

# A Survey of Task-driven Heterogeneous Feature Embedding and Selection

Hedong YAN,  
Computer Science,  
Hong Kong Baptist University  
Supervisor: Yiu-ming Cheung

- Background
- Related Works
  - Feature Embedding
  - Feature Selection
- Methodology
- Futural Plan

# Background

- Heterogeneous data widely exists in reality, such as user information, EHR, and surveys.
- Heterogeneous data is critical for many tasks in the real world.



CTR

VS.

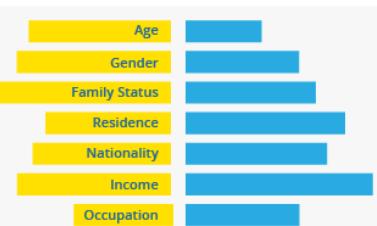


Conversion



Laboratory

**CTR and Conversion  
Rate Prediction**

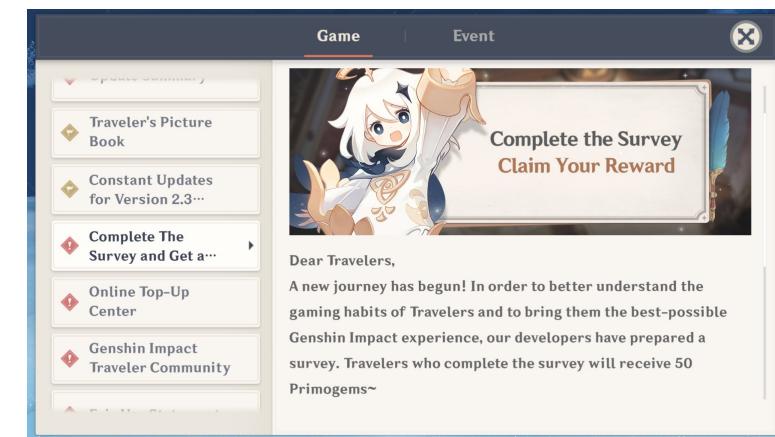


Demographic



Vital Sign

**Disease Progression  
Prediction**



**Survey Analysis**

# Problem

- Traditional deep learning models for homogeneous features can not be directly applied to heterogeneous data. Not much attention has been paid to describing how DNN can be designed for heterogeneous datasets.



Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huurdlest gelburn"? Kijf – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huurdlest gelburn"? Kijf – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Research Question 1: How to design embedding modules for heterogeneous dataset?



Video

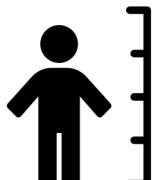


Voice

**Homogeneous Features**



Interval



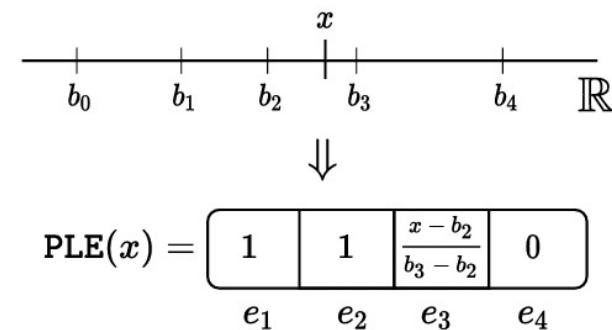
Ratio

**Heterogeneous Features**

# Related works

| Feature Scale | Encoder               | Input            | Output   |
|---------------|-----------------------|------------------|--|
| Nominal       | One-hot               | [1,2,3]          | [[1,0,0],[0,1,0],[0,0,1]]  |
|               | Binary                | [1,2,3]          | [[0,0],[0,1],[1,0]]  |
|               | Dumpy                 | [1,2,3]          | [[1,0],[0,1],[0,0]]  |
|               | Count                 | [1,1,3]          | [[2],[2],[1]]  |
|               | Simple                | [1,2,3]          | $[\left[\frac{2}{3}, -\frac{1}{3}, -\frac{1}{3}\right], \left[-\frac{1}{3}, \frac{2}{3}, -\frac{1}{3}\right], \left[-\frac{1}{3}, -\frac{1}{3}, \frac{2}{3}\right]]$ |
| Ordinal       | Ordinal               | [1,2,3]          | [1,2,3]  |
|               | Rank-hot              | [1,2,3]          | [[1,0,0],[1,1,0],[1,1,1]]  |
|               | Gray                  | [1,2,3]          | [[0,0],[0,1],[1,1]]  |
| Continuous    | Bins + One-hot        | [0.11,0.27,0.34] | [[1,0,0],[0,1,0],[0,0,1]]  |
|               | Piece-wise linear [1] | [0.11,0.27,0.34] | [[0.1,0,0],[1,0.2,0],[1,1,0.1]]  |

- Can we use existing encoders to transform the heterogeneous feature into homogeneous features?



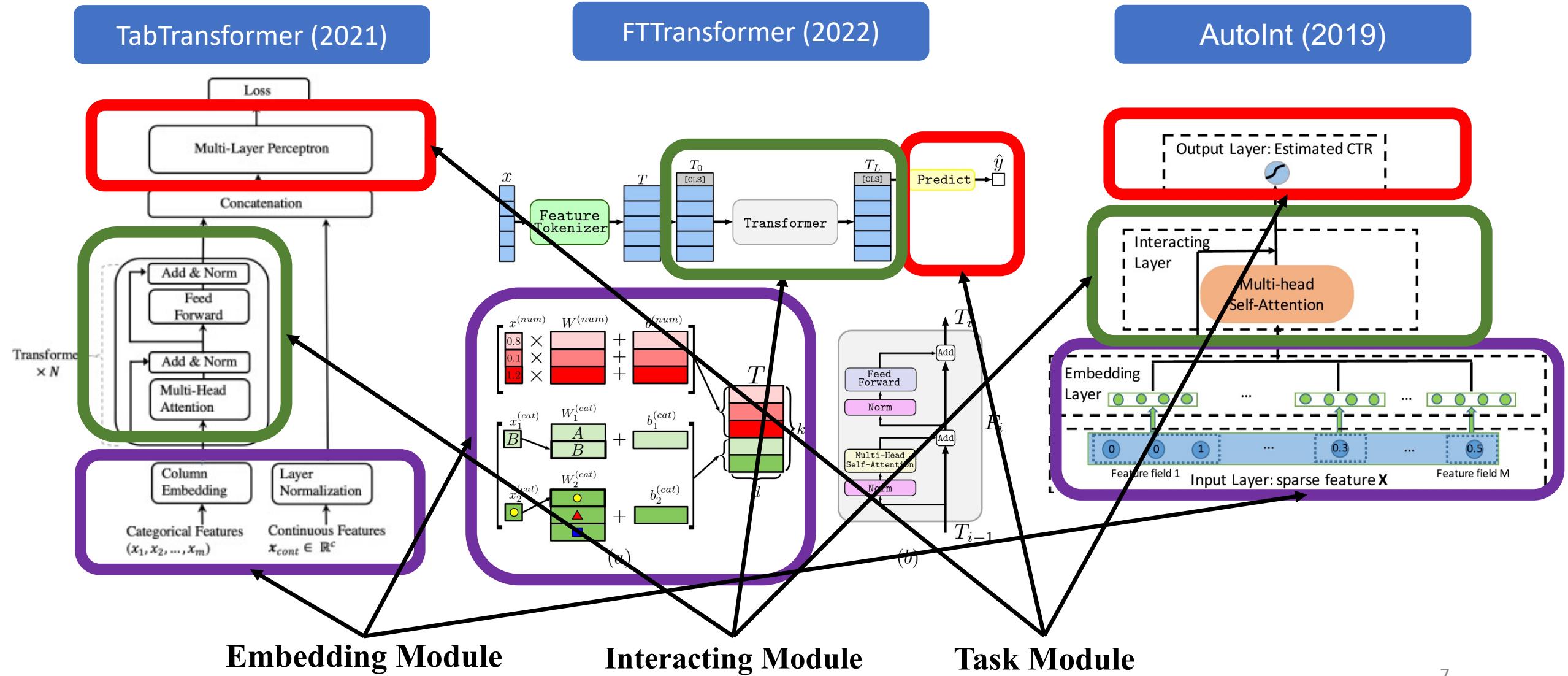
# Related works

## Heterogeneous embedding for different models

- Transformer-based model
  - TabTransformer, FTTransformer, AutoInt, ILEAHE
- MLP-based model
  - DeepFM, DANETs, DVN v2
- Diffusion-based model
  - TabDDPM
- Graph-based model
  - T2G-Former

