

Chapter 11

Memory Networks



Memory network is a powerful extension of attention models. The memory network models have shown initial successes in natural language processing such as question answering [60, 92, 112, 141, 166]. In particular, memory networks use external memory components to assist the deep neural networks in remembering and storing information. Various memory network based models have been proposed, such as [92, 112, 141]. In healthcare, memory networks can be valuable due to their capacities in utilizing medical knowledge and patient history, as shown in [134].

11.1 Original Memory Networks

Memory Networks [166] and Differentiable Neural Computers (DNC) [60] proposed to use external memory components to assist the deep neural networks in remembering and storing information. After that, various memory network models [92, 112, 141] have been proposed. The original memory networks are proposed for solving question answering in the natural language processing domain [92, 112, 141, 166].

At a high level, a memory network consists of memory m and the following four learning components.

1. **Input feature map** I converts the incoming input to the internal feature representation.
2. **Generalization** component G updates old memories given new input. We call this generalization as there is an opportunity for the network to compress and generalize its memories at this stage for some intended future use.
3. **Output feature map** O produces a new output in the feature representation space, given the new input and the current memory state.
4. **Response** component R converts the output embeddings into the response format desired. For example, a text description or a particular class label.

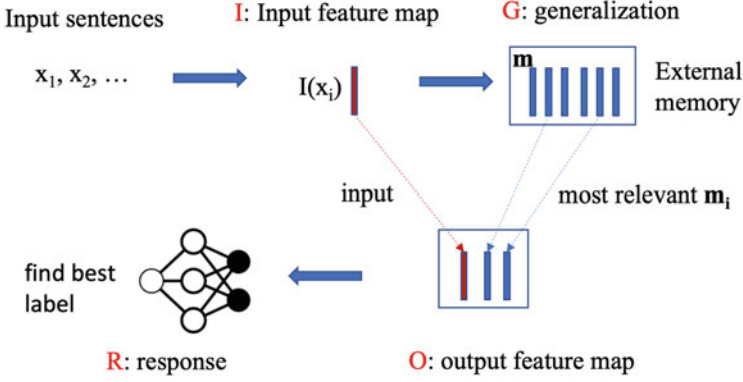


Fig. 11.1 Components of memory networks

Next, we illustrate the interaction between those components. Given an input x such as a sequence of words or sequence of clinical events, the flow of the memory network is given as follows.

1. First, we will convert x to an internal feature representation $I(x)$. Here the representation of inputs and memories can use all kinds of encodings such as a bag of words or multi-hot embeddings, RNN embedding.
2. Second, we store $I(x)$ into the memory, or update existing memory m_i with generalization and compression. Mathematically, it denotes as $m_i = G(m_i, I(x), m) \forall i$. In the simplest form, G function just stores $I(x)$ into a memory slot m_i in memory bank m .
3. Third, we compute the output feature o given the new input $I(x)$ and the current memory bank m . This component tries to find the best matching memory slots based on some matching function m : $o = O(I(x), m)$. For example, Fig. 11.1 retrieves two most relevant memory slots to combine with input $I(x)$ as the output feature map.
4. Finally we map output o to give final response $r = R(o)$. We can choose any output decoders—for example, softmax for classification or RNN for sequence generation.

We need to specify the inference mechanism in O and R modules to implement the memory network with neural networks. One simple way is to use the input as a query to score all memory slots and retrieve the ones with highest scores: $o_1 = \arg \max_i S_O(x, m_i)$ where S_O is the output scoring function that computes the similarity between input x and memory slot m_i . This can be generalized to more than one memory slot. The response module will score the input and most relevant memory slots to generate the response. For example, if the response is a single word from a set W , the response will be $r = \arg \max_{w \in W} S_R([x, o_1], w)$ where S_R is the response scoring function between a word w and the combination of input x and most relevant memory o_1 . In [166], the similarity function is

$S(x, y) = \Phi_x(x)^\top U^\top U \Phi_y(y)$ where $\Phi_x(x), \Phi_y(y)$ are the input embeddings and U is the linear embedding matrix. Finally, the model is trained in a fully supervised fashion using the ranking loss of the output module and response module:

$$\underbrace{\sum_{\tilde{f} \neq o_1} \max(0, \gamma - S_O(x, o_1) + s_O(x, \tilde{f}))}_{\text{output module}} + \underbrace{\sum_{\tilde{r} \neq r} \max(0, \gamma - S_R([x, o_1], r) + s_R([x, o_1], \tilde{r}))}_{\text{response module}}$$

where \tilde{f} and \tilde{r} are the other choices than the correct label, and γ is the margin. In the paper [166], the top two output o_1 and o_2 are generated, which leads to a slightly more complicated expression. One major limitation of the original memory network is the lack of end-to-end training due to the arg max operations which break the gradient flows which is needed for backpropagation.

11.2 End-to-End Memory Networks

The original memory network [166] uses the arg max function to find the optimal memory slots and the best output. Later on, the authors [141] extended this “hard” attention with “soft” attention via softmax function, which leads to memory networks that can be trained in an end-to-end fashion.

This end-to-end memory network involves input $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ stored in the memory bank and a query \mathbf{q} . And this network outputs an answer \mathbf{a} . For example, input \mathbf{x}_i can be all the historical patient records in the EHR database. In contrast, query \mathbf{q} is the record of the current patient, and the answer \mathbf{a} can be the medication recommended for the current patient. Two different memory representations are constructed for the input: one input memory representation called the *keys*, and one output memory representation called the *values*. Figure 11.2 illustrates the architecture of the end-to-end memory network.

Input Memory (Keys) All the input $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are embedded via matrix \mathbf{A} into a set of vectors $\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_n$. For example, a simple linear embedding can be applied $\mathbf{m}_i = \mathbf{A}\mathbf{x}_i$. The query \mathbf{q} is also embedded using a different matrix \mathbf{B} to produce the embedded vector \mathbf{u} . For example, $\mathbf{u} = \mathbf{B}\mathbf{q}$. Then the similarity matching score between query embedding \mathbf{u} and input embedding \mathbf{m}_i can be computed as inner product followed by softmax:

$$p_i = \text{Softmax}(\mathbf{u}^\top \mathbf{m}_i).$$

Here p_i is the matching score between the query \mathbf{u} and input \mathbf{m}_i . Once we calculated the matching scores from the query to all input embeddings, denoted by p_1, p_2, \dots, p_n , we can retrieve the output embeddings with highest matching scores.

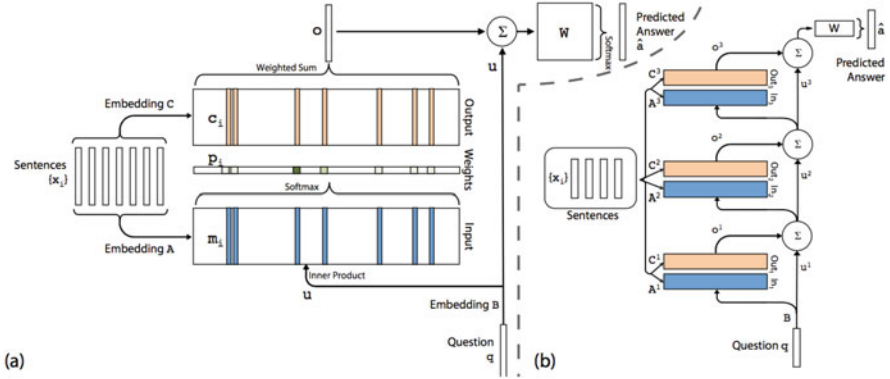


Fig. 11.2 End-to-end memory networks: (a) a single layer version, (b) multi-layer version

Output Memory (Value) Similar to input memory, the output memory are another set embedding vectors transformed from input x_1, x_2, \dots, x_n via embedding matrix C via $c_i = Cx_i$. The response vector is the weighted sum of all output embedding c_1, c_2, \dots, c_n , namely

$$o = \sum_i p_i c_i.$$

Generating Final Prediction Both the response vector o and the query vector u are used to generate the final prediction:

$$\hat{a} = \text{Softmax}(W(o + u)).$$

With this architecture, all the embedding matrix A, B, C and W are jointly trained using the cross-entropy loss between ground truth label a (a one-hot vector) and the prediction \hat{a} .

