory measurement tables (labevents,  
72 d\_labitems), microbiology culture table (microbiologyevents), provider order tables  
73 (poe, poe\_detail), medication administration tables (emar, emar\_detail), medication  
74 prescription tables (prescriptions, pharmacy), and hospital billing information tables  
75 (diagnoses\_icd, d\_icd\_diagnoses, procedures\_icd, d\_icd\_procedures, services).
  - 76 (b) The **ICU** module contains information collected from the clinical information system  
77 (BIDMC: MetaVision (iMDSoft)) used within the ICU. Documented data includes intra-  
78 venous administrations, ventilator settings, and other charted items. Data documented  
79 in the icu module includes intravenous and fluid inputs (inputevents), ingredients of the  
80 aforementioned inputs (ingredientevents), patient outputs (outputevents), procedures  
81 (procedureevents), information documented as a date or time (datetimeevents), and  
82 other charted information (chartevents).

---

<sup>1</sup><https://mimic.mit.edu/>

83 2. **MIMIC IV-ED**Johnson et al. (2023) - The *ED* module of MIMIC IV-ED focuses on emer-  
84 gency department patients and encompasses information regarding reasons for admission,  
85 triage assessments, vital signs, and medication reconciliation. The `subject_id` and `hadm_id`  
86 identifiers within MIMIC-IV-ED allow for linkage with other MIMIC-IV modules.

87 3. **MIMIC IV-Note** Johnson et al. (2023) - The *Note* module contains deidentified free-text  
88 clinical notes for hospitalized patients..

89 4. **MIMIC IV-CXR** Johnson et al. (2023) - The *CXR* module of MIMIC IV-CXR provides  
90 lookup tables that establish connections between patient identifiers and MIMIC-CXR  
91 `study_id` and `dicom_id`, facilitating the analysis of patient chest x-rays in conjunction  
92 with clinical data from other MIMIC-IV modules.

93 **+ How is the data arranged within each modules and for what purpose?**

94 The data within each module is structured in the form of tables, as MIMIC is a well-organized  
95 relational database. Each table within a module represents a specific type of data. Within each  
96 table, the data is organized into rows and columns. Each row corresponds to a particular patient  
97 or event, while each column represents a specific variable or attribute corresponding to that row.  
98 This organized structure allows researchers to efficiently extract customized datasets tailored to their  
99 research inquiries and facilitates the construction of machine learning models.

100 **+ Can the modules be linked together to create specific dataset for specific tasks?**

101 **Yes.** The tables within a module can be connected to other tables within the same module or across  
102 different modules using unique identifiers

103 **+ Explain in detail the tables presented in each module?**

104

105 **MIMIC IV**

106

107 **Hosp Module**

Table 1: Description of *Hosp* module Tables with detailed information about the features

Table	Description	Features
omr	The Online Medical Record (OMR) table contains miscellaneous information from the EHR	subject_id, chartdate, seq_num, result_name, result_value
provider	The provider table lists deidentified provider identifiers used in the database	provider_id
admission	Detailed information about hospital stays	subject_id, hadm_id, admittance, dischtime, deathtime, admission_type, admit_provider_id, admission_location, discharge_location, insurance, language, marital_status, race, edregtime, edouttime, hospital_expire_flag
diagnoses_icd	Billed ICD-9/ICD-10 diagnoses for hospitalizations	subject_id, hadm_id, seq_num, icd_code, icd_version
drgcodes	Billed diagnosis-related group (DRG) codes for hospitalizations	subject_id, hadm_id, drg_type, drg_code, description, drg_severity, drg_mortality
emar	The Electronic Medicine Administration Record (eMAR); barcode scanning of medications at the time of administration	subject_id, hadm_id, emar_id, emar_seq, poe_id, pharmacy_id, enter_provider_id, charttime, medication, event_txt, scheduletime, storetime
emar_detail	Supplementary information for electronic administrations recorded in eMAR	subject_id, emar_id, emar_seq, parent_field_ordinal, administration_type, pharmacy_id, barcode_type, reason_for_no_barcode, complete_dose_not_given, dose_due, dose_due_unit, dose_given, dose_given_unit, will_remainder_of_dose_be_given, product_amount_given, product_unit, product_code, product_description, product_description_other, prior_infusi, n_rate, infusion_rate, infusion_rate_adjustment, infusion_rate_adjustment_amount, infusion_rate_unit, route, infusion_complete, completion_interval, new_iv_bag_hung, contin, ed_infusion_in_other_location, restart_interval, side, site, non_formulary_visual_verification
hpcsevents	Billed events occurring during the hospitalization. Includes CPT codes	subject_id, hadm_id, chartdate, hcpcs_cd, seq_num, short_description
labevents	Laboratory measurements sourced from patient-derived specimens	labevent_id, subject_id, hadm_id, specimen_id, itemid, order_provider_id, charttime, storetime, value, valueenum, valueuom, ref_range_lower, ref_range_upper, flag, priority, comments
microbiologyevents	Microbiology cultures	microevent_id, subject_id, hadm_id, micro_specimen_id, order_provider_id, chartdate, charttime, spec_itemid, spec_type_desc, test_seq, storedate, storetime, test_itemid, test_n, me, org_itemid, org_name, isolate_num, quantity, ab_itemid, ab_name, dilution_text, dilution_comparison, dilution_value, interpretation, comments
patients	Patients' gender, age, and date of death if information exists	subject_id, gender, anchor_age, anchor_year, anchor_year_group, dod
pharmacy	Formulary, dosing, and other information for prescribed medications	subject_id, hadm_id, pharmacy_id, poe_id, starttime, stoptime, medication, proc_type, status, entertime, verifiedtime, route, frequency, disp_sched, infusion_type, sliding_scale, lockout_interval, basal_rate, one_hr_max, doses_per_24_hrs, duration, duration_interval, expiration_value, expiration_unit, expirationdate, dispensation, fill_quantity