# Related works

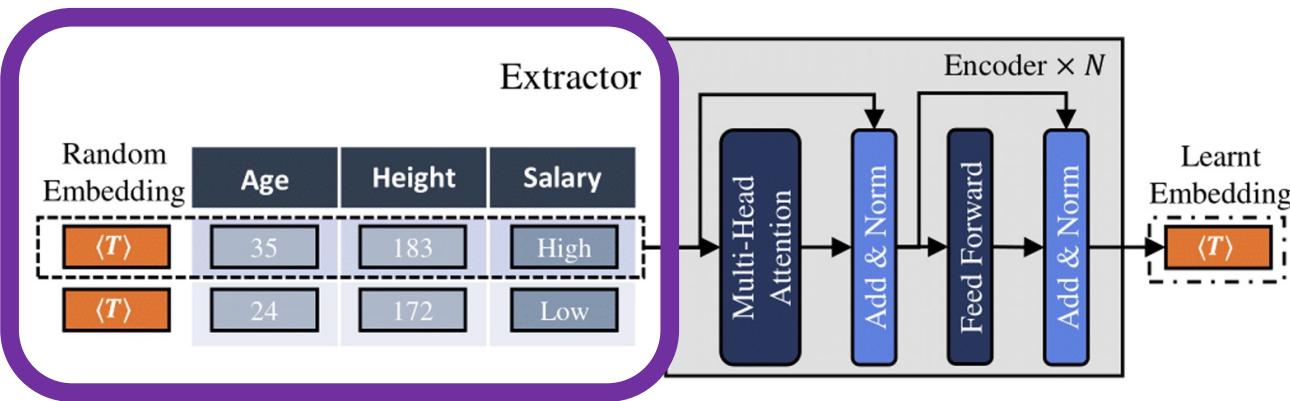
## Transformer-based model



# Related works

## Transformer-based model

ILEAHE (2023)



Categorical: Dictionary embedding

Numerical: 2-layer perceptron

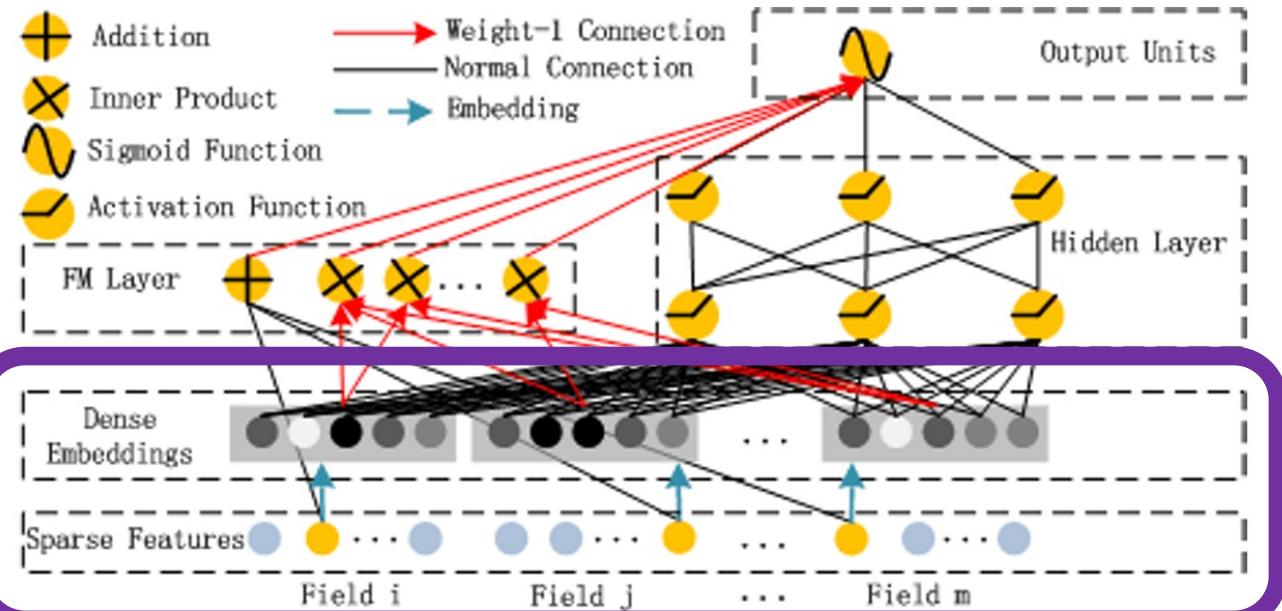
### Heterogeneous Embedding Modules

- TabTransformer
  - Categorical: Dictionary embedding
  - Continuous: None
- FTTransformer
  - Categorical: Dictionary embedding
  - Continuous: Linear
- AutoInt
  - Categorical: One-hot + linear
  - Continuous: Linear
- ILEAHE
  - Categorical: Dictionary embedding
  - Continuous: 2-layer perceptron

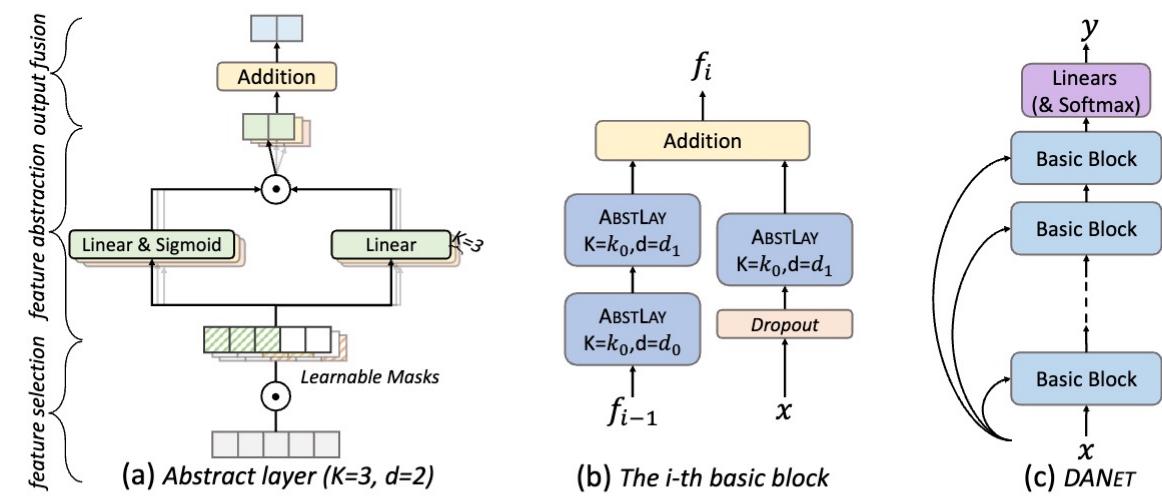
# Related works

## MLP-based model

DeepFM (2017)