Multiple Layers Multiple memory layers can be stacked:

- The input to layer $k + 1$ can be the sum between output o^k and the input u^k from layer k , namely, $u^{k+1} = o^k + u^k$;
- Each layer will have its own embedding matrices A^k, B^k, C^k ;
- Finally the prediction per layer can be computed as $\hat{a} = \text{Softmax}(Wu^{k+1})$.

The authors proposed two ways to connect the parameters across layers.

- **Adjacent:** Output embedding matrix of layer k is the input embedding matrix for layer $k + 1$: $A^{k+1} = C^k$.
- **Layer-wise:** Tying input and output embeddings across all layers, i.e., $A^1 = A^2 = \dots = A^k$ and $C^1 = C^2 = \dots = C^k$. Then a linear mapping H is added between input layers: $u^{k+1} = Hu^k + o^k$.

11.3 Self-Attention and Transformer

Self-attention is a more efficient strategy for computing attention scores among the original input without any dependency on RNN models. Standard attention models heavily rely on RNN as the annotation vectors \mathbf{h}_i are the results of the encoder RNN and the hidden states \mathbf{s}_i are also outputs of another RNN. Because of the sequential nature of RNN, the model training process leads to significant computational and optimization challenges. More specifically, attention models need all the hidden states \mathbf{h}_1 to \mathbf{h}_T of the entire input sequence to compute the context vector \mathbf{c}_i . Such a training process is intrinsically sequential within an input example, limiting its ability to optimize training in parallel, which is computationally more efficient.

To address this parallel training challenge, *Transformer* model is proposed, which introduces a new type of attention without RNN [160]. The insight comes from a simple question “why don’t we compute context vectors over \mathbf{x}_i directly”. We actually can, and the way to achieve this is through the *self-attention* mechanism. Conceptually, the self-attention model replaces the RNNs from the encoder and decoder completely with an attention mechanism with the input to itself.

The goal of transformer model is to “transform” an input embedding \mathbf{x}_i into a better embedding \mathbf{x}'_i that captures the contexts (or correlation) of \mathbf{x}_i with all the other input. To achieve that, a query, keys, and values are introduced similar to an end-to-end memory network [142]. Here a query is a specific input \mathbf{x}_i , while keys and values are all input $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. There are two steps involved:

- **Similarity search:** The idea is to compute similarity between query \mathbf{x}_i and all input represented by $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d_x \times n}$. Instead of using original input \mathbf{X} , a linear embedding layer is introduced: $\mathbf{K} = \mathbf{W}^K \mathbf{X} \in \mathbb{R}^{d_K \times n}$ where $\mathbf{W}^K \in \mathbb{R}^{d_K \times d_x}$. Similarly, we linearly transform the query $\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i \in \mathbb{R}^{d_K}$. A scaled dot product is computed as the similarity score vector $\mathbf{e}_i = \frac{\mathbf{K}^\top \mathbf{q}_i}{\sqrt{d_K}}$. The scaled dot product is very similar to a dot product $\mathbf{K}^\top \mathbf{q}_i$ but with additional $1/\sqrt{d_K}$ factor. The idea of $1/\sqrt{d_K}$ factor is to alleviate the concern that the softmax function of a high dimensional input may have an extremely small gradient. Like the standard attention mechanism, the attention weight is normalized by $\mathbf{a}_i = \text{softmax}(\mathbf{e}_i) \in \mathbb{R}^n$.
- **Value retrieval:** Based on the similarity score \mathbf{a}_i , we can retrieve or recombine the input to generate the output embedding $\mathbf{x}'_i = \mathbf{W}^V \mathbf{X} \mathbf{a}_i$.

Note that all weight matrices $\mathbf{W}^K, \mathbf{W}^V, \mathbf{W}^Q \in \mathbb{R}^{d_K \times d_x}$ are of the same size. The reason to have separate weight matrices is to avoid dependency and maximize potential parallelism in training. The self-attention mechanism eliminates any sequential dependency; hence the training can be done in parallel more easily.

Besides the self-attention mechanism, the Transformer model introduces a few other innovations, namely, **multi-head attention** and **positional embedding**.

- The idea of multi-head attention is to compute multiple \mathbf{a}_i for the same input \mathbf{x}_i with different sets of weight matrices, which is analogous to multiple filters in

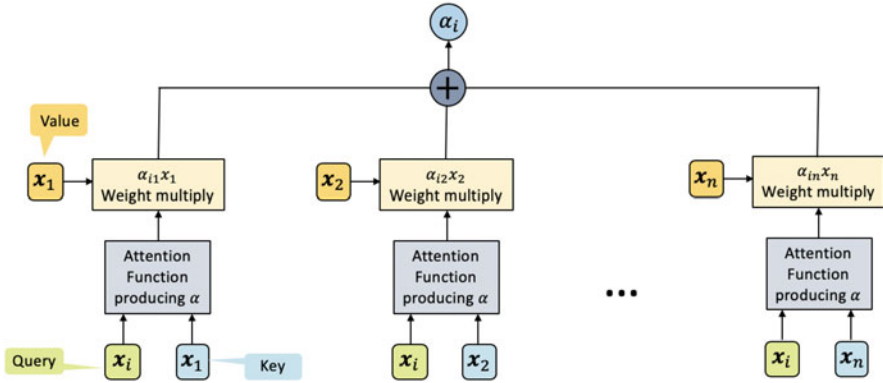


Fig. 11.3 Concept illustration of self-attention mechanism

CNN models. The benefit having multiple attention scores instead of one is to provide more robust estimate of the embedding.

- The positional embedding adds some features related to the position of x_i in the input sequence. This is important for natural language process application as a sequence corresponds to a sentence, and x_i is a word in the sequence. Therefore, where the word x_i appears in the sentence (at the beginning or the end) can be important (Fig. 11.3).

Now we understand the intuition and key components of the Transformer model. Next, we describe the overall algorithm, which follows an encoder and decoder architecture like a sequence-to-sequence model. Figure 11.4 illustrates the Transformer model.

Transformer Encoder The encoder can have multiple multi-head self-attention modules. The input embeddings (with optional positional embedding) are the input for the first layer. For subsequent layers, the output of the previous layer will be its input. The same input set will be used within each layer as keys, queries, and values to compute self-attention. The self-attention output will send through a feed-forward network layer with a skipped connection.