	poe	Orders made by providers relating to patient care	poe_id, poe_seq, subject_id, hadm_id, ordertime, order_type, order_subtype, transaction_type, discontinued_of_poe_id, discontinued_by_poe_id, order_provider_id, order_status
	poe_detail	Supplementary information for orders made by providers in the hospital	poe_id, poe_seq, subject_id, field_name, field_value
	prescriptions	Prescribed medications	subject_id, hadm_id, pharmacy_id, poe_id, poe_seq, order_provider_id, starttime, stoptime, drug_type, drug, formulary_drug_cd, gsn, ndc, prod_strength, form_rx, dose_val_rx, dose_unit_rx, form_val_disp, form_unit_disp, doses_per_24_hrs, route
	procedures_icd	Billed procedures for patients during their hospital stay	subject_id, hadm_id, seq_num, chartdate, icd_code, icd_version
110	services	The hospital service(s) that cared for the patient during their hospitalization	subject_id, hadm_id, transfertime, prev_service, curr_service
	transfers	Detailed information about patients' unit transfers	subject_id, hadm_id, transfer_id, eventtype, careunit, intime, outtime
	d_hcpcs	Dimension table for hcpcsevents; provides a description of CPT codes	code, category, long_description, short_description
	d_icd_diagnoses	Dimension table for diagnoses_icd; provides a description of ICD-9/ICD-10 billed diagnoses	icd_code, icd_version, long_title
	d_icd_procedures	Dimension table for procedures_icd; provides a description of ICD-9/ICD-10 billed procedures	icd_code, icd_version, long_title
	d_labitems	Dimension table for labevents provides a description of all lab items	itemid, label, fluid, category

## 111 **ICU module**

Table 2: Description of *ICU* module Tables with detailed information about the features

112	Table	Description	Features
	caregiver	The caregiver table lists deidentified provider identifiers used in the ICU module	caregiver_id
	d_items	Dimension table describing itemid. Defines concepts recorded in the events table in the ICU module	itemid, label, abbreviation, linksto, category, unitname, param_type, lownormalvalue, highnormalvalue
	chartevents	Charted items occurring during the ICU stay. Contains the majority of information documented in the ICU	subject_id, hadm_id, stay_id, caregiver_id, charttime, storetime, itemid, value, valuenum, valueuom, warning
	datetimeevents	Documented information which is in a date format (e.g., date of last dialysis)	subject_id, hadm_id, stay_id, caregiver_id, charttime, storetime, itemid, value, valueuom, warning
113	icustays	Tracking information for ICU stays including admission and discharge times	subject_id, hadm_id, stay_id, first_careunit, last_careunit, intime, outtime, los
	Ingredientevents	Ingredients of continuous or intermittent administrations including nutritional and water content	subject_id, hadm_id, stay_id, caregiver_id, starttime, endtime, storetime, itemid, amount, amountuom, rate, rateuom, orderid, linkorderid, statusdescription, originalamount, originalrate
	inpuvents	Information documented regarding continuous infusions or intermittent administrations	subject_id, hadm_id, stay_id, caregiver_id, starttime, endtime, storetime, itemid, amount, amountuom, rate, rateuom, orderid, linkorderid, ordercategoryname, secondaryordercategoryname, ordercomponenttypedescription, ordercategorydescription, patientweight, totalamount, totalamountuom, isopenbag, statusdescription, originalamount, originalrate

114	outputevents	Information regarding patient outputs including urine, drainage, and so on	subject_id, hadm_id, stay_id, caregiver_id, charttime, storetime, itemid, value, valueuom
	procedureevent	Procedures documented during the ICU stay (e.g., ventilation), though not necessarily conducted within the ICU (e.g., x-ray imaging)	subject_id, hadm_id, stay_id, caregiver_id, starttime, endtime, storetime, itemid, value, valueuom, location, locationcategory, orderid, linkorderid, ordercategoryname, ordercategorydescription, patientweight, isopenbag, continueinnextdept, statusdescription, originalamount, originalrate

## 115 MIMIC IV-ED

Table 3: Description of *ED* module Tables with detailed information about the features

116	Table	Description	Features
	diagnosis	The diagnosis table provides billed diagnoses for patients. Diagnoses are determined after discharge from the emergency department	subject_id, stay_id, seq_num, icd_code, icd_version, icd_title
	edstays	The edstays table is the primary tracking table for emergency department visits. It provides the time the patient entered the emergency department and the time they left the emergency department	subject_id, hadm_id, stay_id, intime, outtime, gender, race, arrival_transport, disposition
117	medrecon	On admission to the emergency departments, staff will ask the patient what current medications they are taking. This process is called medicine reconciliation, and the medrecon table stores the findings of the care providers	subject_id, stay_id, charttime, name, gsn, ndc, etc_rn, etccode, etcdescription
	pyxis	The pyxis table provides information for medicine dispensations made via the Pyxis system	subject_id, stay_id, charttime, med_rn, name, gsn_rn, gsn
	triage	The triage table contains information about the patient when they were first triaged in the emergency department	subject_id, stay_id, temperature, heartrate, resprate, o2sat, sbp, dbp, pain, acuity, chiefcomplaint
	vitalsign	Patients admitted to the emergency department have routine vital signs taken every 1-4 hours. These vital signs are stored in the vitalsign table	subject_id, stay_id, charttime, temperature, heartrate, resprate, o2sat, sbp, dbp, rhythm, pain

## MIMIC IV-CXR

Table 4: Description of CXR module Tables with detailed information about the features

Table	Description	Features
cxr_record_list	Lists all records in the MIMIC-CXR database	subject_id, study_id, dicom_id

## MIMIC IV-Note

Table 5: Description of Note module Tables with detailed information about the features

Table	Description	Features
discharge	Discharge summaries for hospitalizations	note_id, subject_id, hadm_id, note_type, note_seq, charttime, storetime, text
discharge_detail	Auxiliary information for discharge summaries	note_id, subject_id, field_name, field_value, field_ordinal
radiology	Radiology report	note_id, subject_id, hadm_id, note_type, note_seq, charttime, storetime, text
radiology_detail	Auxiliary information for radiology notes	note_id, subject_id, field_name, field_value, field_ordinal
cxr_record_list	Lists all records in the MIMIC-CXR database	subject_id, study_id, dicom_id

### + Can/How the dataset be/are created from the MIMIC database?

The MIMIC database is a comprehensive clinical research database that encompasses various types of data, such as patient admissions, ICU records, triage information, bedside health records, X-rays, and clinician medical notes. It offers researchers the flexibility to create custom datasets tailored to their specific research tasks.

