



电子科技大学

University of Electronic Science and Technology of China

Reproduced Analysis for “Integrated multimodal artificial intelligence framework for healthcare applications”

Author: 黄其涵, Adviser: 毛晓伟

February 26, 2024



Catalogue

► Intro

► Methods

► Results



Background

Artificial intelligence (AI) systems hold great promise to improve healthcare over the next decades. Specifically, AI systems leveraging multiple data sources and input modalities are poised to become a viable method to deliver more accurate results and deployable pipelines across a wide range of applications [1].

[1] Soenksen, Luis R., et al. "Integrated multimodal artificial intelligence framework for healthcare applications." *NPJ digital medicine* 5.1 (2022): 149.



Challenges of Single-source AI Systems

The specific challenges faced by AI systems relying on single data sources:

1. Accuracy limitations in complex cases.
2. Lack of generalizability across diverse patient populations and conditions.
3. Reduced applicability in multifaceted healthcare scenarios.



Contributions of HAIM System

The contributions of Holistic AI in Medicine (HAIM) can be summarized as follows:

1. This work introduced the HAIM framework, **integrating multiple data sources and modalities into a cohesive system for constructing AI applications in healthcare.**
2. This research shows that **multiple-source or models developed within the HAIM framework consistently outperform single-source models** across various healthcare applications, including diagnostics and prognostics. It also emphasizes **the utility of multimodal data in enhancing the accuracy and reliability of AI systems in medicine.**
3. **Utilizing Shapley values to quantify the contributions of each data modality and source**, this provide insights into the heterogeneity of data modality importance.



Catalogue

► Intro

► Methods

► Results



Datasets

MIMIC-IV: The Medical Information Mart for Intensive Care IV (MIMIC-IV) is a large database containing de-identified health-related data associated with over forty thousand patients who stayed in critical care units of the Beth Israel Deaconess Medical Center between 2008 and 2019 [2].

MIMIC-CXR-JPG: The MIMIC Chest X-ray JPG (MIMIC-CXR-JPG) Database v2.0.0 is a large publicly available dataset of chest radiographs in JPG format with structured labels derived from free-text radiology reports. The MIMIC-CXR-JPG dataset is wholly derived from MIMIC-CXR, providing JPG format files derived from the DICOM images and structured labels derived from the free-text reports [3].

[2] Johnson, A., Bulgarelli, L., Pollard, T., Horng, S., Celi, L. A., & Mark, R. (2021). MIMIC-IV (version 1.0). PhysioNet. <https://doi.org/10.13026/s6n6-xd98>.

[3] Johnson, A., Lungren, M., Peng, Y., Lu, Z., Mark, R., Berkowitz, S., & Horng, S. (2019). MIMIC-CXR-JPG - chest radiographs with structured labels (version 2.0.0). PhysioNet. <https://doi.org/10.13026/8360-t248>.



Datasets

HAIM-MIMIC-MM: A unified multimodal dataset created by combining the MIMIC-IV v1.0 and the MIMIC-CXR-JPG database v2.0.0. This dataset mainly uses `subject_id`, `hadm_id`, `stay_id` as the primary key to create links in the tables of each database:

- **subject_id:** Unique identifier for a patient. This ID ensures patient anonymity while allowing researchers to link different hospital visits or admissions for the same individual.
- **hadm_id:** Unique identifier for a patient's hospital admission. This ID is used to differentiate between multiple admissions for the same patient, facilitating the study of each separate hospital stay.
- **stay_id:** Unique identifier for a patient's stay in a specific ward or ICU during a hospital admission. It allows for the analysis of patient care at a more granular level, focusing on individual segments of the hospital stay.