Transformer Decoder The decoder can also have multiple layers. Each layer is similar to the encoder but with different inputs. The input will be the output embedding of the decoder so far. Specifically, the input of the first layer will be the word embeddings generated so far. And subsequent layers will take the output of the previous layer as its input. Within each layer, a **masked multi-head self-attention** is computed with an appropriate mask so the attention will not be computed over the unseen words (i.e., those words to the right of the current word). For example, if the decoder output generates 5 tokens so far while the entire sequence has 20 tokens, the masked attention will not compute attention with the token 6–20 as they are not yet seen. The output is sent to another multi-head attention layer with the input embedding produced by the encoder. Finally, a feed-forward network layer

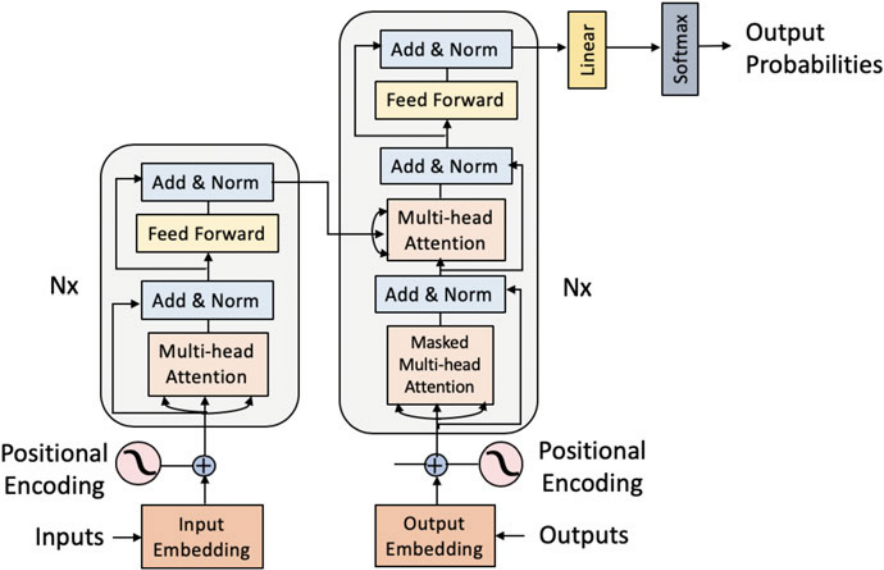


Fig. 11.4 The Transformer model architecture: The left side is the encoder and the right side is the decoder

like in encoder is applied at the end. The prediction task at position i is to predict next token at position $i + 1$ for all positions.

11.4 BERT: Pre-training of Deep Bidirectional Transformers

One important application of the Transformer is the BERT model for supporting natural language processing (NLP) tasks such as language understanding and question answering [37]. BERT stands for Bidirectional Encoder Representations from Transformers. The key ideas behind BERT are (1) context-specific pre-training and (2) masked language model.

Context Specific Pre-training Embedding vectors are the basis for deep learning models. In many tasks, people find it is useful to pre-train embedding from large datasets. For example, Word2Vec [110] and GloVe [117] are often used to generate pre-trained embeddings on large text corpus based on co-occurrence information between words. However, the same word will have the same embedding no matter what the contexts are. For example, although “bank account” and “river bank” have completely different meanings, the word “bank” will have the same embedding. BERT changed this assumption by introducing the context-specific embeddings, where the embedding of word x is dynamically computed based on the context of the sentence where x is in. BERT can be trained in an unsupervised fashion on a

large corpus without any specific target label. If the labeled data exist, BERT also allows fine-tuning using the labeled data.

Masked Language Model Context-specific language models such as RNN are trained either from left to right or right to left. Bi-directional models are achieved by training 2 models from both directions (left to right and right to left) and concatenating the hidden states shown in bi-directional RNN. BERT introduces a new mechanism to achieve a bi-directional language model by randomly masking words. More specifically, by masking $k\%$ of input words (e.g., $k = 15$), the BERT model tries to predict the mask words with a modified Transformer decoder. The details of the masking procedure are more involved:

- 80% of the time we replace the masked word with a special token named [MASK];
- 10% of the time we replace the word with a random word from the entire vocabulary;
- 10% of the time we keep the original word unchanged. The purpose is to guide the BERT model towards this actual word.

The BERT model provides a power embedding method based on pre-training and masked language model, which have shown great performance on many real-world applications.

11.5 Case Study: Doctor2Vec—Doctor Recommendation for Clinical Trial Recruitment

Problem Massive electronic health records (EHRs) enable the success of learning accurate patient representations to support various predictive health applications. In contrast, doctor representation was not well studied, although doctors play pivotal roles in healthcare. How to construct the right doctor representations? How to use doctor representation to solve important health analytic problems? The authors of [8] study the doctor representation problem for the *clinical trial recruitment* application, which is about identifying the right doctors to help conduct the trials based on the trial description and patient EHR data of those doctors.

Method They propose Doctor2Vec, which simultaneously learns (1) doctor representations from EHR data and (2) trial representations from the description and categorical information about the trials. In particular, Doctor2Vec utilizes a *dynamic memory network* where the doctor's experience with patients is stored in the memory bank. The network will dynamically assign weights based on the trial representation via an attention mechanism. Figure 11.5 illustrates the overall framework of Doctor2Vec. Since each doctor sees a diverse set of patients, doctor representation should be dynamically constructed for a given trial instead of a static embedding vector staying the same for all trials. The authors achieved that by

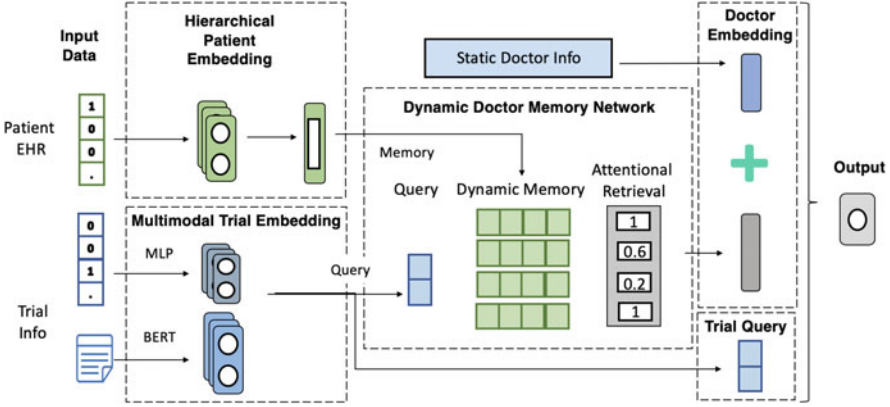


Fig. 11.5 Doctor2Vec Framework. (1) Hierarchical patient embedding: We obtain patient embeddings from patient visits using a Bi-LSTM with attention module. Unstructured text and categorical data are modeled using MLP and BERT, respectively. (2) Multimodal clinical trial information embedding: The obtained clinical trial embeddings from text and categorical data are concatenated together to form the query vector for the memory network. (3) Memory network module: This query vector is used to attend over the memory bank consisting of patient embeddings. The attention vector is used to obtain the final doctor embedding. The doctor embedding \mathbf{Doc}_{emb} is combined with clinical trial embedding $\mathbf{Q}_{emb}(\mathbf{I})$ and static information about the doctors \mathbf{Doc}_{static} to predict the enrollment rate of the clinical trial (output)

a dynamic memory network where patients are stored as memory vectors of the doctor. Using a trial embedding as a query, they fetch the relevant patient vectors from the memory bank and dynamically assemble a doctor representation for this trial.

Inspired by Weston et al. [166], four memory components **I**, **G**, **O**, **R** are proposed which mimics the architecture of modern computer in storing and processing information.