For example, if the objective is to predict **heart failure**, relevant cohorts related to heart failure can be extracted from tables like admission, patient, diagnoses\_icd, and d\_icd\_diagnoses in the hosp module. Additional features associated with heart failure can be obtained by linking tables from the ICU module and ED module. Once the cohort and their corresponding heart-related features are extracted, they undergo pre-processing and cleaning before being represented in either a time series or non-time series format, depending on the prediction task. This allows for the creation of suitable datasets for predictive modeling. Similarly, researchers can curate a wide range of task-specific datasets based on their specific needs.

### What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?

A dataset derived from the MIMIC contains patient health data. The data can be patient demography (Age, gender, ethnicity, language etc.), ICU details, X-ray images or even Clinician notes. It differs depending on the the intended prediction task.

### How many instances are there in total (of each type, if appropriate)?

Dataset is extracted from the MIMIC database based on the intended task and count of instances depend on the dataset extracted.

For instance, If we intend to create a *MIMIC IV ED dataset* by linking ED, hosp and ICU modules then the dataset will have 425087 instances. Similarly several complex datasets can be created and the instance of the dataset vary depending of the prediction task/requirements.

**Does the dataset/database contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).**

154 The MIMIC-IV database is a subset of deidentified electronic health records (EHRs) obtained from  
155 patients admitted to BIDMC between 2008 and 2019. It is a curated collection that has undergone  
156 validation and quality assurance by a team of interdisciplinary experts. The database includes a  
157 diverse range of patients and diagnoses, making it suitable for various research purposes. However, it  
158 is important to acknowledge that the dataset is not comprehensive, as it represents a subset of the  
159 overall patient population. Researchers should be mindful of potential biases inherent in the dataset  
160 and employ appropriate methods to address them when conducting analyses or studies<sup>1</sup>.

161 ***What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or features?***  
162 ***In either case, please provide a description.***

163 MIMIC IV (hosp and ICU module) and MIMIC-IV-ED (ED module) consists of raw unprocessed text,  
164 Date time and number data in comma separated format of the patients admitted to the hospital, ICU  
165 and ED. Whereas, MIMIC-CXR and MIMIC-Note contains images of the Chest X-Rays and free-text  
166 clinical notes for hospitalized patients respectively. Table ?? provides detailed feature information of  
167 the data.

168 ***Is there a label or target associated with each instance? If so, please provide a description.***

169 The choice of target variable in the MIMIC dataset depends on the specific prediction task at hand.  
170 For example, if the goal is to predict the length of stay in the ICU, the *los* attribute in the *icustay* table  
171 can serve as the target variable. On the other hand, if the objective is to predict *in-hospital mortality*,  
172 the *hospital\_expire\_flag* in the *admissions* table can be used as the target variable. **The selection of**  
173 **the target variable is contingent upon the specific prediction task being undertaken.**

174 ***Are there recommended data splits (e.g., training, development/validation, testing)?***

175 No

176 ***Are there any errors, sources of noise, or redundancies in the database? If so, please provide a***  
177 ***description.***

178 Our analysis of the MIMIC IV dataset has revealed several biases and inconsistencies that researchers  
179 should be aware of,

- 180 1. **Inconsistencies in patient details:** Patient language is inconsistently recorded, with only  
181 English being specified while other languages are marked as '?' or unknown.
- 182 2. **Inconsistencies in in-hospital expiry information:** The admission table contains multiple  
183 reports of the same patient's death, leading to inconsistencies.
- 184 3. **Vagueness in insurance coverage information:** The dataset lacks definitive information  
185 about insurance coverage, limiting researchers' ability to draw conclusions on insurance  
186 choices.
- 187 4. **Inconsistencies in hospital admit and discharge timestamps:** The admission table exhibits  
188 inconsistencies in the recorded timestamps, and there are also *missing values for death time*.
- 189 5. **Potential representation bias in the dataset:** The database owners acknowledge the  
190 potential for bias, particularly since the data is derived from a single hospital system and  
191 may not be **representative of the entire population**.

192 The data Johnson et al. (2023) in the database is collected during routine clinical practice, reflecting  
193 the specific practices of the hospital. It is important to note that there may be implausible values  
194 present in the database due to the archival process<sup>1</sup>. Therefore, caution should be exercised when  
195 using the data, and researchers should be mindful of the dataset limitations and potential biases.

196 We strongly recommend that researchers adhere to best practice guidelines Goldberger et al. (2000)  
197 when analyzing the data

198 ***Does the database contain data that might be considered confidential (e.g., data that is protected by***  
199 ***legal privilege or by doctor– patient confidentiality, data that includes the content of individuals’***  
200 ***non-public communications)?***

201 Yes, the MIMIC IV dataset includes medical records of patients, encompassing confidential personal  
202 and health-related information. However, the dataset is constructed with patient privacy as a priority,  
203 and all data within the database undergoes de-identification processes to comply with Health Insurance  
204 Portability and Accountability Act (HIPAA) regulations.



205 ***Does the database identify any subpopulations (e.g., by age, gender)?***

206 Yes. Database (specifically admission and patient tables) has patient demographic data such as age,  
207 gender, ethnicity, language, insurance and marital status.

208 **MIMIC IV Distribution statistics**

Table 6: Admission distribution statistics

Description	Value
Total records	180,733
Male	47%
Female	53%
Min Age	18
Max Age	91
Predominant Ethnicity	White (67.2%)

Table 7: Patient distribution statistics

Description	Value
Total records	299,712
Male	47%
Female	53%
Min Age	18
Max Age	91

Table 8: ED table distribution statistics

Description	Value
Total records	299712
Male	46%
Female	54%
Predominant Ethnicity	White (58%)
Predominant Disposition	Home

209 ***Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly***  
210 ***(i.e., in combination with other data) from the database? If so, please describe how.***

211 No, all data in the database is de-identified in accordance with HIPAA regulations.

212 ***Does the dataset contain data that might be considered sensitive in any way (e.g., data that***  
213 ***reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union***  
214 ***memberships, or locations; financial or health data; biometric or genetic data; forms of government***  
215 ***identification, such as social security numbers; criminal history)?***

216 Yes, database recorded demographic information like ethnicity, gender, age, marital status, language  
217 and insurance

218 **+ Does researchers have to take any important measures to handle the data with care?**

219 To ensure patient privacy, researchers are required to comply data usage agreements mandated in  
220 Goldberger et al. (2000), Johnson et al. (2023) and obtain the necessary approvals and certifications  
221 before accessing the dataset. Researchers working with healthcare-related data have a responsibility to  
222 handle the data carefully and ethically, taking measures to prevent any potential harm or dissatisfaction.  
223 While the data is de-identified in accordance with HIPAA regulations, it is crucial to treat the data  
224 with respect and caution, following best practices. Additionally, the collection of patient information  
225 and the creation of the research resource have been approved by the Institutional Review Board of  
226 the Beth Israel Deaconess Medical Center.