Datasets

The final HAIM data unit object ->

```
1 class Patient_ICU(object):
2     def __init__(self, admissions, demographics, transfers, core,
3                 diagnoses_icd, drgcodes, emar, emar_detail, hcpcsevents,
4                 labevents, microbiologyevents, poe, poe_detail,
5                 prescriptions, procedures_icd, services, procedureevents,
6                 outputevents, inputevents, icustays, datetimeevents,
7                 chartevents, cxr, imcxr, noteevents, dsnotes, ecgnotes,
8                 echonotes, radnotes):
9
10    ## CORE
11    self.admissions = admissions # Patient admissions information
12    self.demographics = demographics # Patient demographics data
13    self.transfers = transfers # Patient transfer data within hospital
14    self.core = core # Core patient information
```



Datasets

```
1 ## HOSP
2 self.diagnoses_icd = diagnoses_icd # Diagnoses in ICD format
3 self.drgcodes = drgcodes # DRG codes for billing and reimbursements
4 self.emar = emar # Electronic Medication Administration Record
5 self.emar_detail = emar_detail
6 # Detailed medication administration data
7 self.hcpcsevents = hcpcsevents
8 # Healthcare Common Procedure Coding events
9 self.labevents = labevents # Laboratory test results
10 self.microbiologyevents = microbiologyevents
11 # Microbiology test results
12 self.poe = poe # Physician Order Entry records
13 self.poe_detail = poe_detail # Detailed Physician Order Entry records
14 self.prescriptions = prescriptions # Prescribed medications
15 self.procedures_icd = procedures_icd # Procedures coded in ICD format
16 self.services = services # Services provided to the patient
```



Datasets

```
1 ## ICU
2 self.procedureevents = procedureevents # Procedures performed in ICU
3 self.outputevents = outputevents # Outputs like fluids measured
4 self.inputevents = inputevents # Inputs like medications administered
5 self.icustays = icustays # Information on each ICU stay
6 self.datetimeevents = datetimeevents # Timestamped events
7 self.charthevents = charthevents # Charted events like vital signs
8
9 ## CXR
10 self.cxr = cxr # Chest X-ray data
11 self.imcxr = imcxr # Image data from chest X-rays
12
13 ## NOTES
14 self.noteevents = noteevents # Clinical notes
15 self.dsnotes = dsnotes # Discharge summaries or other clinical notes
16 self.ecgnotes = ecgnotes # Electrocardiogram notes
17 self.echonotes = echonotes # Echocardiogram notes
18 self.radnotes = radnotes # Radiology reports and notes
```



Data Preprocessing

The generation of embeddings from input modalities encompasses a diverse range of data types. These modalities and their respective embeddings include:

- **Tabular Data:** Demographics (E_{de}) are represented as tabular data.
- **Structured Time-Series Events:**
 - Chart events (E_{ce}),
 - Laboratory events (E_{le}),
 - Procedure events (E_{pe}).
- **Unstructured Free Text:**
 - Radiological notes (E_{radn}),
 - Electrocardiogram notes (E_{ecgn}),
 - Echocardiogram notes (E_{econ}).
- **Single-Image Vision:**
 - Visual probabilities (E_{vp}),
 - Visual dense-layer features (E_{vd}).



Data Preprocessing

- Multi-Image Vision:

- Aggregated visual probabilities (E_{vmp}),
- Aggregated visual dense-layer features (E_{vmd}).

Table 1. General characteristics of the HAIM-MIMIC-MM database.

Characteristic	MIMIC-IV-MM
# Samples	34537
# Demographic Variables	6
# Chart Event Variables	9
# Laboratory Event Variables	23
# Procedure Event Variables	10
# X-ray Variables	1
# Text Note Variables	3



DataFrame Embeddings

The following data frames contain embeddings extracted from various patient data, where each prefix in the column names represents the type of embedding:

- `df_demographics_embeddings_fusion`: Contains embeddings from demographic information of patients. Prefix `de_` stands for *demographics embeddings*.
- `df_ts_ce_embeddings_fusion`: Contains time-series feature embeddings extracted from chart events of patients. Prefix `ts_ce_` stands for *time series chart events*.
- `df_ts_le_embeddings_fusion`: Contains time-series feature embeddings extracted from laboratory events of patients. Prefix `ts_le_` stands for *time series lab events*.
- `df_ts_pe_embeddings_fusion`: Contains time-series feature embeddings extracted from procedure events of patients. Prefix `ts_pe_` stands for *time series procedure events*.
- `df_vision_dense_embeddings_fusion`: Contains dense feature embeddings extracted from single chest X-ray images. Prefix `vd_` stands for *vision dense*.