1. **Input Memory Representation.** This layer converts the patient representations to the input representation. They pass all the patient representations through a dense layer to obtain the input representations.

$$\mathbf{I}_f(\mathbf{k}) = \mathbf{W}_i \mathbf{I}(\mathbf{k}) + \mathbf{b}_i$$

where $\mathbf{I}(\mathbf{k})$ and $\mathbf{I}_f(\mathbf{k})$ are input and final embedding of patient k .

2. **Generalization.** Typically generalization can be referred to as the process of updating memory representation for the memory bank. In our case, they use the patient representations to initialize the memory representation \mathbf{M}_d , which is the combination of all the patient representations. They then apply an LSTM layer to update the memory via multiple iterations.

$$\mathbf{M}_d = \text{LSTM}(\mathbf{I}_f(1), \dots, \mathbf{I}_f(k)) \quad (11.1)$$

3. **Output.** In this step, the final output memory representation is generated. They calculate the relevance between trial embedding $\mathbf{Q}_{\text{emb}}(\mathbf{l})$ and doctor embedding \mathbf{M}_d to obtain $\mathbf{A}(\mathbf{k})$ as the attention vector over patient representations.

$$\mathbf{A}(\mathbf{k}) = \text{softmax}[\mathbf{Q}_{\text{emb}}(\mathbf{l})^T \mathbf{M}_d] \quad (11.2)$$

4. **Response.** In this step, they obtain the final $\mathbf{Doc}_{\text{emb}}$ using the patient embeddings and attention weights over the patients.

$$\mathbf{Doc}_{\text{emb}} = \sum \mathbf{A}(\mathbf{k}) \mathbf{I}_f(\mathbf{k}) \quad (11.3)$$

They use the doctor representation, composed of patients, and the clinical trial representation to obtain a final context vector. Besides dynamic doctor embedding $\mathbf{Doc}_{\text{emb}}$, they also include static information about doctors in the final embedding, such as their educational history, length of practice, length of practice into the feature vector. The resulting final embedding vectors are then fed into a fully connected layer and passed through a softmax to obtain class labels.

$$\mathbf{Y} = \text{Softmax}([\mathbf{Doc}_{\text{emb}}; \mathbf{Q}_{\text{emb}}(\mathbf{l}); \mathbf{Doc}_{\text{static}}])$$

where the input to Softmax are concatenation of dynamic doctor embedding $\mathbf{Doc}_{\text{emb}}$, trial query embedding $\mathbf{Q}_{\text{emb}}(\mathbf{l})$ and static doctor embedding $\mathbf{Doc}_{\text{static}}$.

The final prediction target is the trial enrollment rate for each doctor. The enrollment rate of a doctor is the number of patients enrolled by a doctor to the trial. And the enrollment rate is normalized by min and max values within each trial. Then the enrollment rate category is obtained by binning the continuous enrollment rate. They divide the continuous enrollment scores into five discrete classes ranging at 0–0.2, 0.2–0.4, 0.4–0.6, 0.6–0.8, 0.8–1.0. The five enrollment categories are used labels for classification.

Result Doctor2Vec is validated on large real-world trials and EHR data, including 2609 trials, 25K doctors and 430K patients.

They consider the following baselines.

1. Median Enrollment (Median). The current industry standard considers each therapeutic area's median enrollment rate as the estimated rate for all trials in that area.
2. Logistic Regression (LR). They combine the medication, diagnosis, procedure codes, and clinical trial information to create feature vectors and then apply LR to predict the enrollment rate category.
3. Random Forest (RF). They combine the medication, diagnosis, and procedure codes and clinical trial information to create feature vectors and then pass it to RF to predict the enrollment rate category.

Table 11.1 Doctor2Vec achieves the best performance on both metrics in predicting actual enrollment rate (regression task) and rate categories (classification task) compared to state-of-the-art baselines. Results of ten independent runs

	PR-AUC	R^2 Score
Median	0.571 ± 0.014	0.54 ± 0.072
LR	0.672 ± 0.041	0.314 ± 0.082
RF	0.731 ± 0.034	0.618 ± 0.034
AdaBoost	0.747 ± 0.002	0.684 ± 0.146
MLP	0.761 ± 0.019	0.762 ± 0.049
LSTM	0.792 ± 0.034	0.780 ± 0.621
DeepMatch	0.735 ± 0.068	0.821 ± 0.073
Doctor2Vec	0.861 ± 0.021	0.841 ± 0.072

The bold values indicate the best performance numbers across all methods

4. AdaBoost. They combine the medication, diagnosis, procedure codes, and clinical trial information to create feature vectors and then apply the AdaBoost classifier to predict the enrollment rate categories.
5. Multi-layer Perceptron (MLP). They use MLP to process doctor features. In this case, they obtain the doctor features by converting all the visit vectors associated with a doctor to a count vector of different diagnoses, medication, procedure codes. They convert categorical information of clinical trials to multi-hot vectors and obtain TF-IDF features from text information of clinical trials.
6. Long Short-Term Memory Networks (LSTM). They process all the temporally ordered visit vectors associated with a doctor using an LSTM. The embedding obtained from LSTM is concatenated with embedding obtained from categorical and text information of clinical trials to predict enrollment rate.
7. DeepMatch [56] In this model, the doctors’ features are obtained by collecting the top 50 most frequent medical codes and passed through an MLP layer to obtain an embedding vector. This embedding is concatenated with embedding obtained from categorical and text information of clinical trials via MLP and TF-IDF to predict enrollment rate finally.

They conducted experiments for both classification (e.g., predict enrollment rate category) and regression (e.g., predict actual rate) tasks. Results are provided in Table 11.1. From the results, they observe that Doctor2Vec achieved the best performance in both settings.

In particular, Doctor2Vec demonstrated improved performance over the best baseline LSTM by up to 8.7% in PR-AUC. In the actual rate prediction task, Doctor2Vec gains 2.4% relative improvement in R^2 over the best baseline DeepMatch.

11.6 Case Study: Medication Recommendation

Problem How to select effective and safe medications? Can we perform medication recommendations based on medical history from EHR data and relevant

Table 11.2 Statistics of the data

# Patients	6350
# Clinical events	15,016
# Diagnosis	1958
# Procedure	1426
# Medication	145
Avg # of visits	2.36
Avg # of diagnosis	10.51
Avg # of procedure	3.84
Avg # of medication	8.80
Max # of diagnosis	128
Max # of procedure	50
Max # of medication	55
# Medication in DDI knowledge base	123
# DDI types in knowledge base	40

medical knowledge base such as drug-drug interactions (DDI) database? A good medication recommendation needs to perform recommendation based on longitudinal patient history and considers drug safety in their modeling, especially adverse drug-drug interactions (DDI). Authors [134] proposed a memory network based approach called GAMENet for recommending medication.

Data The experiment is conducted using EHR data from MIMIC-III [81]. Here patients with more than one visit are included in the dataset. They used DDI knowledge from TWOSIDES dataset [152]. In this work, the Top-40 severity DDI types are transformed to ATC Third Level for integrating with MIMIC-III data. The dataset statistics are summarized in Table 11.2.

Method GAMENet model is a graph augmented memory network that embeds multiple knowledge graphs into the memory bank. GAMENet also enables attention-based memory search using query generated from longitudinal patient records. In particular, GAMENet not only stores the EHR graph and the drug-drug interaction graph as facts in Memory Bank (MB), but also inserts patient history to the Dynamic Memory (DM). Next, we present a high-level description of different modules in GAMENet, which is illustrated in Fig. 11.6.