227 **A.3 Collection Process**

228 ***How was the data associated with each instance acquired? Was the data directly observable (e.g.,***  
229 ***raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived***  
230 ***from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If the data***  
231 ***was reported by subjects or indirectly inferred/derived from other data, was the data validated/veri-***  
232 ***fied? If so, please describe how.***

233 The data was extracted from the hospital databases of the Beth Israel Deaconess Medical Center  
234 (BIDMC) specifically for patients admitted to the intensive care units. A comprehensive patient list

was compiled, including all medical record numbers associated with ICU or emergency department admissions from 2008 to 2019. To ensure the reliability of the database, a multidisciplinary team of scientists and clinicians thoroughly evaluated MIMIC-IV during its development, conducting code reviews and documenting any identified issues using a ticket system Johnson et al. (2023).

***What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses or sensors, manual human curation, software programs, software APIs)? How were these mechanisms or procedures validated?***

MIMIC-IV is derived from two distinct database systems within the hospital setting: a customized electronic health record (EHR) used across the entire hospital and a specialized clinical information system called *MetaVision (iMDSoft)* specifically designed for the intensive care units at the Beth Israel Deaconess Medical Center (BIDMC).

To ensure the accuracy and reliability of the MIMIC-IV dataset, a diverse team of scientists and clinicians conducted a comprehensive evaluation during its development, which included code reviews and the systematic documentation of identified issues using a ticket system Johnson et al. (2023).

***Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.***

Over a period of 11 years, from 2008 - 2019

***Were any ethical review processes conducted (e.g., by an institutional review board)?***

Yes, the collection of patient information and creation of the research resource was reviewed by the Institutional Review Board at the Beth Israel Deaconess Medical Center.

***Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?***

Data is collected from hospital EHR and ICU specific clinical information system at the BIDMC called *CareVue and MetaVision (iMDSoft)*.

***Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted? If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation.***

Unknown, however the MIMIC data is deidentified<sup>1</sup>, and patient identifiers were removed according to the Health Insurance Portability and Accountability Act (HIPAA) Safe Harbor provision Johnson et al. (2023).

#### **A.4 Preprocessing/cleaning/labeling**

***Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remaining questions in this section.***

The data within the MIMIC-IV database underwent reorganization to enhance its suitability for retrospective data analysis Johnson et al. (2023). This involved denormalizing tables, eliminating audit trails, and consolidating the data into a smaller number of tables. The primary objective of this process was to simplify the retrospective analysis of the database. Notably, no data cleaning procedures were applied to ensure that the dataset accurately represents real-world clinical data<sup>1</sup>.

To protect patient privacy, patient identifiers were removed in compliance with HIPAA regulations. Random ciphers were used to replace patient identifiers, resulting in deidentified integer values for patients, hospitalizations, and ICU stays. Structured data underwent filtering using look-up tables and allow lists. Additionally, dates and times were randomly shifted into the future by a specific number of days. Consequently, the data for each individual patient remains internally consistent Johnson et al. (2023).

***Is the software that was used to preprocess/clean/label the data available? If so, please provide a link or other access point.***

Unknown. However, authors have stated that the application of a free-text deidentification algorithm was used as a measure to remove personally identifiable information (PHI) from the free-text data, if needed.

## A.5 Uses

*Has the database/dataset been used for any tasks already? If so, please provide a description.*

Yes, MIMIC database is one of the most widely used CRD. It has been widely used for below types of works,

1. Prediction tasks like,

- (a) **Readmission** Assaf and Jayousi (2020), Chen et al. (2022), Partovi et al. (2022), Rojas et al. (2018), Thacker (2023) (30, 60, 90, 120 and custom days) - Predict patients at risk of readmission early in the health care process(helps to prioritize care towards such patients preventing mortality and readmission).
- (b) **Mortality** Ahmad et al. (2021), Caicedo-Torres and Gutierrez (2019), Feng et al. (2018), Kong et al. (2020), Lin et al. (2019) - Predict the likelihood of patients dying.
  - i. **In-hospital** Chen et al. (2021), Ding et al. (2021), Theis et al. (2021), Yang et al. (2023b) - Predict the likelihood of patient dying in hospital while they are admitted (helpful to identify high risk patients early on to provide medical interventions).
  - ii. **Short term** Gao et al. (2021), Hou et al. (2020), Lu et al. (2021), Luo et al. (2022), Zhang et al. (2022a) - Predict short-term mortality (typically within 2-3 days) after ICU admission
  - iii. **Long term** Ahmad et al. (2021), Caicedo-Torres and Gutierrez (2019), Kong et al. (2020), Liu et al. (2021a) - Predict long-term mortality (typically within 30 days to 1 year) after hospital discharge.
- (c) **Length of stay (LOS)** Geethamani and Rangaraj, Geethamani and Rangaraj, Liu et al. (2023), Shu et al. (2023), Wang et al. (2022) - Predict the length of stay of each admission( typically predicting > 3 and 7 days stay. Custom days is also being predicted.)
- (d) **Phenotype label and ICD-9/10 code grouping** - Helpful in tasks like disease prediction, outcome analysis, treatment recommendation and customized treatments.
  - i. **Phenotype labelling** Dong et al. (2022), Singh et al. (2020), Yang et al. (2023a), Zhang et al. (2022b)classify patients into specific groups based on their diagnoses, procedures, medications, and other clinical variables.
  - ii. **Grouping ICD 9/10 codes** Huang et al. (2019), Li and Yu (2020), Li et al. (2018)into different categories based on patient diagnosis to classify the disease.

2. Prediction for specific **health ailments** like,

- (a) Heart failure Ali et al. (2022), Li et al. (2021)
- (b) Chronic Kidney Disease (CKD) Sun et al. (2022), Yue et al. (2022)
- (c) Chronic obstructive pulmonary disease (COPD) Liu et al. (2021b), Rojas et al. (2018)
- (d) Coronary artery disease (CAD) Yang et al. (2023b), Ye et al. (2023)
- (e) Sepsis Böck et al. (2022), Yue et al. (2022)
- (f) Cancer Kurniati et al. (2018), Magna et al. (2020)
- (g) Ventilation failure Geri et al. (2021), Sayed et al. (2021)

*Is there a repository that links to any or all papers or systems that use the database/dataset? If so, please provide a link or other access point.*

No, however the owners Johnson et al. (2023) have provided the repository <https://github.com/MIT-LCP/mimic-code> where the code and other discussions related to the database are hosted.

*Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks,*

334 *financial harms)? If so, please provide a description. Is there anything a dataset consumer could*  
335 *do to mitigate these risks or harms?*

336 The data available in the database reflects the idiosyncrasies of routine clinical practice, as stated by  
337 the owners. The archival process may have introduced implausible values and potential bias into the  
338 data. Therefore, it is important for researchers to follow best practice guidelines when using the data  
339 for analysis or other purposes Johnson et al. (2023).

340 *Are there tasks for which the dataset should not be used? If so, please provide a description.*

341 Unknown, the owners of the database did not provide clear information in the documentation.

## 342 **A.6 Distribution**

343 *+ Is the data publicly available? How and where can it be accessed (e.g., website, GitHub)?*

344 Yes. The MIMIC-IV data is accessible to the public through the PhysioNet<sup>2</sup>. To gain access,  
345 individuals need to become a PhysioNet credentialed user and agree to the data use agreement. Once  
346 granted access, users can download the complete set of file or select specific subsets that align with  
347 their requirements.

348 *Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or*  
349 *under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a*  
350 *link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well*  
351 *as any fees associated with these restrictions.*

352 Access to the MIMIC-IV data is granted through a license agreement called the *Data Use Agreement*  
353 (*DUA*), which outlines the terms and conditions for data usage. To obtain access, users are required to  
354 complete an online course on the ethical use of human subjects research data and obtain a certificate  
355 of completion. With the certificate, users can then apply for dataset access through the PhysioNet<sup>2</sup>.  
356 The application process involves agreeing to the DUA terms and providing details about the intended  
357 use of the data.

## 358 **A.7 Maintenance**

359 *Is the database maintained? Who will be supporting/hosting/maintaining the database?*

360 Yes, MIMIC-IV is maintained by the Laboratory for Computational Physiology at the Massachusetts  
361 Institute of Technology (MIT) and BIDMC. They provide ongoing support and maintenance for the  
362 database

363 *How can the owner/curator/manager of the database be contacted (e.g., email address)?*

364 For private issue, they can be contacted at `mimic-support@physionet.org` and for issues related  
365 to patient health information (PHI) `phi-report@physionet.org` is being used<sup>1</sup>.

366 *Will the database be updated (e.g., to correct labeling errors, add new instances, delete instances)?*  
367 *If so, please describe how often, by whom, and how updates will be communicated to dataset*  
368 *consumers (e.g., mailing list, GitHub)?*

369 Yes, the MIMIC-IV database is regularly updated by the MIT Laboratory for Computational Physi-  
370 ology team. The latest version, v2.2, has been released, which includes updates from the previous  
371 version, v1.0. The frequency of future updates is unknown, but any information regarding up-  
372 dates can be found on the official website <https://physionet.org/content/mimiciv/2.2/>  
373 and <https://github.com/MIT-LCP/mimic-code>.

374 *Will older versions of the database continue to be supported/hosted/maintained? If so, please*  
375 *describe how. If not, please describe how its obsolescence will be communicated to dataset*  
376 *consumers.*

377 Previous versions of the database will continue to be supported and maintained, however it is not  
378 explicitly stated whether they might have any further updates by the owners.

---

<sup>2</sup><https://physionet.org/content/mimiciv/2.2/>

379 *If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them*  
 380 *to do so? If so, please provide a description.*

381 Content for the MIMIC website and documentation is hosted publicly on GitHub: <https://github.com/MIT-LCP/mimic-website>. To raise a problem or to suggest an improvement, new issue can be  
 382 created at: <https://github.com/MIT-LCP/mimic-website/issues>. To take part in discussion  
 383 <https://github.com/MIT-LCP/mimic-code/discussions> can be used.

384 + *What is the Data life cycle of MIMIC database?*

- 386 1. Data Acquisition
- 387 2. Data Archive
- 388 3. Data Preparation
- 389 4. Data loading

390 constitutes the life cycle if MIMIC CRD.

### 391 **Data Acquisition**

392 Data is collected from the source which may be internal, external or both.

Table 9: Data Acquisition

Internal (In-hospital) data	External data source
ICU/MICU/SICU/CCU/CVICU/NICU data (vitals, trends, anomalies) Chart details (Fluids, medications etc.) Demographics (age, gender, ethnicity, language, marital status, religion, insurance etc.) Lab reports Billing details Physician notes Provider order entries etc.	Social security death Index etc

393 **Data Archive** Data collected from the source is Archived before proceeding with data preparation for  
 394 later use.

### 395 **Data Preparation**

396 To ensure compliance with HIPAA regulations, measures such as deidentification, date shifts, and  
 397 format conversions are applied to the archival data. The data is then reorganized into a more suitable  
 398 format for retrospective analysis, which involves consolidating tables, denormalizing data, and  
 399 removing audit trails. It's important to note that no data cleaning procedures were performed to  
 400 maintain the authenticity of the real-world clinical dataset. Feedback from users will be considered  
 401 for further iterations, and the final version of the data will be loaded into the database.

### 402 **Data loading to database**

403 Final version of the data is then loaded to the database which is built on a PostgreSQL relational  
 404 database management system and is hosted on a secure server infrastructure. The data can either be  
 405 downloaded locally or accessed on the cloud via BigQuery, AWS or GCS

## 406 **B Risk Prediction Task Analysis**

407 We analyzed prevalent HML prediction models using the MIMIC dataset across all demograph-  
 408 ics. Access to the MIMIC IV v2.0 data <https://physionet.org/content/mimiciv/2.2/> is  
 409 restricted to PhysioNet credentialed users only.

## B.1 Credential Access process

Researchers can gain physionet Credential access by following the below,

1. Submit their personal details for physionet’s review <sup>3</sup>
2. Complete CITI Data or Specimens Only Research <sup>4</sup> training and/or sign a Data Use Agreement<sup>5</sup>

Once PhysioNet has reviewed the personal details, CITI training completion, and DUA agreement, they will grant access to the MIMIC IV database. You can retrieve all the necessary files from the files section of the MIMIC IV v2.0 <sup>6</sup>.