DataFrame Embeddings

- `df_vision_predictions_embeddings_fusion`: Contains predictive feature embeddings extracted from single chest X-ray images. Prefix `vp_` stands for *vision predictions*.
- `df_vision_multi_dense_embeddings_fusion`: Contains dense feature embeddings accumulated from multiple chest X-ray images. Prefix `vmd_` stands for *vision multi dense*.
- `df_vision_multi_predictions_embeddings_fusion`: Contains predictive feature embeddings accumulated from multiple chest X-ray images. Prefix `vmp_` likely stands for *vision multi predictions*.
- `df_ecgnotes_embeddings_fusion`: Contains embeddings extracted from ECG notes. Prefix `n_ecg_` stands for *notes ECG*.
- `df_echonotes_embeddings_fusion`: Contains embeddings extracted from echocardiogram notes. Prefix `n_ech_` stands for *notes echocardiogram*.
- `df_radnotes_embeddings_fusion`: Contains embeddings extracted from radiology reports. Prefix `n_rad_` stands for *notes radiology*.



Modeling

Length of Stay Modeling

This work framed the prediction of patient discharge within the next 48 hours as a binary classification problem: discharged alive within 48 hours (1) or otherwise (0). Each sample in this predictive task corresponds to a single patient-admission EHR time point where an X-ray image was obtained ($N = 45,050$).

```
1 Main function:
2     # load the HAIM embedding file
3     fname = 'data/cxr_ic_fusion.csv'
4     df = pd.read_csv(fname, skiprows=[45051, 45052])
5     # Tagging and labelling
6     df_alive_small48 = df[((df['img_length_of_stay'] < 48) & (df['
    death_status'] == 0))]
7     df_alive_big48 = df[((df['img_length_of_stay'] >= 48) & (df['
    death_status'] == 0))]
8     df_death = df[(df['death_status'] == 1)]
9     df_alive_small48['y'] = 1
10    df_alive_big48['y'] = 0
11    df_death['y'] = 0
12    df = pd.concat([df_alive_small48, df_alive_big48, df_death], axis
    =0)
```




Modeling

Length of Stay Modeling

```
1 Main function:
2 .....
3 # Remove unnecessary columns
4 df = df.drop(['img_id', 'img_charttime', 'img_deltacharttime', '
5             'discharge_location', 'img_length_of_stay',
6             'death_status'], axis=1)
7
8 # Get all types of data sources, 'de_', 'vd_', 'vp_', 'vmd_',...
9 data_type_dict = get_data_dict(df)
10
11 # Combine and arrange all types of input samples
12 all_types_experiment = get_all_dtypes()
13
14 # Number of Cases: 2047
15 # Run multithreaded execution task
16 results = parallel_run(all_types_experiment, data_type_dict, df,
17                         'lengthOfStay', start_index=22)
```



Modeling

Length of Stay Modeling

This function is designed to run a grid search with cross-validation on a set of hyperparameters for the XGBoost classifier, optimizing for the ROC AUC metric.

```
1 def run_xgb(x_train, y_train, x_test):
2     cv_folds = 5
3     gs_metric = 'roc_auc'
4     param_grid = {'max_depth': [5, 6, 7, 8],
5                   'n_estimators': [200, 300],
6                   'learning_rate': [0.3, 0.1, 0.05],}
7
8     est = xgb.XGBClassifier(verbosity=0, scale_pos_weight=(len(
9         y_train) - sum(y_train)) / sum(y_train), seed=42,
10                           tree_method='hist', gpu_id=0, eval_metric
11                           = 'logloss')
12
13     gs = GridSearchCV(estimator=est, param_grid=param_grid, scoring=
14                       gs_metric, cv=cv_folds)
15     gs.fit(x_train, y_train)
16     .....
```