- *Input medical embedding*: A visit \mathbf{x}_t consists of $[\mathbf{c}_d^t, \mathbf{c}_p^t, \mathbf{c}_m^t]$ where each \mathbf{c}_*^t is a multi-hot vector at the t th visit. The multi-hot vector \mathbf{c}_*^t is binary encoded showing the existence of each medical codes recorded at the t th visit. A linear embedding is applied on the input vectors. We derive medical embeddings for $\mathbf{c}_d^t, \mathbf{c}_p^t$ for diagnosis codes and procedure codes separately at the t th visit as follows:

$$\mathbf{e}_*^t = \mathbf{W}_{*,e} \mathbf{c}_*^t$$

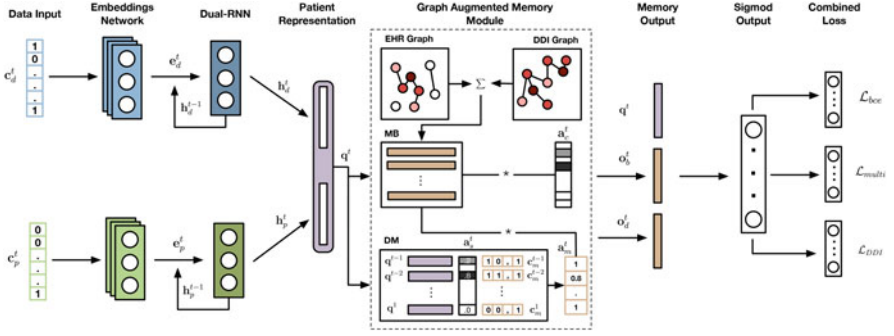


Fig. 11.6 GAMENet: At the t th visit, there are the multi-hot input vectors $\mathbf{c}_d^t, \mathbf{c}_p^t$ for diagnosis and procedure, respectively, which are mapped to embedding $\mathbf{e}_d^t, \mathbf{e}_p^t$ correspondingly. Then two separate RNN generates current hidden states on diagnosis \mathbf{h}_d^t and procedure \mathbf{h}_p^t based on embedding $\mathbf{e}_d^t, \mathbf{e}_p^t$. We use concatenated $\mathbf{h}_d^t, \mathbf{h}_p^t$ as query \mathbf{q}^t (a.k.a. patient representation) to output \mathbf{o}_b^t by reading from Memory Bank (MB) \mathbf{M}_b . Meantime, the Dynamic Memory (DM) stores key-value form history information along time and can be used to generate output \mathbf{o}_d^t . Finally, query and memory outputs are concatenated to make final medication recommendation. A combined loss is used to balance prediction performance and DDI

where $\mathbf{W}_{*,e} \in \mathbb{R}^{|\mathcal{C}_*| \times d}$ is the embedding matrix to learn and $*$ can be either d for diagnosis or p for procedure. Thus a visit \mathbf{x}_t is transformed to $\hat{\mathbf{x}}_t = [\mathbf{e}_d^t, \mathbf{e}_p^t, \mathbf{c}_m^t]$ where $\mathbf{e}_d^t, \mathbf{e}_p^t$ are diagnosis and procedure embeddings at time t , and \mathbf{c}_m^t is the multi-hot vectors for medication at time t .

- *Patient representation:* The two RNNs are used for modeling diagnosis sequences and procedure sequences of patients, separately. Thus, the RNNs accept all patient history to produce hidden states to generate patient representation as a query \mathbf{q}^t . Formally, for each input vector in transformed clinical history $[\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_t]$, we retrieve $\mathbf{e}_d, \mathbf{e}_p$ and utilize RNN to encode visit-level diagnosis and procedure embeddings respectively as follows:

$$\begin{aligned} \mathbf{h}_d^t &= \text{RNN}_d(\mathbf{e}_d^1, \mathbf{e}_d^2, \dots, \mathbf{e}_d^t) \\ \mathbf{h}_p^t &= \text{RNN}_p(\mathbf{e}_p^1, \mathbf{e}_p^2, \dots, \mathbf{e}_p^t) \end{aligned}$$

Finally, we can convert hidden states into a query vector as $\mathbf{q}^t = f([\mathbf{h}_d^t, \mathbf{h}_p^t])$ where we concatenate hidden diagnosis state \mathbf{h}_d^t and procedure state \mathbf{h}_p^t as the input patient health state. $f(\cdot)$ is a fully connected neural network with one hidden layer.

- *Static memory bank:* To leverage drug knowledge, we construct a graph augmented memory module that not only embeds and stores the EHR graph and the DDI graph as facts in Memory Bank (MB). Here graph convolution neural

networks (GCN)¹ are used to construct the knowledge embedding vectors for medications to be stored in the memory bank. More formally, EHR graph and DDI graph are represented by two adjacency matrices A_e, A_d , respectively. A_e incorporates medication co-prescriptions (i.e., edges between drugs that are often co-prescribed together) while A_d encodes drug-drug interactions (i.e., two drugs with DDI are connected). Following the GCN procedure [88], each A_* is pre-processed as follows:

$$\tilde{A}_* = \tilde{D}^{-\frac{1}{2}}(A_* + I)\tilde{D}^{-\frac{1}{2}}$$

where \tilde{D} is a diagonal matrix such that $\tilde{D}_{ii} = \sum_j A_{ij}$ and I is an identity matrix. Then we applied a two-layer GCN on A_e and A_d , respectively. The output M_b is generated as a weighted sum of the two graph embeddings.

$$\begin{aligned} Z_1 &= \tilde{A}_e \tanh(\tilde{A}_e W_{e1}) W_1 \\ Z_2 &= \tilde{A}_d \tanh(\tilde{A}_d W_{e2}) W_2 \\ M_b &= Z_1 - \beta Z_2 \end{aligned} \quad (11.4)$$

where $W_{e1}, W_{e2} \in \mathbb{R}^{|\mathcal{C}_m| \times d}$ are embedding matrices for EHR graph and DDI graph (each contains $|\mathcal{C}_m|$ number of d -dimensional vectors), $W_1, W_2 \in \mathbb{R}^{d \times d}$ are parameter matrices. All W_* are updated during the training phase. Then, medication/node embeddings $Z_1, Z_2 \in \mathbb{R}^{|\mathcal{C}_m| \times d}$ are computed using GCN. Finally we combine node embeddings Z_1 and Z_2 into $M_b \in \mathbb{R}^{|\mathcal{C}_m| \times d}$ where β is a weighting variable to fuse both graphs. The combined embedding M_b for each medication will be stored into the memory bank.

- *Dynamic memory bank* stores patient history to Dynamic Memory (DM) as key-value pairs. The keys are the previous patient representation query $q^{t'}, t' < t$, while the values are the corresponding multi-hot medication vectors $c_m^{t'}, t' < t$ that are used in the past. Specifically, we can incrementally insert key-value pair after each visit step and treat M_d^t as a vectorized indexable dictionary as follows $M_d^t = \{q^{t'} : c_m^{t'}\}_{1}^{t-1}$ where M_d^t is empty when $t = 1$. For clarity, we use $M_{d,k}^t = [q^1, q^2, \dots, q^{t-1}] \in \mathbb{R}^{d \times |t-1|}$ to denote the key vectors and $M_{d,v}^t = [c_m^1; c_m^2; \dots; c_m^{t-1}] \in \mathbb{R}^{|t-1| \times |\mathcal{C}_m|}$ to denote the value vectors at t th visit.
- *Medication recommendation*: We first query two memory banks MB and DM with current patient representation q^t to retrieve output o_b^t, o_d^t as follows:

$$o_b^t = M_b^T \overbrace{\text{Softmax}(M_b q^t)}^{a_c^t}$$

¹GCN will be covered in Chap. 10.

$$\mathbf{o}_d^t = \mathbf{M}_b^\top (\mathbf{M}_{d,v}^t)^\top \underbrace{\text{Softmax}((\mathbf{M}_{d,k}^t)^\top \mathbf{q}^t)}_{\mathbf{a}_s^t}$$

Then we concatenate query \mathbf{q}^t and their output from memory banks and perform the final classification

$$\hat{\mathbf{y}}_t = \sigma([\mathbf{q}^t, \mathbf{o}_b^t, \mathbf{o}_d^t])$$

Results

The performance of GAMENet is compared with the following baselines.