DANETs (2017)



Categorical: One-hot encoder + linear

Numerical: Linear

$$x_k^i = \frac{\sum_{j=1}^n \mathbb{I}_{x_j^i=x_k^i} * y_j + ap}{\sum_{j=1}^n \mathbb{I}_{x_j^i=x_k^i} + a}$$

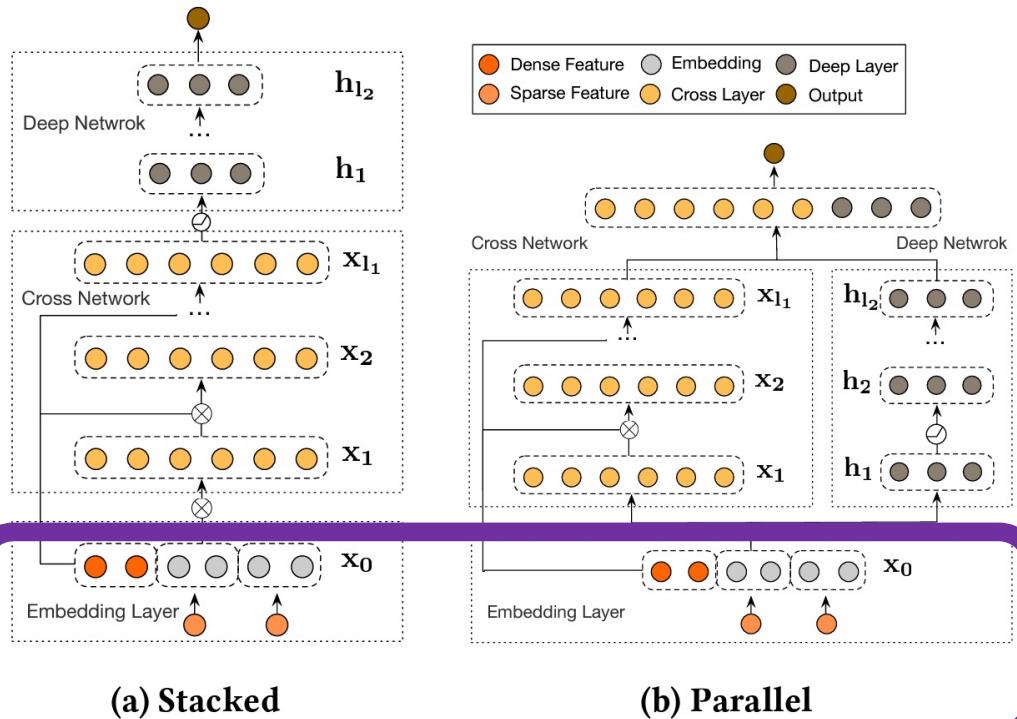
Categorical: Target Statistic

Numerical: None

# Related works

## MLP-based model

DVN v2 (2020)



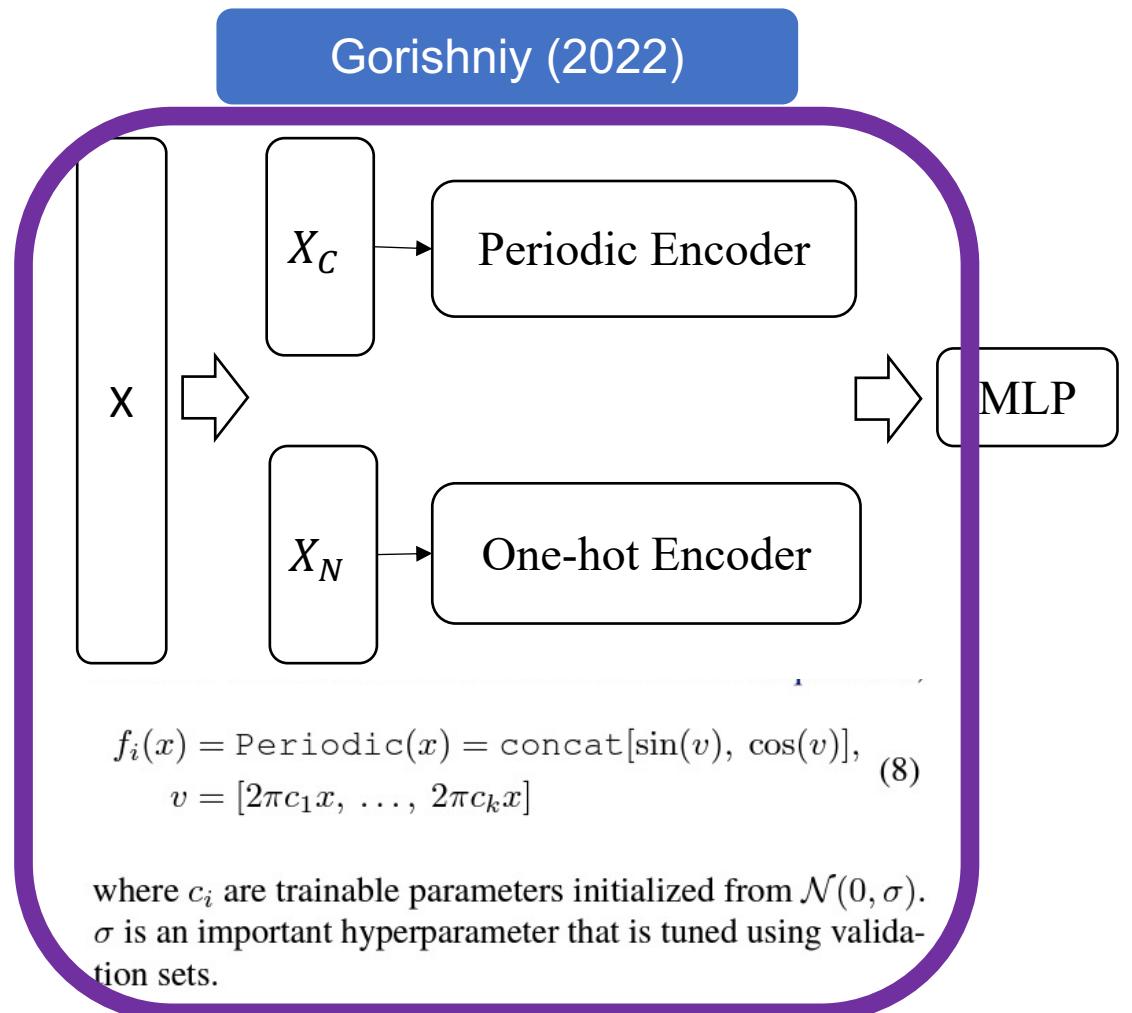
(a) Stacked

(b) Parallel

Categorical: Dictionary embedding

Numerical: None

Gorishniy (2022)