## B.2 Prediction task-specific dataset retrieval

For our analysis, we retrieved datasets for various prediction tasks, including in-hospital mortality, heart failure, chronic kidney disease (CKD), 30-day readmission, and length of stay (LOS) for heart failure. We followed an established pipeline Gupta et al. (2022) for most tasks, except for sepsis mortality, where we directly retrieved the dataset without a specific pipeline to evaluate and authenticate the results obtained from both the strategies.

## B.3 Additional Analysis Information

Python 3 software was used for the analysis, leveraging its analytical and statistical libraries. We examined the data distribution, including 33 different ethnic records, and categorized patients into Asian, Black, Hispanic/Latino, Other, and White ethnic subgroups following the previous fairness works Meng et al. (2022), Rösli et al. (2021). Our analysis uncovered disparities in patient treatment based on ethnicity and insurance status, impacting patient outcomes. Statistical tests, such as chi-square and ANOVA, confirmed a significant association between ethnicity and prediction outcome, *highlighting the importance of considering model performance across all demographics for fair, generalizable, and data-aware HML models.*

## C Example Usage & Reproducibility

The complete code for analyzing the length of ICU stay for heart failure patients over 7 days is available at <https://github.com/Trustworthy-HML-Models/Unmasking-Data-Bias-for-Data-Aware-Modeling>, and other prediction tasks can be reproduced by substituting the respective datasets.

```
df_new = pd.read_csv('Task-specific-dataset.csv')
```

Distribution of the dataset is analysed to understand the dataset better and the ethnic grouping is performed

```
racial_groupings = {
    'White': ['white', 'white - brazilian', 'white - eastern european',
              'white - other european', 'white - russian', 'portuguese'],
    'Black': ['black/african', 'black/african american', 'black/cape
              verdean', 'black/caribbean island', 'south american'],
    'Hispanic/Latino': ['hispanic or latino', 'hispanic/latino -
                        central american', 'hispanic/latino - columbian', 'hispanic/
                        latino - cuban', 'hispanic/latino - dominican', 'hispanic/
                        latino - guatemalan', 'hispanic/latino - honduran', 'hispanic/
                        latino - mexican', 'hispanic/latino - puerto rican', 'hispanic/
                        latino - salvadoran'],
```

<sup>3</sup>Personal detail submission - Available at <https://physionet.org/settings/credentialing/>

<sup>4</sup>CITI Data or Specimens Only Research - <https://physionet.org/content/mimiciv/view-required-training/2.2/>

<sup>5</sup>DUA - Available at <https://physionet.org/content/mimiciv/view-dua/2.2/>

<sup>6</sup>MIMIC IV v2.0 - Available at <https://physionet.org/content/mimiciv/2.2/>

```

453     'Asian': ['asian', 'asian - asian indian', 'asian - chinese', '
454             asian - korean', 'asian - south east asian'],
455     'Other': ['native hawaiian or other pacific islander', 'other', '
456             patient declined to answer', 'unable to obtain', 'unknown', '
457             american indian/alaska native', 'multiple race/ethnicity']
458 }
459
460 # Replace the original race values with the new groupings
461
462 df_new['ethnicity'] = df_new['ethnicity'].str.strip().str.lower()
463
464
465 def get_race_sub_group(race):
466     for group, races in racial_groupings.items():
467         if race in races:
468             return group
469     return 'Unknown'
470
471 df_new['race_sub_group'] = df_new['ethnicity'].apply(
472     get_race_sub_group)

```

473 Further detailed analysis of patient treatment based on ethnic and Insurance demographics is done  
474 and they are visualized for better understanding the insights of the results. This analysis uncovered  
475 disparities in patient treatment based on ethnicity and insurance status, impacting patient outcomes.  
476 Statistical tests, such as chi-square and ANOVA is performed to identify the association between  
477 ethnicity and prediction outcome

```

478
479 # Create a contingency table
480 contingency_table = pd.pivot_table(df, values='count', index=['
481     race_sub_group', 'insurance'], columns='label', fill_value=0)
482
483 # Perform the chi-square test
484 chi2, p, _, _ = chi2_contingency(contingency_table)
485
486 # Print the contingency table
487 print("Contingency Table:")
488 print(contingency_table)
489
490 # Print the test statistic and p-value
491 print("Chi-square test statistic:", chi2)
492 print("p-value:", p)
493
494 # Perform ANOVA
495 from scipy import stats
496 result = stats.f_oneway(
497     df[df['race_sub_group'] == 'Asian']['count'],
498     df[df['race_sub_group'] == 'Black']['count'],
499     df[df['race_sub_group'] == 'Hispanic/Latino']['count'],
500     df[df['race_sub_group'] == 'Other']['count'],
501     df[df['race_sub_group'] == 'White']['count']
502 )
503
504 # Print the ANOVA test result
505 print("ANOVA test result:")
506 print("F-value:", result.statistic)
507 print("p-value:", result.pvalue)

```

508 Test results confirmed a significant association between ethnicity and prediction outcome, *highlighting*  
509 *the importance of considering model performance across all demographics for fair, generalizable,*  
510 *and data-aware HML models*