- **Nearest** will simply recommend the same combination medications at previous visit for current visit (i.e., $\hat{Y}_t = Y_{t-1}$)
- **Logistic Regression (LR)** is a logistic regression with L2 regularization. Here we represent the input data by sum of one-hot vector.
- **LEAP** [174] is an instance-based medication combination recommendation method.
- **RETAIN** [25] can provide a sequential prediction of medication combination based on a two-level neural attention model that detects influential past visits and significant clinical variables within those visits.
- **DMNC** [95] is a recent work of medication combination prediction via memory augmented neural network.

Table 11.3 compares the performance on accuracy and safety issues. Results show GAMENet has the highest score among all baselines for Jaccard, PR-AUC, and F1.

Table 11.3 Performance comparison of GAMENet and baseline methods

Methods	DDI rate	Δ DDI rate %	Jaccard	PR-AUC	F1	Avg # of med.
Nearest	0.0791	+1.80%	0.3911	0.3805	0.5465	14.77
LR	0.0786	+1.16%	0.4075	0.6716	0.5658	11.42
Leap	0.0532	−31.53%	0.3844	0.5501	0.5410	14.42
RETAIN	0.0797	+2.57%	0.4168	0.6620	0.5781	16.68
DMNC	0.0949	+22.14%	0.4343	0.6856	0.5934	20.00
GAMENet	0.0749	−3.60%	0.4509	0.6904	0.6081	14.02

The bold values indicate the best performance numbers across all methods

11.7 Case Study: Pre-training of Graph Augmented Transformers for Medication Recommendation

Problem G-BERT is another medication recommendation algorithm that extends the GAMENet model. Recall GAMENet utilizes memory networks for medication recommendation [134]. However, there are two limitations with GAMENet, which are addressed by G-BERT, a graph augmented Transformer model [133].

1. **Selection bias:** Data from patients who only have one hospital visit were discarded from training by GAMENet since it requires more than one visit to build a temporal prediction model like GAMENet. As a result, the training population is biased towards more severe patients with multiple inpatient visits.
2. **Lack of hierarchical knowledge:** Medical ontologies such as a diagnosis ontology follow hierarchical structures that were not utilized in GAMENet.

Method G-BERT is a model that first derives the initial embedding of medical codes from medical ontology using graph neural networks (**Ontology embedding**). Then, G-BERT constructs a variation of the BERT model on data from all patients of single or multiple visits (**Visit embedding**). Finally, we add a prediction layer to fine-tune the medication recommendation model (**Fine-tuning**). Figure 11.7 illustrates the model architecture of G-BERT.

Ontology Embedding

We constructed ontology embedding from diagnosis ontology \mathcal{O}_d and medication ontology \mathcal{O}_m . Since the medical codes in raw EHR data can be considered leaf nodes in these ontology trees, we can enhance the medical code embedding using graph neural networks (GNNs) to integrate these codes' ancestors' information.

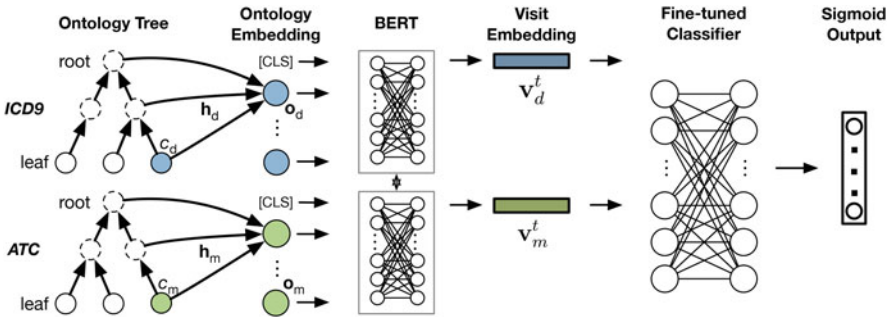


Fig. 11.7 G-BERT architecture: It consists of three parts: ontology embedding, visit embedding and classification fine-tuning. First, we derive ontology embedding for medical codes from medical ontology using graph attention networks. Second, we input diagnosis and medication ontology embeddings separately to the BERT models. Third, we concatenate the mean of all previous visit embeddings and the current visit embedding as the final feature vector to train a classifier for medication recommendation

Here we perform a two-stage procedure with a specially designed GNN for ontology embedding.

To start, we assign an initial embedding vector to every medical code $c_* \in \mathcal{O}_*$ with a learnable embedding matrix $\mathbf{W}_e \in \mathbb{R}^{|\mathcal{O}_*| \times d}$ where d is the embedding dimension. Here $\overline{\mathcal{C}}_*$ and \mathcal{C}_* correspond to the set of non-leaf nodes and that of the leaf nodes, respectively.

Bottom-Up Stage For each non-leaf node $c_* \in \overline{\mathcal{C}}_*$, we obtain its enhanced medical embedding $\mathbf{h}_{c_*} \in \mathbb{R}^d$ as follows:

$$\mathbf{h}_{c_*} = g(c_*, ch(c_*), \mathbf{W}_e)$$

where $g(\cdot, \cdot, \cdot)$ is an aggregation function which accepts the target medical code c_* , its direct children $ch(c_*)$ and initial embedding matrix. Intuitively, the aggregation function can fuse information into the target node c_* from its direct children.

Top-Down Stage For the leaf nodes, we use the embedding matrix of non-leaf nodes $\mathbf{H}_e \in \mathbb{R}^{|\overline{\mathcal{O}}_*| \times d}$ in a top-down fashion to obtain their embeddings $c_* \in \mathcal{C}_*$.

$$\mathbf{o}_{c_*} = g(c_*, pa(c_*), \mathbf{H}_e)$$

where $g(\cdot, \cdot, \cdot)$ accepts the target medical code c_* and its parents $pa(c_*)$.

The option for the aggregation function $g(\cdot, \cdot, \cdot)$ is flexible, including *sum*, *mean*. Here we choose the one from graph attention networks (GAT) [161]. For example, for the top-down stage, the aggregation function is

$$g(c_*, p(c_*), \mathbf{H}_e) = \left\|_{k=1}^K \sigma \left(\sum_{j \in \mathcal{N}_{c_*}} \alpha_{c_*, j}^k \mathbf{W}^k \mathbf{h}_j \right) \right\|$$

where $\|$ represents concatenation of the K -dimensional output from the multi-head attention mechanism, σ is a nonlinear activation function, $\mathbf{W}^k \in \mathbb{R}^{m \times d}$ is the k -th weight matrix for input transformation, and $\alpha_{c_*, j}^k$ are the corresponding k -th normalized attention coefficients computed as follows:

$$\alpha_{c_*, j}^k = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}^k \mathbf{h}_{c_*} \parallel \mathbf{W}^k \mathbf{h}_j]))}{\sum_{k \in \mathcal{N}_{c_*}} \exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}^k \mathbf{h}_{c_*} \parallel \mathbf{W}^k \mathbf{h}_k]))} \quad (11.5)$$

where \mathcal{N}_{c_*} is the neighbors of c_* which is $ch(c_*) \cup c_*$ in bottom-up stage and $pa(c_*) \cup c_*$ in top-down stage, respectively. $\mathbf{a} \in \mathbb{R}^{2d/K}$ is a learnable weight vector and LeakyReLU is a nonlinear activation function introduced in Sect. 4.1.1.