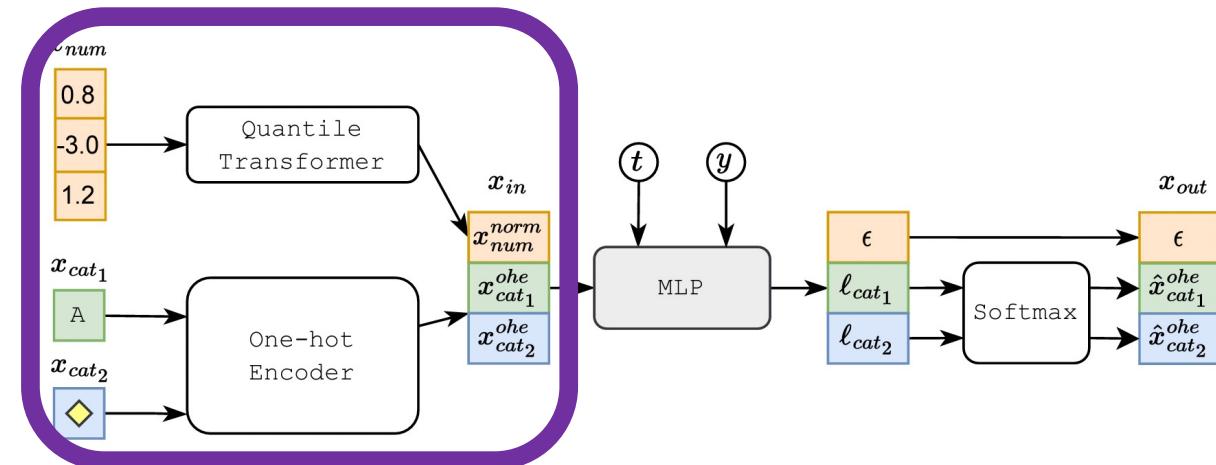
Categorical: One-hot

Numerical: Periodic encoder

# Related works

## Diffusion-based model

TabDDPM (2023)



## Graph-based

T2G-Former (2023)

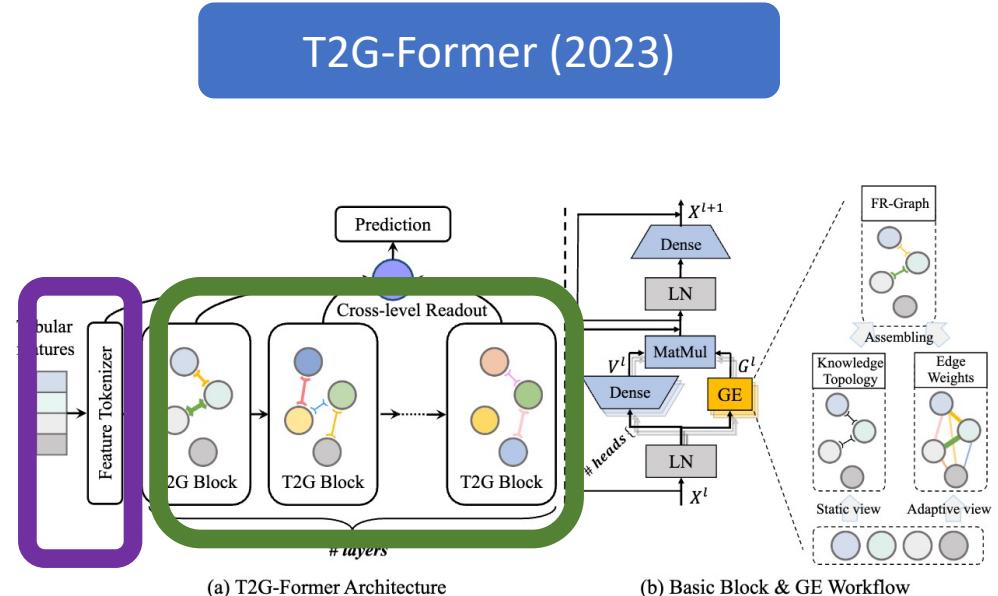


Figure 2: (a) The architecture of T2G-FORMER for tabular learning. Each T2G block builds an FR-Graph for a feature level and performs selective interaction. A global readout node collects salient features from each layer to form tabular semantics. (b) Illustrating a basic block in Sec. and GE in Sec. .

- Use diffusion procedure to optimize the parameters
- Categorical: One-hot
- Numerical: Quantile Gaussian Normalization

- Add graph blocks to model the features' interaction
- Categorical: Dictionary embedding
- Numerical: Linear