Modeling

Mortality Prediction Modeling

This work predicted if a patient would pass away within the next 48 hours, treating it as a binary classification: expiring within 48 hours (1) or not (0). For patients not marked for expiration, the class label is 0.

```
1 Main function:
2 .....
3 fname = 'data/cxr_ic_fusion.csv'
4 df = pd.read_csv(fname, skiprows=[45051, 45052])
5
6 df_death_small48 = df[((df['img_length_of_stay'] < 48) & (df['
  death_status'] == 1))]
7 df_alive_big48 = df[((df['img_length_of_stay'] >= 48) & (df['
  death_status'] == 0))]
8 df_death_big48 = df[((df['img_length_of_stay'] >= 48) & (df['
  death_status'] == 1))]
9
10 df_death_small48['y'] = 1
11 df_alive_big48['y'] = 0
12 df_death_big48['y'] = 0
13 .....
```



Modeling Pathology Diagnosis Modeling

This work targeted predicting 10 common chest pathologies (e.g., fractures, lung lesions) using the HAIM framework, which showed superior performance over image-only methods. The dataset for this study, derived from MIMIC-CXR-JPG v2.0.0, labeled pathologies as present (1), absent (0), or inconclusive (-1); we excluded inconclusive cases. Excluding unstructured radiology notes to prevent overfitting, we utilized multimodal inputs for binary classification of each pathology. Sample sizes varied by pathology, with totals ranging from 557 (Fracture) to 18,571 (Cardiomegaly).



Catalogue

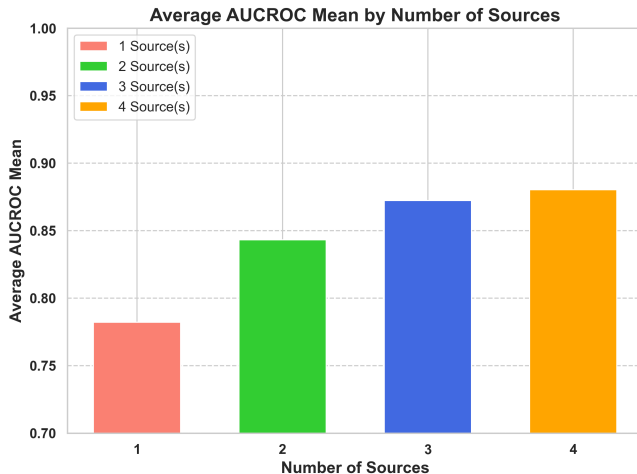
► Intro

► Methods

► Results

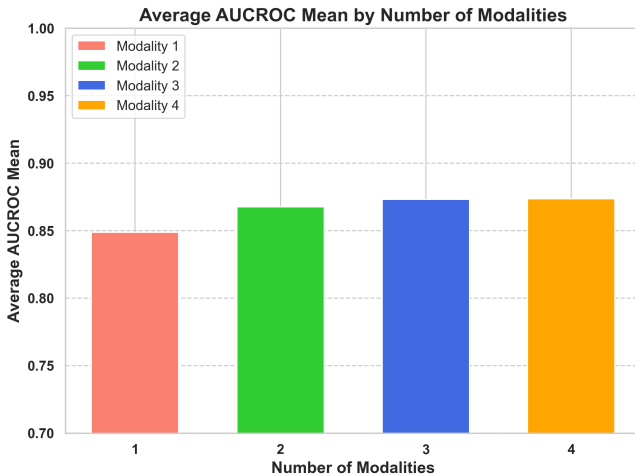


Multi-source vs. Single-source Length of Stay Modeling





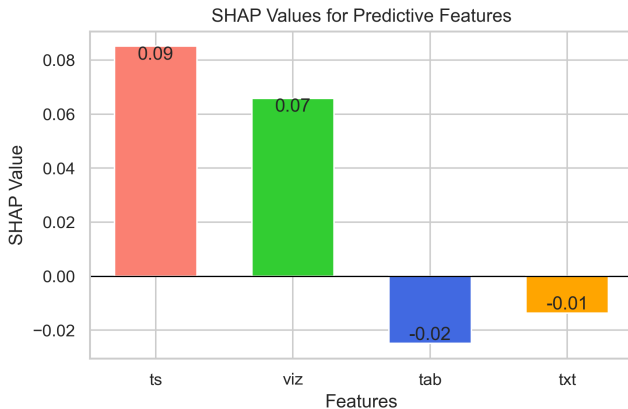
Multimodality vs. Single Modality Length of Stay Modeling





Quantification by Shapley Values

Length of Stay Modeling





Q&A

Thanks for Listening