As shown in Fig. 11.7, ontology hierarchies for diagnosis and medications are constructed based on ICD and Anatomical Therapeutic Chemical (ATC) Classification, respectively.

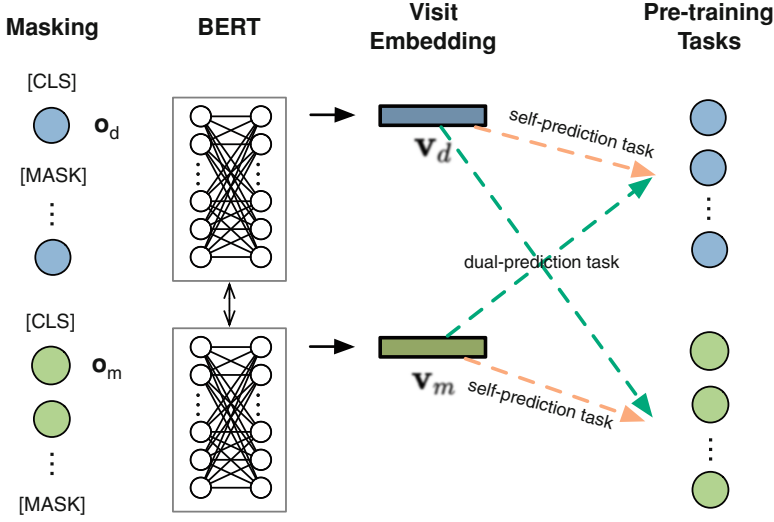


Fig. 11.8 Graphical illustration of pre-training procedure. We firstly randomly mask the input medical codes using a $[MASK]$ symbol. **Orange arrow**: self-prediction task takes \mathbf{v}_m or \mathbf{v}_d as input to restore the original medical codes with the same type. **Green arrow**: dual-prediction task takes one type of visit embedding such as \mathbf{v}_m or \mathbf{v}_d and tries to predict the other type of medical codes

Visit Embedding

Similar to BERT, a multi-layer Transformer architecture is used to compute the visit embeddings. More specifically, the Transformer model maps a sequence of tokens (ontology embeddings of medical codes) to another sequence of embeddings (Fig. 11.8).

To generate a visit embedding summarizing all medical codes within the visit, a special token $[CLS]$ is introduced to separate all the medical codes between visits. And the output embedding for $[CLS]$ corresponds to the overall embedding for the entire visit. Formally, the visit embedding $\mathbf{v}_*^t \in \mathbb{R}^d$ for a patient at the t -th visit is the first element of the output from the Transformer model:

$$\mathbf{v}_*^t = \text{Transformer}(\{[CLS]\} \cup \{\mathbf{o}_{c_*}^t | c_* \in \mathcal{C}_*^t\})[0] \quad (11.6)$$

where $[0]$ means the first element of the output sequence, $[CLS]$ is a special token as in BERT indicating the first position of each visit of a specific type, visit type $*$ can be diagnosis and medication. In fact, it is more reasonable to use Transformers as encoders (multi-head attention based architecture) than RNN for visit embedding since medical codes within one visit are unordered. Also, the positional embedding is not used in G-BERT as the unordered nature of medical codes within a visit.

Pre-training Before predicting medications, we need to learn better embeddings for medical codes using data from all visits. Many existing temporal models, such

as RNN, require two or more visits per patient to train a model. However, there are many patients with only a single visit, which before G-BERT are mostly ignored in model training.

The pre-training strategy of G-BERT utilizes the data of those single visits. In particular, G-BERT modified the original pre-training tasks, i.e., Masked language model task and Next Sentence prediction task to **self-prediction task** and **dual-prediction task**. The idea to conduct these tasks is to visit embedding, absorb enough information about what it is made of, and predict.

Thus, for the **self-prediction task**, we want the visit embedding v_* to recover *what it is made of*, i.e., the input medical codes C_* limited by the same type for each visit as follows:

$$\begin{aligned}\mathcal{L}_{se}(\mathbf{v}_*, C_*^{(n)}) &= -\log p(C_*^{(n)} | \mathbf{v}_*) \\ &= -\sum_{c_* \in C_*^{(n)}} \log p(c_* | \mathbf{v}_*) + \sum_{c_* \in \{C_* \setminus C_*^{(n)}\}} \log p(c_* | \mathbf{v}_*)\end{aligned}$$

where $C_*^{(n)}$ is the medical codes set of the n -th patient, $* \in \{d, m\}$ indicating diagnosis and medication respectively. And we minimize the binary cross entropy loss \mathcal{L}_{se} . For instance, assume that the n -th patient takes 10 different medications out of total 100 medications which means $|C_m^{(n)}| = 10$ and $|C_m| = 100$. In such case, we instantiate and minimize $\mathcal{L}_{se}(\mathbf{v}_m, C_m^{(n)})$ to produce high probabilities among 10 taken medications captured by $-\sum_{c_* \in C_m^{(n)}} \log p(c_* | \mathbf{v}_m)$ and lower the probabilities among 90 other medications captured by $\sum_{c_* \in \{C_m \setminus C_m^{(n)}\}} \log p(c_* | \mathbf{v}_m)$. In practice, $\text{Sigmoid}(f(\mathbf{v}_*))$ is applied after a fully connected neural network $f(\cdot)$ with one hidden layer. G-BERT also uses specific symbol [MASK] to randomly replace the 15% original medical codes $c_* \in C_*$.

The **dual-prediction task** predicts different types of medical codes based on the other types of codes. For example, we want to predict multiple medications given only the diagnosis codes. Inversely, we can also predict unknown diagnosis given the medication codes. Formally, predicting diagnoses based on medication embedding $p(C_d | \mathbf{v}_m)$ and predicting medications based on diagnosis embedding $p(C_m | \mathbf{v}_d)$ lead the following loss:

$$\mathcal{L}_{du} = -\log p(C_d | \mathbf{v}_m) - \log p(C_m | \mathbf{v}_d)$$

The probability scores $p(C_d | \mathbf{v}_m)$ and $p(C_m | \mathbf{v}_d)$ are learned from two DNN models with a hidden layer $\text{Sigmoid}(f_1(\mathbf{v}_m))$, $\text{Sigmoid}(f_2(\mathbf{v}_d))$. This is a direct adaptation of the next sentence prediction task.