# Related works

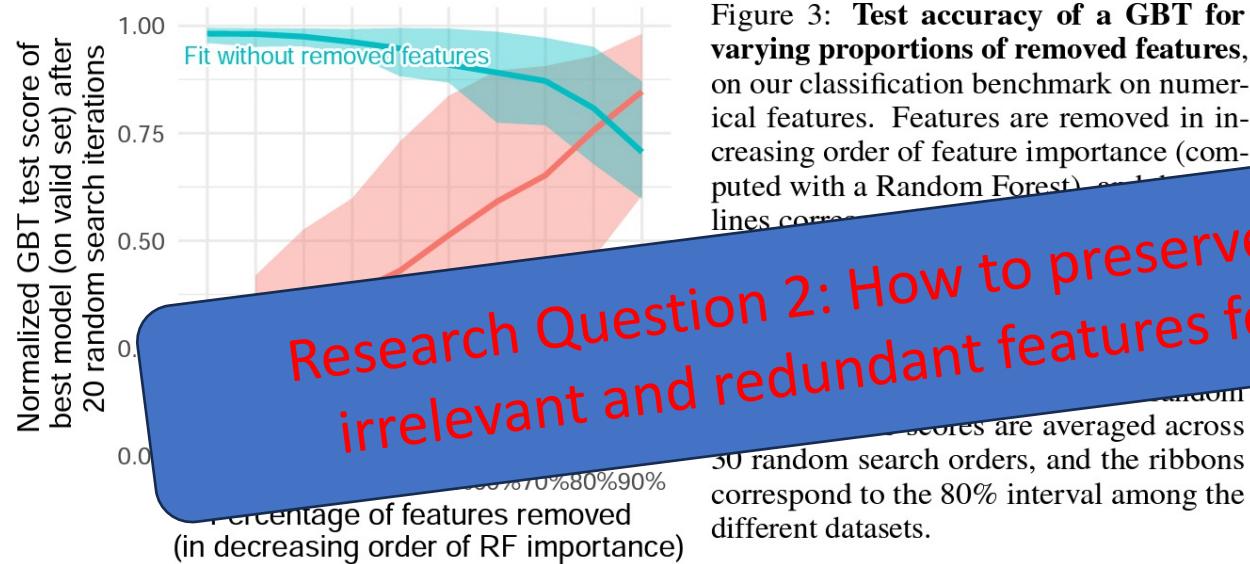
## Section Conclusion

|                  | GE ↑         | CH ↑         | CA ↓         | HO ↓         | AD ↑         | OT ↑         | HI ↑         | FB ↓         | SA ↑         | CO ↑         | MI ↓         | Avg. Rank |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-----------|
| CatBoost         | 0.692        | 0.861        | 0.430        | 3.093        | 0.873        | 0.825        | 0.727        | 5.226        | 0.924        | 0.967        | <b>0.741</b> | 3.6 ± 2.9 |
| XGBoost          | 0.683        | 0.859        | 0.434        | 3.152        | <b>0.875</b> | 0.827        | 0.726        | 5.338        | 0.919        | 0.969        | 0.742        | 4.6 ± 2.7 |
| MLP              | 0.665        | 0.856        | 0.486        | 3.109        | 0.856        | 0.822        | 0.727        | 5.616        | 0.913        | 0.968        | 0.746        | 8.5 ± 2.6 |
| MLP-LR           | 0.679        | 0.861        | 0.463        | 3.012        | 0.859        | 0.826        | 0.731        | 5.477        | 0.924        | 0.972        | 0.744        | 5.5 ± 2.7 |
| MLP-Q-LR         | 0.682        | 0.859        | 0.433        | 3.080        | 0.867        | 0.818        | 0.724        | <b>5.144</b> | 0.924        | 0.974        | 0.745        | 5.1 ± 1.9 |
| MLP-T-LR         | 0.673        | 0.861        | 0.435        | 3.099        | 0.870        | 0.821        | 0.727        | 5.409        | 0.924        | 0.973        | 0.746        | 5.1 ± 1.7 |
| MLP-PLR          | <b>0.700</b> | 0.858        | 0.453        | <b>2.975</b> | 0.874        | <b>0.830</b> | <b>0.734</b> | 5.388        | <b>0.924</b> | 0.975        | 0.743        | 3.0 ± 2.4 |
| ResNet           | 0.690        | 0.861        | 0.483        | 3.081        | 0.856        | 0.821        | 0.734        | 5.482        | 0.918        | 0.968        | 0.745        | 6.7 ± 3.3 |
| ResNet-LR        | 0.672        | 0.862        | 0.450        | 2.992        | 0.859        | 0.822        | 0.733        | 5.415        | 0.923        | 0.971        | 0.743        | 5.6 ± 2.7 |
| ResNet-Q-LR      | 0.674        | 0.859        | 0.427        | 3.066        | 0.868        | 0.815        | 0.729        | 5.309        | 0.923        | 0.976        | 0.746        | 4.7 ± 2.0 |
| ResNet-T-LR      | 0.683        | 0.862        | <b>0.425</b> | 3.030        | 0.872        | 0.822        | 0.731        | 5.471        | 0.923        | 0.975        | 0.744        | 4.1 ± 1.9 |
| ResNet-PLR       | 0.691        | 0.861        | 0.443        | 3.040        | <b>0.874</b> | 0.825        | 0.734        | 5.400        | 0.924        | 0.975        | 0.743        | 3.2 ± 1.3 |
| Transformer-L    | 0.668        | 0.861        | 0.455        | 3.188        | 0.860        | 0.824        | 0.727        | 5.434        | 0.924        | 0.973        | 0.743        | 5.9 ± 2.2 |
| Transformer-LR   | 0.666        | 0.861        | 0.446        | 3.193        | 0.861        | 0.824        | 0.733        | 5.430        | 0.924        | 0.973        | 0.743        | 5.2 ± 2.2 |
| Transformer-Q-LR | 0.690        | 0.857        | <b>0.425</b> | 3.143        | 0.868        | 0.818        | 0.726        | 5.471        | <b>0.924</b> | 0.975        | 0.744        | 4.4 ± 2.2 |
| Transformer-T-LR | 0.686        | 0.862        | <b>0.423</b> | 3.149        | 0.871        | 0.823        | 0.733        | 5.515        | 0.924        | <b>0.976</b> | 0.744        | 3.7 ± 2.2 |
| Transformer-PLR  | 0.686        | <b>0.864</b> | 0.449        | 3.091        | 0.873        | 0.823        | 0.734        | 5.581        | <b>0.924</b> | 0.975        | 0.743        | 3.9 ± 2.5 |

- The key to handling the heterogeneous features is the **embedding layer**
- Resnet and Transformer is not better than MLP with suitable heterogeneous embedding

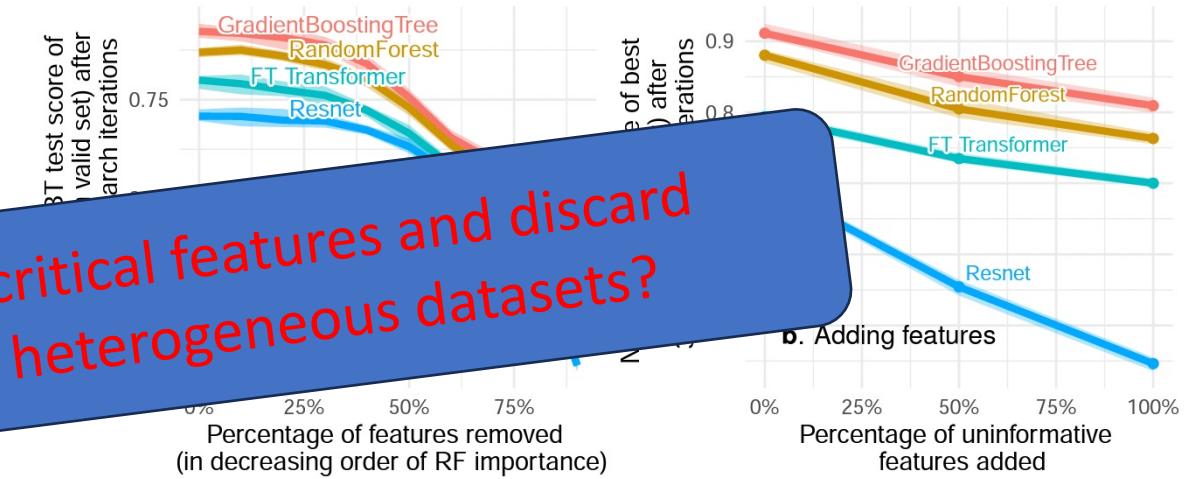
# Problem

- Feature selection is proved critical for heterogeneous datasets.



Research Question 2: How to preserve critical features and discard irrelevant and redundant features for heterogeneous datasets?

Heterogeneous datasets contain many uninformative features.



MLP-like architectures are not robust to uninformative features.

# Related works

| Name                    | Loss   |
|-------------------------|--|
| LASSO                   | $\min_w loss(w; X, y) + \alpha \ w\ _1$  |
| Group LASSO             | $\min_w loss(w; X, y) + \alpha \sum_{i=1}^g h_i \ w_{G_i}\ _2$                                 |
| Sparse Group LASSO      | $\min_w loss(w; X, y) + \alpha \ w\ _1 + (1 - \alpha) \sum_{i=1}^g h_i \ w_{G_i}\ _2$          |
| Tree-guided Group LASSO | $\min_w loss(w; X, y) + \alpha \sum_{i=0}^d \sum_{j=1}^{n_i} h_j^i \ w_{G_i}\ _2$              |
| Graph LASSO             | $\min_w loss(w; X, y) + \alpha \ w\ _1 + (1 - \alpha) \sum_{i,j} M(i,j)(w_i - w_j)^2$          |
| GFLASSO                 | $\min_w loss(w; X, y) + \alpha \ w\ _1 + (1 - \alpha) \sum_{i,j} A(i,j)(w_i - sign(i,j)w_j)^2$ |

- Can we combine the existing feature selection approaches with state-of-the-art models for heterogeneous datasets?