## References

- Fahad Shabbir Ahmad, Liaqat Ali, Hasan Ali Khattak, Tahir Hameed, Iram Wajahat, Seifedine Kadry, and Syed Ahmad Chan Bukhari. A hybrid machine learning framework to predict mortality in paralytic ileus patients using electronic health records (ehrs). *Journal of Ambient Intelligence and Humanized Computing*, 12:3283–3293, 2021.
- Aliya Ali, Saleha Yurf Asghar, Ali Danish Khan Yousafzai, Ali Haider Bangash, Rabia Mohsin, Arshiya Fatima, Saiqa Zehra, Ayesha Khalid Khan, Ali Haider Shah, Syed Mohammad Mehmood Abbas, et al. Prediction of in-hospital mortality among heart failure patients: An automated machine learning analysis of mimic-iii database. *American Heart Journal*, 254:261, 2022.
- Rasha Assaf and Rashid Jayousi. 30-day hospital readmission prediction using mimic data. In *2020 IEEE 14th International Conference on Application of Information and Communication Technologies (AICT)*, pages 1–6. IEEE, 2020.
- Markus Böck, Julien Malle, Daniel Pasterk, Hrvoje Kukina, Ramin Hasani, and Clemens Heitzinger. Superhuman performance on sepsis mimic-iii data by distributional reinforcement learning. *PLoS One*, 17(11):e0275358, 2022.
- Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.
- William Caicedo-Torres and Jairo Gutierrez. Iseu: Visually interpretable deep learning for mortality prediction inside the icu. *Journal of biomedical informatics*, 98:103269, 2019.
- Han Chen, Shu-Rong Gong, and Rong-Guo Yu. Association between normalized lactate load and mortality in patients with septic shock: an analysis of the mimic-iii database. *BMC anesthesiology*, 21(1):1–8, 2021.
- Qifan Chen, Yang Lu, Charmaine Tam, and Simon Poon. Outcome-oriented predictive process monitoring to predict unplanned icu readmission in mimic-iv database. 2022.
- Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*, 2018.
- Ning Ding, Cuirong Guo, Changluo Li, Yang Zhou, and Xiangping Chai. An artificial neural networks model for early predicting in-hospital mortality in acute pancreatitis in mimic-iii. *BioMed research international*, 2021, 2021.
- Hang Dong, Víctor Suárez-Paniagua, Huayu Zhang, Minhong Wang, Arlene Casey, Emma Davidson, Jiaoyan Chen, Beatrice Alex, William Whiteley, and Honghan Wu. Ontology-based and weakly supervised rare disease phenotyping from clinical notes. *arXiv preprint arXiv:2205.05656*, 2022.
- Mengling Feng, Jakob I McSparron, Dang Trung Kien, David J Stone, David H Roberts, Richard M Schwartzstein, Antoine Vieillard-Baron, and Leo Anthony Celi. Transthoracic echocardiography and mortality in sepsis: analysis of the mimic-iii database. *Intensive care medicine*, 44:884–892, 2018.
- Qiaoyan Gao, Dandan Wang, Pingping Sun, Xiaorong Luan, Wenfeng Wang, et al. Sentiment analysis based on the nursing notes on in-hospital 28-day mortality of sepsis patients utilizing the mimic-iii database. *Computational and Mathematical Methods in Medicine*, 2021, 2021.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92, 2021.
- D Geethamani and R Rangaraj. Heterogeneous multi-model ensemble based length of stay prediction on mimic iii.
- Guillaume Geri, Loic Ferrer, Nam Tran, Leo A Celi, Matthieu Jamme, Joon Lee, and Antoine Vieillard-Baron. Cardio-pulmonary-renal interactions in icu patients. role of mechanical ventilation, venous congestion and perfusion deficit on worsening of renal function: Insights from the mimic-iii database. *Journal of critical care*, 64:100–107, 2021.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23): e215–e220, 2000.



560 Mehak Gupta, Brennan Gallamoza, Nicolas Cutrona, Pranjal Dhakal, Raphael Poulain, and Rahmatollah  
561 Beheshti. An extensive data processing pipeline for mimic-iv. In *Machine Learning for Health*, pages  
562 311–325. PMLR, 2022.

563 Nianzong Hou, Mingzhe Li, Lu He, Bing Xie, Lin Wang, Rumin Zhang, Yong Yu, Xiaodong Sun, Zhengsheng  
564 Pan, and Kai Wang. Predicting 30-days mortality for mimic-iii patients with sepsis-3: a machine learning  
565 approach using xgboost. *Journal of translational medicine*, 18(1):1–14, 2020.

566 Jinmiao Huang, Cesar Osorio, and Luke Wicent Sy. An empirical evaluation of deep learning for icd-9 code  
567 assignment using mimic-iii clinical notes. *Computer methods and programs in biomedicine*, 177:141–153,  
568 2019.

569 Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard,  
570 Benjamin Moody, Brian Gow, Li-wei H Lehman, et al. Mimic-iv, a freely accessible electronic health record  
571 dataset. *Scientific data*, 10(1):1, 2023.

572 Guilan Kong, Ke Lin, and Yonghua Hu. Using machine learning methods to predict in-hospital mortality of  
573 sepsis patients in the icu. *BMC medical informatics and decision making*, 20:1–10, 2020.

574 Angelina Prima Kurniati, Geoff Hall, David Hogg, and Owen Johnson. Process mining in oncology using the  
575 mimic-iii dataset. In *Journal of Physics: Conference Series*, volume 971, page 012008. IOP Publishing, 2018.

576 Fei Li and Hong Yu. Icd coding from clinical text using multi-filter residual convolutional neural network. In  
577 *proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8180–8187, 2020.

578 Fuhai Li, Hui Xin, Jidong Zhang, Mingqiang Fu, Jingmin Zhou, and Zhexun Lian. Prediction model of in-  
579 hospital mortality in intensive care unit patients with heart failure: machine learning-based, retrospective  
580 analysis of the mimic-iii database. *BMJ open*, 11(7):e044779, 2021.

581 Min Li, Zhihui Fei, Min Zeng, Fang-Xiang Wu, Yaohang Li, Yi Pan, and Jianxin Wang. Automated icd-9 coding  
582 via a deep learning approach. *IEEE/ACM transactions on computational biology and bioinformatics*, 16(4):  
583 1193–1202, 2018.

584 Ke Lin, Yonghua Hu, and Guilan Kong. Predicting in-hospital mortality of patients with acute kidney injury in  
585 the icu using random forest model. *International journal of medical informatics*, 125:55–61, 2019.

586 Dongliang Liu, Yiyang Tang, and Qian Zhang. Admission hyperglycemia predicts long-term mortality in  
587 critically ill patients with subarachnoid hemorrhage: a retrospective analysis of the mimic-iii database.  
588 *Frontiers in Neurology*, 12:678998, 2021a.

589 Tao Liu, Haochen Xuan, Lili Wang, Xiaoqun Li, Zhihao Lu, Zhaoxuan Tian, Junhong Chen, Chaofan Wang,  
590 Dongye Li, and Tongda Xu. The association between serum albumin and long length of stay of patients with  
591 acute heart failure: A retrospective study based on the mimic-iv database. *Plos one*, 18(2):e0282289, 2023.

592 Taotao Liu, Qinyu Zhao, and Bin Du. Effects of high-flow oxygen therapy on patients with hypoxemia after  
593 extubation and predictors of reintubation: a retrospective study based on the mimic-iv database. *BMC*  
594 *Pulmonary Medicine*, 21(1):1–15, 2021b.

