

9.6. [Learning Notes] TieNet: Text-Image Embedding Network for common Thorax Disease Classification and Reporting in Chest X-rays

1. Aim: Extract the image & text representations

① multiple label disease prediction (= disease classification)

② use only image as input, output disease classification and a preliminary report together.

② LSTM:

$$H = (\vec{h}_1, \dots, \vec{h}_T) \quad [dh \times T]$$

\vec{h}_t : hidden layer in LSTM $[dh \times 1]$

dh : LSTM state space of dimension dh .

\vec{X} : the convolutional activations from the transition layer. $[D \times D \times C]$

\vec{a}_t : subsequent attention map.

$$\vec{h}_t = \text{LSTM}([w_t, a_t * X], h_{t-1}) \quad [dh \times 1]$$

2. Method

1. End-to-End Trainable CNN-RNN Model

Input 1: Text Report $\rightarrow \vec{w}_t$

$$S = \{\vec{w}_1, \dots, \vec{w}_T\}, \vec{w}_t \in \mathbb{R}^r \quad [dw \times T]$$

\vec{w}_t : vector standing for a dw dimensional word embedding for the t -th word in report $[dw \times 1]$

V : size of vocabulary

T : length of the report.

Input 2: Image $\rightarrow \vec{X}$

2. Attention Encoded Text Embedding

$$\text{① weights } G = \text{softmax}(\vec{w}_{s2} \tanh(\vec{u}_{s1} \vec{H})) \quad [r \times T]$$

r : number of global attentions we want to extract from the sentence.

\vec{u}_{s1} : $[s\text{-by-}dh]$ matrix

\vec{w}_{s2} : $[r\text{-by-}s]$ matrix

s : a hyperparameter governing the dimensionality.

$$\text{② embedding matrix } \vec{M} = \vec{G} \vec{H} \quad [r \times dh]$$

\vec{g}^i : how much each hidden state contributes to the final embedding representation of \vec{M} .

③ \hat{X}_{AETE} . $\boxed{1 \times d_h} \rightarrow \text{Test}$

To provide a final global text embedding of the sentences in report, the AETE executes max-over-r pooling across M , producing an embedding vector \hat{X}_{AETE} with size d_h .

3. Saliency Weighted Global Average Pooling

$$\hat{X}_{sw-gap}(c) = f(\vec{a}_c, \vec{g}_c, X) \quad \boxed{1 \times C}$$

\hat{X}_{sw-gap} representing the global visual information, guided by both text- and visual-based attention.

4. Joint Learning.

① Combine both text & image to produce final classification.

$$\hat{X} = [\hat{X}_{AETE} ; \hat{X}_{sw-gap}]$$

use a Fully-connected Layer to produce the output for multi-label classification

④ Balance the large difference between image loss & text loss.

5. Image Auto-Annotation

→ mine image classification labels.

Test: ① R: reports only.

Input: only reports. (no image).

Output: generated reports.
(image classification label)

② I+R: image + report pairs

Input: Image + report pairs.

Output: image classification label.

Train: All use Chest X-ray 14 (Report+Image)

6. Automatic Classification & Reporting of Thorax Diseases.

→ only have images, classify label.

Train: Input: use both text & image.

Test: Input: one image

Output: generated text/report.

② image classification label.
(multiple label prediction)

9.6. Explainable Prediction of Medical codes from Clinical Text.

medical codes $\left\{ \begin{array}{l} \text{diagnoses} \\ \text{treatment (procedures)} \end{array} \right.$

Aim - predict medical codes from clinical text. (notes).
(ICD codes)
↓
labels Document.

Method. classification: Documents - labels

Application: decision support setting:

- ① explain why it predict each other &
- ② explain what parts of the text are most relevant to each code

Method -

1. Convolutional Architecture:
to translate docs to $H: d_e \times n$ dimensional matrix.
for pre-label attention mechanism.

Input: $X = [x_1, x_2, \dots, x_n] \quad [d_e \times N]$

n : length of document.

x_i : d_e : dimensional pre-trained embeddings for each word in the document

Output: $H: [R^{d_e \times N}]$

d_e : size of the filter output.

n : length of document.

2. Attention. (Instead Pooling to reduce the matrix to a vector)

aim: to assign multiple labels (i.e. medical codes) for each document

② An additional benefit: it selects the k -grams from the text that are most relevant to each predicted label

Input: $H: [R^{d_e \times N}]$

Output:

① resulting vector: $H^T u_e^i \quad (u_e^i \in R^{d_e})$
 $[R^{n \times 1}]$

② attention vector a :

$$a_e = \text{Softmax}(H^T u_e^i) \quad [R^{1 \times N}]$$

② vector representations for each label

$$v_e = \sum_{n=1}^N a_{e,n} h_n \quad [d_e \times 1]$$

baseline model: compute a single vector v for all labels.

$$v_i = \max_n h_{n,i}$$

$$[d_e \times 1]$$

3. Classification

aim: compute a probability for label l .

using linear layer & a sigmoid transformation.

$$\hat{y}_l = \sigma(\beta_l^T v_c + b_l)$$

Output: 1. (a probability)

? z_l : a max-pooled vector,

Variant	Description	CAM
·	VID-CAM	

4. Datasets

1. ICD codes (labels) : 8921 {diagnose 6818
procedure 2003}

documents (discharge summaries) $\approx 50,000$

training	47714
validation	1632
testing	3372

4. Training:

minimize the binary cross-entropy loss.

L_2 norm of the ~~model~~ model weights,
using Adam optimizer.

CAM Done

2. Secondary Evaluations:
Mimic II - 50

ICD top codes (labels) : 50

documents (discharge summaries)

training	8067
	1574
	1730

3. Mimic II Full

labels : ~~20,533~~ 5031

Documents (discharge summaries)

Training	20,533
Testing	2282

Optimization:

5. Embedding label descriptions.

Aim: Due to the dimensionality of the label space, many codes are rarely observed in the labeled data. To improve performance on these codes, we add the following regularizing objective to our loss.

$$L(x, y) = L_{BCE} + \dots$$

4. Evaluation.

AUC (micro-R, Macro-R)

F₁ (micro-R, Macro-R)

P@n.

4.1. Evaluation of Interpretability

4.2 Results:

Evaluator to judge. (Table 7).