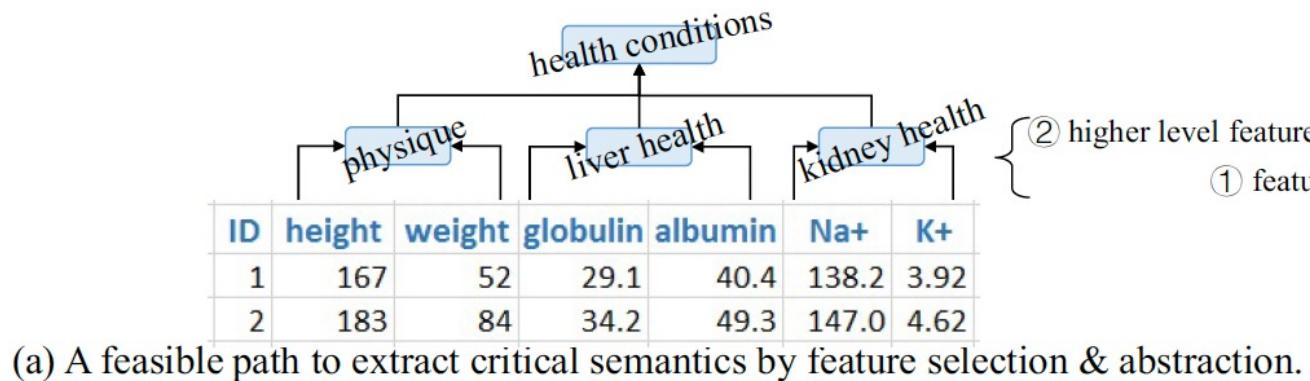
# Related works

## Heterogeneous feature selection approach

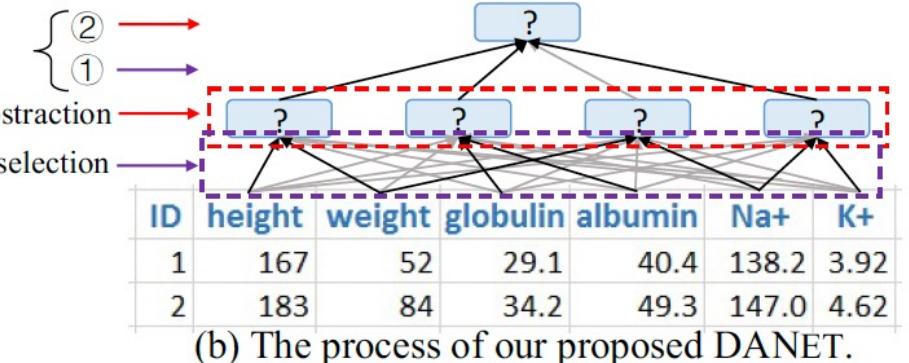
- Mask-based
  - DANETs
- Fuzzy rough set-based FS
  - Fuzzy relation: Hu et al. (2006)
  - Categorical: Wang et al. (2019)
  - Supervised: Yuan et al. (2018), Yuan et al. (2021a)
  - Unsupervised: Yuan et al. (2021b), Zhang et al. (2022)

# Related works

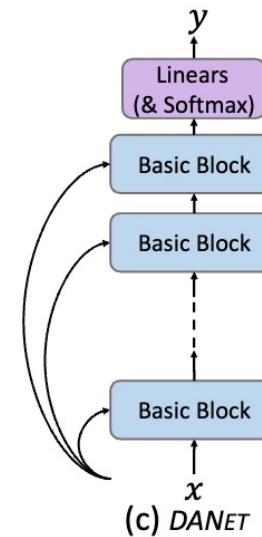
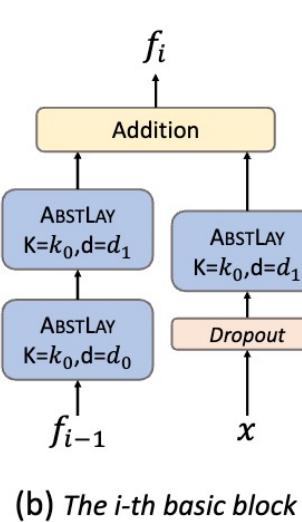
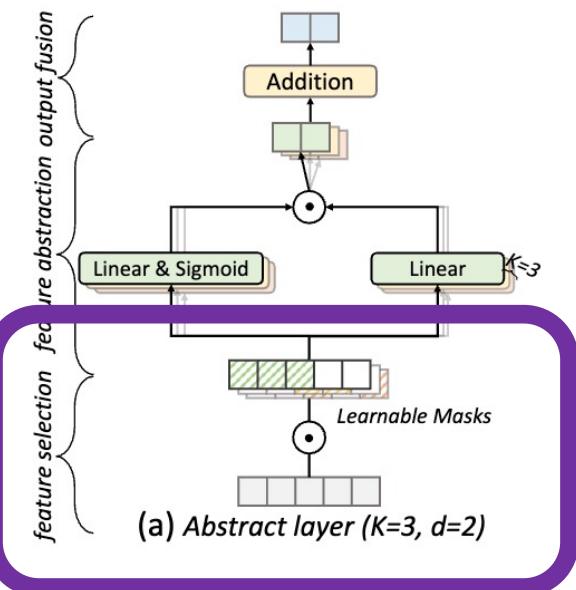
DANETs (2017)



(a) A feasible path to extract critical semantics by feature selection & abstraction.



(b) The process of our proposed DANET.



- Use multiple learnable masks to discard the uninformative features parallelly
- Feature abstraction is used to abstract high level information

# Related works

## Algorithm 1: FMIUFS algorithm.

**Input:**  $IS = \langle U, C \rangle$ , threshold value  $\lambda$ ,  $|C| = m$

**Output:** An ordered feature sequence  $S$

```

1  $S \leftarrow \emptyset$   $S_u \leftarrow C$ ;
2 for  $k \leftarrow 1$  to  $m$  do
   | Calculate the fuzzy relation matrix  $M_{\mathcal{R}_{c_k}}$ ;
   | Calculate the fuzzy entropy  $FE(c_k)$ ;
end
for  $k \leftarrow 1$  to  $m$  do
   for  $s \leftarrow 1$  to  $m$  do
      | Calculate the fuzzy joint entropy  $FE(c_k, c_s)$ ;
      | Calculate the fuzzy mutual information  $FMI(c_k; c_s)$ ;
   end
end
for  $k \leftarrow 1$  to  $m$  do
   | Calculate the fuzzy relevance  $FRel(c_k)$ ;
end
Select feature  $c_{\ell_1}$  so that  $FRel(c_{\ell_1})$  has the maximum value;
16  $S \leftarrow S \cup \{c_{\ell_1}\}$ ,  $S_u \leftarrow S_u - \{c_{\ell_1}\}$ ;
17 while  $|S_u| \neq 0$  do
   for  $l \leftarrow 1$  to  $|S_u|$  do
      for  $s \leftarrow 1$  to  $|S|$  do
         | Calculate the fuzzy redundancy  $FRed(c_l, c_{\ell_s})$ ;
      end
   end
23 Select feature  $c_{\ell_r}$  so that  $FRel(c_{\ell_r}) - \frac{1}{|S|} \sum_{s=1}^{|S|} FRed(c_{\ell_r}, c_{\ell_s})$ 
has the maximum value;
24  $S \leftarrow S \cup \{c_{\ell_r}\}$ ,  $S_u \leftarrow S_u - \{c_{\ell_r}\}$ ;
25 end
26 return  $S$ .

```

## Fuzzy relation

$$r_{ij}^k = \begin{cases} 1, & \text{if } c_k(x_i) = c_k(x_j) \text{ and } c_k \text{ is discrete} \\ 0, & \text{if } c_k(x_i) \neq c_k(x_j) \text{ and } c_k \text{ is discrete} \\ 1 - |c_k(x_i) - c_k(x_j)|, & \text{if } |c_k(x_i) - c_k(x_j)| \leq \epsilon_{c_k} \text{ and } c_k \text{ is continuous} \\ 0, & \text{if } |c_k(x_i) - c_k(x_j)| > \epsilon_{c_k} \text{ and } c_k \text{ is continuous} \end{cases} \quad (14)$$

where  $c_k$  is the measured value of data point  $x$  for feature  $c_k$  and  $\epsilon_{c_k}$  a adaptive fuzzy radius. The  $\epsilon_{c_k}$  is calculated as following,

$$\epsilon_{c_k} = \frac{std(c_k)}{\lambda} \quad (15)$$

where  $std(c_k)$  is standard deviation of the feature values  $c_k$  and  $\lambda$  is a hyper-parameter that is fine-tuned with step 0.1 in the range [0.1,2.0].