595 Zongqing Lu, Jin Zhang, Jianchao Hong, Jiatian Wu, Yu Liu, Wenyan Xiao, Tianfeng Hua, and Min Yang.  
596 Development of a nomogram to predict 28-day mortality of patients with sepsis-induced coagulopathy: an  
597 analysis of the mimic-iii database. *Frontiers in medicine*, 8:661710, 2021.

598 Mengdi Luo, Yang Chen, Yuan Cheng, Na Li, and He Qing. Association between hematocrit and the 30-day  
599 mortality of patients with sepsis: A retrospective analysis based on the large-scale clinical database mimic-iv.  
600 *PloS one*, 17(3):e0265758, 2022.

601 Andrés Alejandro Ramos Magna, Héctor Allende-Cid, Carla Taramasco, Carlos Becerra, and Rosa L Figueroa.  
602 Application of machine learning and word embeddings in the classification of cancer diagnosis using patient  
603 anamnesis. *Ieee Access*, 8:106198–106213, 2020.

604 Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and  
605 fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35, 2021.

606 Chuizheng Meng, Loc Trinh, Nan Xu, James Enouen, and Yan Liu. Interpretability and fairness evaluation of  
607 deep learning models on mimic-iv dataset. *Scientific Reports*, 12(1):7166, 2022.

608 Andi Partovi, Dickson Lukose, and Geoffrey I Webb. Mipy: A framework for benchmarking machine learning  
609 prediction of unplanned hospital and icu readmission in the mimic-iv database. 2022.

610 Alvin Rajkomar, Michaela Hardt, Michael D Howell, Greg Corrado, and Marshall H Chin. Ensuring fairness in  
611 machine learning to advance health equity. *Annals of internal medicine*, 169(12):866–872, 2018.

612 Juan C Rojas, Kyle A Carey, Dana P Edelson, Laura R Venable, Michael D Howell, and Matthew M Churpek.  
613 Predicting intensive care unit readmission with machine learning using electronic health record data. *Annals  
614 of the American Thoracic Society*, 15(7):846–853, 2018.

615 Eliane Rösli, Brian Rice, and Tina Hernandez-Boussard. Bias at warp speed: how ai may contribute to the  
616 disparities gap in the time of covid-19. *Journal of the American Medical Informatics Association*, 28(1):  
617 190–192, 2021.

618 Mohammed Sayed, David Riano, and Jesús Villar. Predicting duration of mechanical ventilation in acute  
619 respiratory distress syndrome using supervised machine learning. *Journal of Clinical Medicine*, 10(17):3824,  
620 2021.

621 Tingting Shu, Jian Huang, Jiewen Deng, Huaqiao Chen, Yang Zhang, Minjie Duan, Yanqing Wang, Xiaofei Hu,  
622 and Xiaozhu Liu. Development and assessment of scoring model for icu stay and mortality prediction after  
623 emergency admissions in ischemic heart disease: a retrospective study of mimic-iv databases. *Internal and  
624 Emergency Medicine*, 18(2):487–497, 2023.

625 AK Singh, Mounika Guntu, Ananth Reddy Bhimireddy, Judy W Gichoya, and Saptarshi Purkayastha. Multi-label  
626 natural language processing to identify diagnosis and procedure codes from mimic-iii inpatient notes. *arXiv  
627 preprint arXiv:2003.07507*, 2020.

628 Chun Sun, Deqing Chen, Xin Jin, Guangtao Xu, Chenye Tang, Xiao Guo, Zhiling Tang, Yixin Bao, Fei Wang,  
629 and Ruilin Shen. Association between acute kidney injury and prognoses of cardiac surgery patients: Analysis  
630 of the mimic-iii database. *Frontiers in Surgery*, 9, 2022.

631 Jared Thacker. *A Machine Learning Pipeline for Readmission Prediction with MIMIC-III*. PhD thesis, Auburn  
632 University, 2023.

633 Julian Theis, William L Galanter, Andrew D Boyd, and Houshang Darabi. Improving the in-hospital mortality  
634 prediction of diabetes icu patients using a process mining/deep learning architecture. *IEEE Journal of  
635 Biomedical and Health Informatics*, 26(1):388–399, 2021.

636 Dongyan Wang, Xiaoyan Guo, Wenwen Xia, Zhijuan Ru, Yihai Shi, and Zhengyu Hu. Effect of admission  
637 serum calcium levels and length of stay in patients with acute pancreatitis: Data from the mimic-iii database.  
638 *Emergency Medicine International*, 2022, 2022.

639 Siyue Yang, Paul Varghese, Ellen Stephenson, Karen Tu, and Jessica Gronsbell. Machine learning approaches  
640 for electronic health records phenotyping: a methodical review. *Journal of the American Medical Informatics  
641 Association*, 30(2):367–381, 2023a.

642 Wei Yang, Hong Zou, Meng Wang, Qin Zhang, Shadan Li, and Hongyin Liang. Mortality prediction among icu  
643 inpatients based on mimic-iii database results from the conditional medical generative adversarial network.  
644 *Heliyon*, 9(2), 2023b.

645 Zixiang Ye, Shuoyan An, Yanxiang Gao, Enmin Xie, Xuecheng Zhao, Ziyu Guo, Yike Li, Nan Shen, Jingyi Ren,  
646 and Jingang Zheng. The prediction of in-hospital mortality in chronic kidney disease patients with coronary  
647 artery disease using machine learning models. *European Journal of Medical Research*, 28(1):1–13, 2023.

648 Suru Yue, Shasha Li, Xueying Huang, Jie Liu, Xuefei Hou, Yumei Zhao, Dongdong Niu, Yufeng Wang, Wenkai  
649 Tan, and Jiayuan Wu. Machine learning for the prediction of acute kidney injury in patients with sepsis.  
650 *Journal of translational medicine*, 20(1):1–12, 2022.

651 HuanRui Zhang, Wen Tian, and YuJiao Sun. The value of anion gap for predicting the short-term all-cause  
652 mortality of critically ill patients with cardiac diseases, based on mimic-iii database. *Heart & Lung*, 55:59–67,  
653 2022a.

654 Jingqing Zhang, Luis Daniel Bolanos Trujillo, Ashwani Tanwar, Julia Ive, Vibhor Gupta, and Yike Guo. Clinical  
655 utility of automatic phenotype annotation in unstructured clinical notes: intensive care unit use. *BMJ Health  
656 & Care Informatics*, 29(1):e100519, 2022b.