Thus, our final pre-training optimization objective can simply be combining the aforementioned losses, as shown in Eq. (11.7).

$$\mathcal{L}_{pr} = \mathcal{L}_{se}(\mathbf{v}_d, C_d^{(n)}) + \mathcal{L}_{se}(\mathbf{v}_m, C_m^{(n)}) + \mathcal{L}_{du} \quad (11.7)$$

Fine Tuning

After obtaining pre-trained visit representation for each visit, we aggregate all the visit embedding and add a prediction layer for the medication recommendation task. To be specific, from pre-training on all visits, we have a pre-trained Transformer encoder, which can then be used to get the visit embedding \mathbf{v}_*^τ at time τ . The known diagnosis codes \mathcal{C}_d^t at the prediction time t is also represented using the same model as \mathbf{v}_*^t . Concatenating the mean of previous diagnoses visit embeddings and medication visit embeddings, also the last diagnoses visit embedding, we built a DNN based prediction layer to predict the recommended medication codes as:

$$\mathbf{y}_t = \text{Sigmoid}(\mathbf{W}_1[(\frac{1}{t} \sum_{\tau < t} \mathbf{v}_d^\tau) || (\frac{1}{t} \sum_{\tau < t} \mathbf{v}_m^\tau) || \mathbf{v}_d^t] + b)$$

where $\mathbf{W}_1 \in \mathbb{R}^{|\mathcal{C}_m| \times 3d}$ is a learnable transformation matrix.

Given the true labels $\hat{\mathbf{y}}_t$ at each time stamp t , the loss function for the whole EHR sequence (i.e. a patient) is

$$\mathcal{L} = -\frac{1}{T-1} \sum_{t=2}^T (\mathbf{y}_t^\top \log(\hat{\mathbf{y}}_t) + (1 - \mathbf{y}_t^\top) \log(1 - \hat{\mathbf{y}}_t)) \quad (11.8)$$

Results We used EHR data from MIMIC-III [81]. We utilize data from patients with both single visits and multiple visits in the training dataset as the pre-training data source (multi-visit data are split into visit slices). In this work, we transform the drug coding from NDC to ATC Third Level. The statistics of the datasets are summarized in Table 11.4.

Different baseline methods include

1. **Logistic Regression (LR)** is logistic regression with L1/L2 regularization. Here we represent sequential multiple medical codes by the sum of the multi-hot vector of each visit.
2. **LEAP** [174] is an instance-based medication combination recommendation method that formalizes the task in a multi-instance and multi-label learning framework. It utilizes an encoder-decoder based model with an attention mechanism to build complex dependency among diseases and medications.

Table 11.4 Statistics of the data (dx for diagnosis, rx for medication)

Stats	Single-visit	Multi-visit
# of patients	30,745	6350
Avg # of visits	1.00	2.36
Avg # of dx	39	10.51
Avg # of rx	52	8.80
# of unique dx	1997	1958
# of unique rx	323	145

3. **RETAIN** [25] makes sequential prediction of medication combination and diseases prediction based on a two-level neural attention model that detects influential past visits and clinical variables within those visits.
4. **GRAM** [28] injects domain knowledge (ICD9 Dx code tree) to **tanh** via attention mechanism.
5. **GAMENet** [134] is the method to recommend accuracy and safe medication based on memory neural networks and graph convolutional networks by leveraging EHR data and Drug-Drug Interaction (DDI) data source. For a fair comparison, we use a variant of GAMENet without DDI knowledge and procedure codes as input renamed as GAMENet[−].
6. **G-BERT** is our proposed model which integrated the GNN representation into Transformer-based visit encoder with pre-training on single-visit EHR data.

To measure the prediction accuracy, we used the Jaccard Similarity Score (Jaccard), Average F1 (F1), and Precision-Recall AUC (PR-AUC). Jaccard is defined as the size of the intersection divided by the size of the union of ground truth set $Y_t^{(k)}$ and predicted set $\hat{Y}_t^{(k)}$.

$$\text{Jaccard} = \frac{1}{\sum_k^N \sum_t^{T_k} 1} \sum_k^N \sum_t^{T_k} \frac{|Y_t^{(k)} \cap \hat{Y}_t^{(k)}|}{|Y_t^{(k)} \cup \hat{Y}_t^{(k)}|}$$

where N is the number of patients in the test set and T_k is the number of visits of the k th patient.

Table 11.5 compares the performance on the medication recommendation task. Incorporating both hierarchical ontology information and pre-training procedure, the end-to-end model G-BERT has more capacity and achieve comparable results with others.

G-BERT performs better than all the baseline methods including the GAMENet model.

Table 11.5 Performance on medication recommendation task

Methods	Jaccard	PR-AUC	F1	# of parameters
LR	0.4075	0.6716	0.5658	–
GRAM	0.4176	0.6638	0.5788	3,763,668
LEAP	0.3921	0.5855	0.5508	1,488,148
RETAIN	0.4456	0.6838	0.6064	2,054,869
GAMENet [−]	0.4401	0.6672	0.5996	5,518,646
GAMENet	0.4555	0.6854	0.6126	5,518,646
G-BERT	0.4565	0.6960	0.6152	3,034,045

The bold values indicate the best performance numbers across all methods

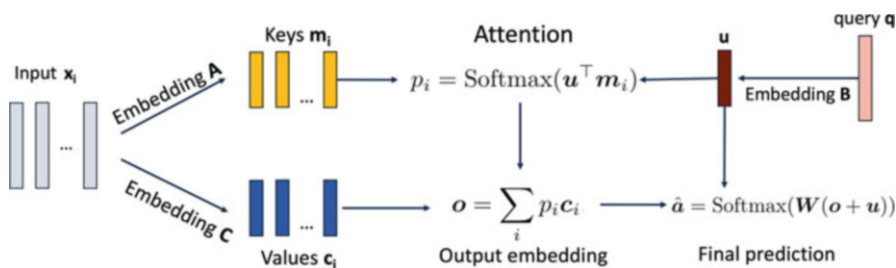


Fig. 11.9 End-to-end memory network exercise

11.8 Exercises

- What is the main difference between the end-to-end memory network and the self-attention/transformer?
- What is the most innovative step in the original memory network model and why?
- What is NOT true about end-to-end memory networks as illustrated in Fig. 11.9?
 - It removes the arg max operation in the original memory networks so that gradient can be computed at all the steps.
 - The size of the embeddings in keys, values, and query are all the same.
 - Multi-layer end-to-end memory network can be created via multiple embedding matrices for keys and values
 - The model parameters are A , B , C and W matrices.
- What is NOT true about self attention?
 - Three versions of embeddings of the input X are produced namely Q , K , V .
 - The main reason to have multiple embedding matrices of X is to be able to apply and learn those embeddings independently in parallel.
 - The temporal dependency in the input sequences is essential in self-attention models.
 - Multiple versions of self-attention are concatenated together to produce more robust embedding, which is called multi-head self attention.
- What are the method ideas used in Transformer encoder? [multiple correct choices]
 - Self attention
 - Positioning encoding
 - Multi-head attention
 - Residual connection

6. What are method ideas used in Transformer decoder? [multiple correct choices]
 - (a) Masked attention
 - (b) Self attention
 - (c) Residual connection
 - (d) Recurrent neural networks
7. What is NOT true about BERT model?
 - (a) BERT provides a masked language model.
 - (b) BERT masks x% of input words and try to predict them with other words.
 - (c) BERT uses a transformer based model.
 - (d) BERT provides static embeddings for all words.
8. What is the healthcare application in Doctor2Vec?
9. What is the healthcare application in GameNet?
10. What are the method ideas behind G-BERT? [multiple correct choices]
 - (a) Medical ontology is used to embed medical concepts.
 - (b) BERT model can be used to pre-train on large number of clinical visits.
 - (c) Temporal dependencies across visits are important in prediction.
 - (d) Multi-task learning is useful in clinical applications.