## Fuzzy entropy

$$FE(B) = -\frac{1}{|U|} \sum_{i=1}^n \log_2 \frac{|[x_i]_B|}{|U|}.$$

$$FE(B|E) = -\frac{1}{|U|} \sum_{i=1}^n \log_2 \frac{|[x_i]_B \cap [x_i]_E|}{|[x_i]_E|}.$$

$$FE(B, E) = -\frac{1}{|U|} \sum_{i=1}^n \log_2 \frac{|[x_i]_B \cap [x_i]_E|}{|U|}.$$

## Fuzzy mutual information

$$FMI(B; E) = -\frac{1}{|U|} \sum_{i=1}^n \log_2 \frac{|[x_i]_B| \times |[x_i]_E|}{|U| \times |[x_i]_B \cap [x_i]_E|}.$$

$$[x_i]_B = \bigcap_{l=1}^h [x_i]_{c_{k_l}}.$$

$$|[x_i]_B| = \sum_{j=1}^n r_{ij}^B = \sum_{j=1}^n \mathcal{R}_B(x_i, x_j).$$

## Fuzzy relevance

$$FRel(c_k) = \frac{1}{m} \sum_{s=1}^m FMI(c_k; c_s).$$

$$FRel(c_{\ell_s}|c) = \frac{FE(c_{\ell_s}|c)}{FE(c_{\ell_s})} FRel(c_{\ell_s}).$$

## Fuzzy redundancy

$$FRed(c, c_{\ell_s}) = FRel(c_{\ell_s}) - FRel(c_{\ell_s}|c)$$

The selected feature subset can minimize the uncertainty of other unselected features.

# Related works

How do they handle the heterogeneous features?

Hu 2016

$$r_{ij}^k = \begin{cases} 1, & \text{iff } f(x_i, a) = f(x_j, a) \text{ and A is discrete, } \forall a \in A \\ 0, & \text{iff } f(x_i, a) \neq f(x_j, a) \text{ and A is discrete, } \forall a \in A \\ f(\|x_i - x_j\|), & \text{if A is continuous} \end{cases}$$

Zhang 2022

$$d_{ij}^k = \begin{cases} 0, & \text{iff } f(x_i, a) = f(x_j, a) \text{ and a is discrete} \\ 1, & \text{iff } f(x_i, a) \neq f(x_j, a) \text{ and a is discrete} \\ |f(x_i, a) - f(x_j, a)|, & \text{if a is continuous} \end{cases}$$

Yuan 2018, 2021

$$r_{ij}^k = \begin{cases} 1, & \text{if } c_k(x_i) = c_k(x_j) \text{ and } c_k \text{ is discrete} \\ 0, & \text{if } c_k(x_i) \neq c_k(x_j) \text{ and } c_k \text{ is discrete} \\ 1 - |c_k(x_i) - c_k(x_j)|, & \text{if } |c_k(x_i) - c_k(x_j)| \leq \epsilon_{c_k} \text{ and } c_k \text{ is continuous} \\ 0, & \text{if } |c_k(x_i) - c_k(x_j)| > \epsilon_{c_k} \text{ and } c_k \text{ is continuous} \end{cases} \quad (14)$$

where  $c_k$  is the measured value of data point  $x$  for feature  $c_k$  and  $\epsilon_{c_k}$  a adaptive fuzzy radius. The  $\epsilon_{c_k}$  is calculated as following,

$$\epsilon_{c_k} = \frac{std(c_k)}{\lambda} \quad (15)$$

where  $std(c_k)$  is standard deviation of the feature values  $c_k$  and  $\lambda$  is a hyper-parameter that is fine-tuned with step 0.1 in the range [0.1,2.0].

Wang 2019

$$r_{ij}^B = \frac{1}{|A|} \text{card}(k \in B : c_k(x_i) = c_k(x_j))$$

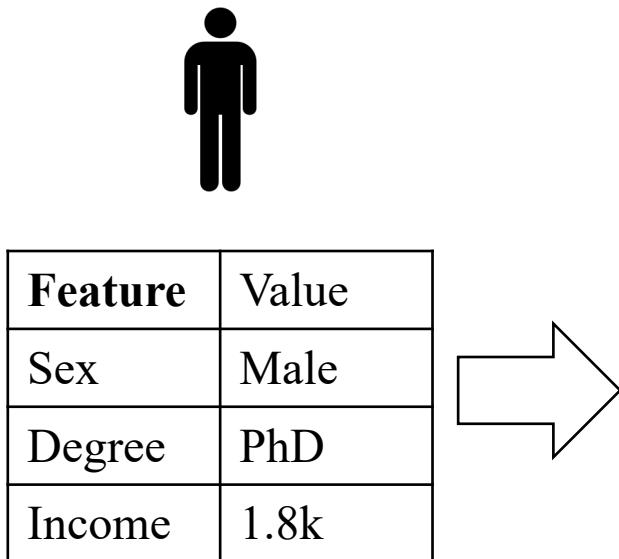
Discretization for continuous features

The goal of those works is to find a feature subset that contains most or all of the information in the original feature set based on **the entropy they defined from the relation function or distance function**.

# Methodology for Embedding

## Motivation

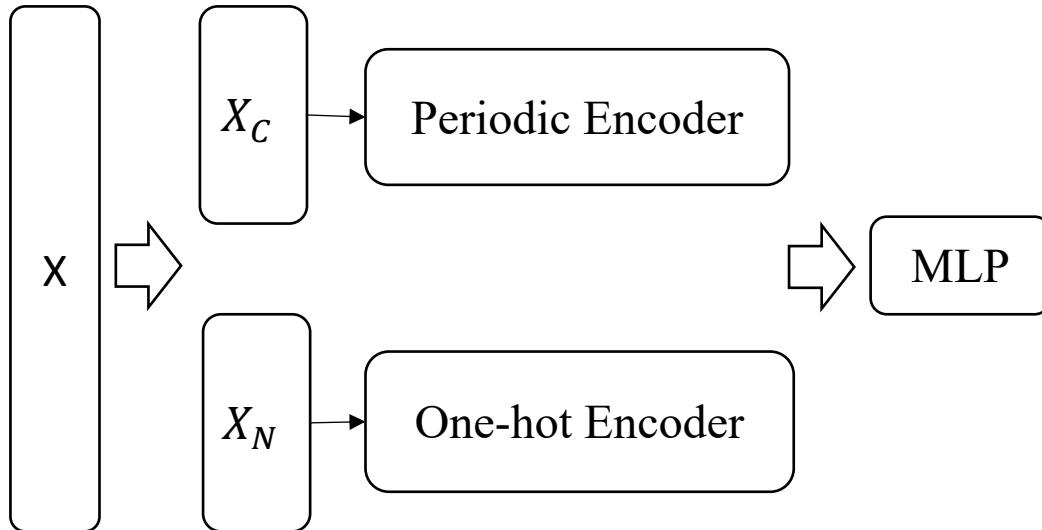
- The existing embedding module did not utilize the information on **ordinal features** and the **global frequency of assignment**



| Occurrence  | Observation | Probability | Coder  |
|-------------|-------------|-------------|--------|
| Male        | 1           | 0.5         | 1/0.5  |
| Female      | 0           | 0.5         | 0      |
| High School | 1           | 0.6         | 1/0.1  |
| Bachelor    | 1           | 0.3         | 1/0.1  |
| PhD         | 1           | 0.1         | 1/0.1  |
| 0-10K       | 1           | 0.1         | 1/0.4  |
| 10k-20k     | 0.8         | 0.4         | 1/0.32 |
| 20-30k      | 0           | 0.4         | 0      |
| >30k        | 0           | 0.1         | 0      |

# Methodology for Embedding

SOTA (2022)

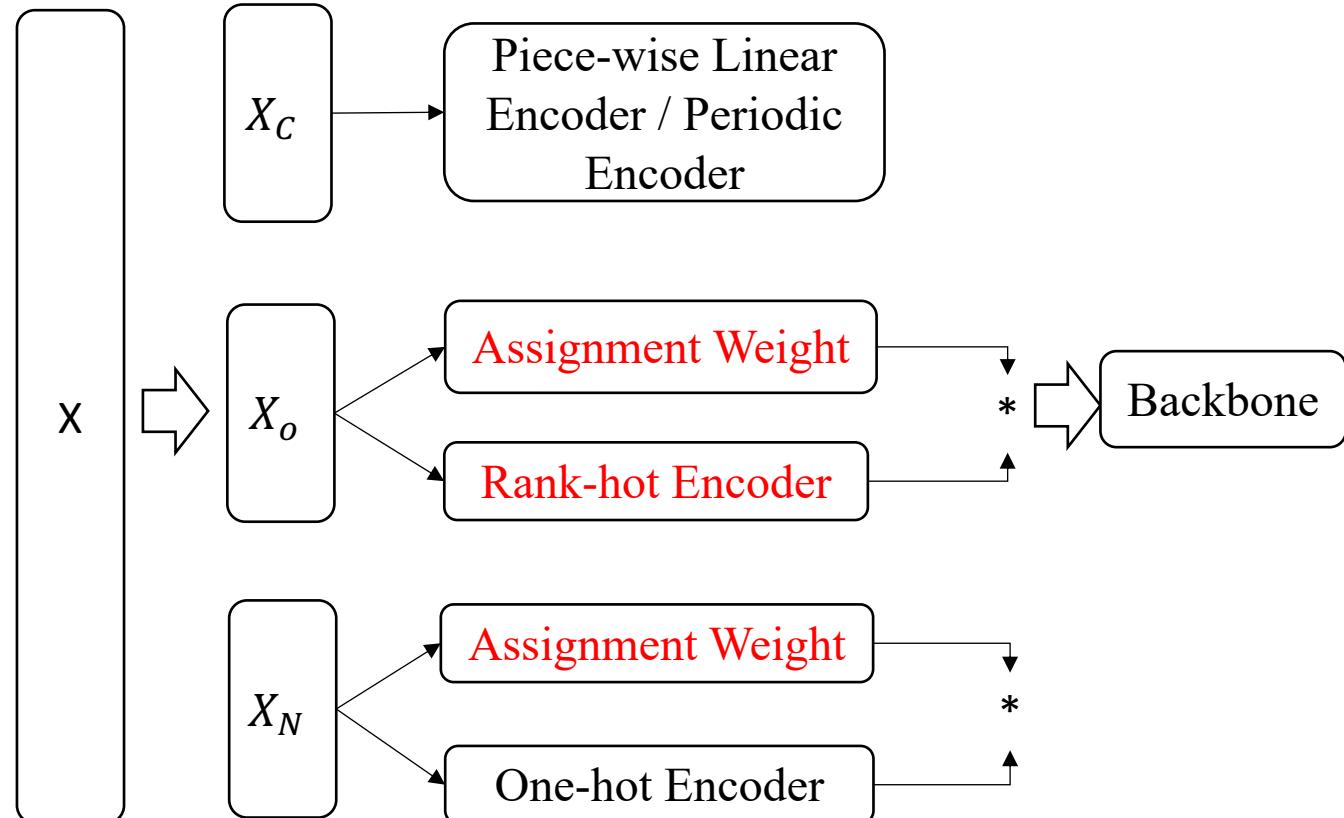


$$f_i(x) = \text{Periodic}(x) = \text{concat}[\sin(v), \cos(v)], \quad (8)$$

$$v = [2\pi c_1 x, \dots, 2\pi c_k x]$$

where  $c_i$  are trainable parameters initialized from  $\mathcal{N}(0, \sigma)$ .  $\sigma$  is an important hyperparameter that is tuned using validation sets.

Our embedding architecture



$$\text{Assignment weight } w_{f=a} = \frac{n}{n_{f=a}}$$

# Methodology for Embedding

## Alert2AKI Dataset

|                      |                                |
|----------------------|--------------------------------|
| <b>Intervention</b>  | AKI Alert or Not               |
| <b>Main-outcome</b>  | AKI Progression in 14 Days     |
| <b>Pre-treatment</b> | EHR Records                    |
| <b>Patients Num</b>  | 6030 in 5 Hospitals (5082/948) |

| SCALE    | NUM |
|----------|-----|
| Nominal  | 9   |
| Ordinal  | 19  |
| Interval | 3   |
| Ratio    | 20  |

We can compare the predicted outcome difference between different treatments for an individual to decide whether a patient should accept the treatment.

## PR-AUC (5 Random Splits)

| MLP Backbone    |              |             |
|-----------------|--------------|-------------|
| HetMLP          | HetMLP_nW    | MLP         |
| .2117±.0009     | .2087±.0164  | .2009±.0329 |
| Resnet Backbone |              |             |
| HetResNet       | HetResnet_nW | Resnet      |
| .2087±.0255     | .2033±.0145  | .1711±.0186 |

Our HetMLP got a **1.43%** performance up compared with SOTA on this dataset. <sup>21</sup>

# Will the patients benefit from the alert?

## Splitting 1

| Metrics                | Num  |
|------------------------|------|
| Patients Num           | 3536 |
| Benefited: AKI=1→AKI=0 | 15   |
| Harmful: AKI=0→AKI=1   | 14   |

## Splitting 2

| Metrics                | Num  |
|------------------------|------|
| Patients Num           | 3552 |
| Benefited: AKI=1→AKI=0 | 8    |
| Harmful: AKI=0→AKI=1   | 2    |

## Splitting 3

| Metrics                | Num  |
|------------------------|------|
| Patients Num           | 3504 |
| Benefited: AKI=1→AKI=0 | 26   |
| Harmful: AKI=0→AKI=1   | 4    |

## Splitting 4

| Metrics                | Num  |
|------------------------|------|
| Patients Num           | 3536 |
| Benefited: AKI=1→AKI=0 | 9    |
| Harmful: AKI=0→AKI=1   | 9    |

The model's prediction is consistent with the conclusion that Alerts did **not** reduce rates of our primary outcome among hospitalized patients with AKI.

# Futural Plan

- **Heterogeneous Embedding**
  - Heterogeneous feature structure and instance structure (such as cluster)
  - Computation Complexity
  - Detailed experiments on more backbone and datasets
- **Heterogeneous Feature Selection**
  - Discover more effective heterogeneous feature distance (Wasserstein etc.)
  - Combine intra-attribute structures and inter-attribute structures