

Project routemap

- Reproduce the results
- Modify the model to improve the performance

Reproduction

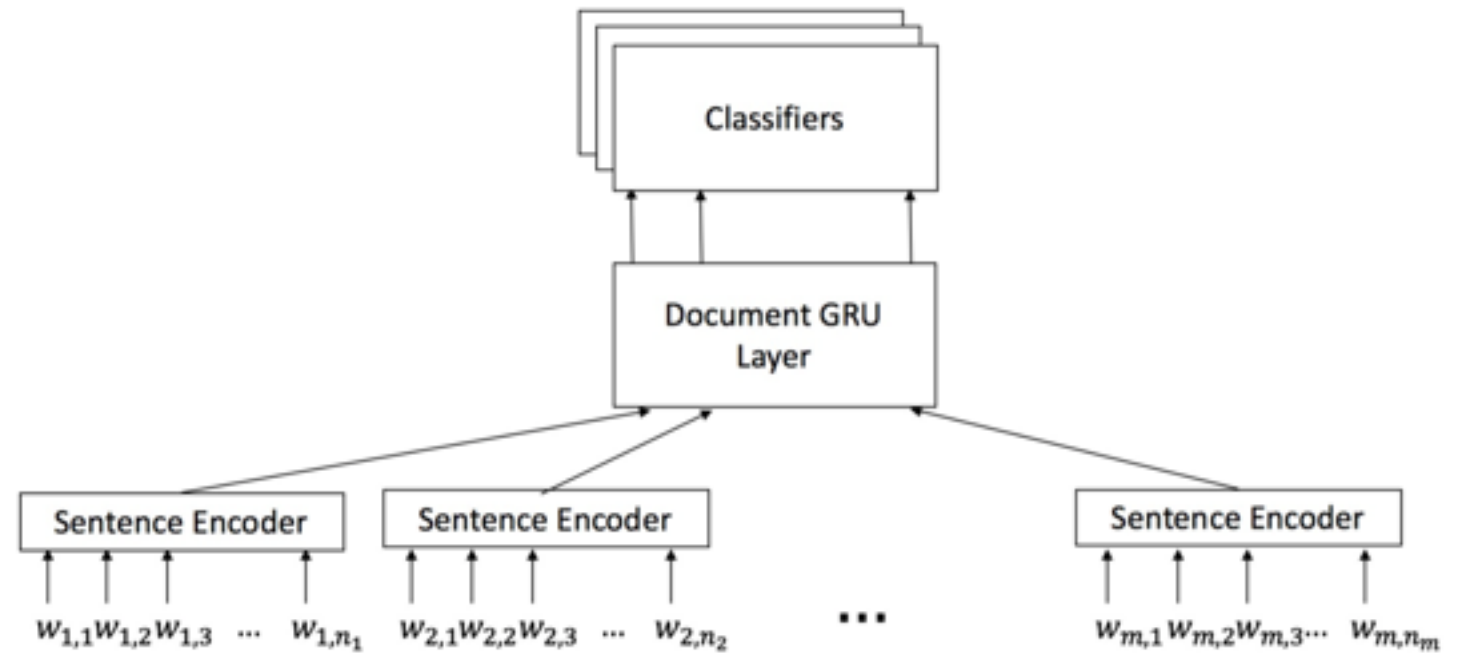
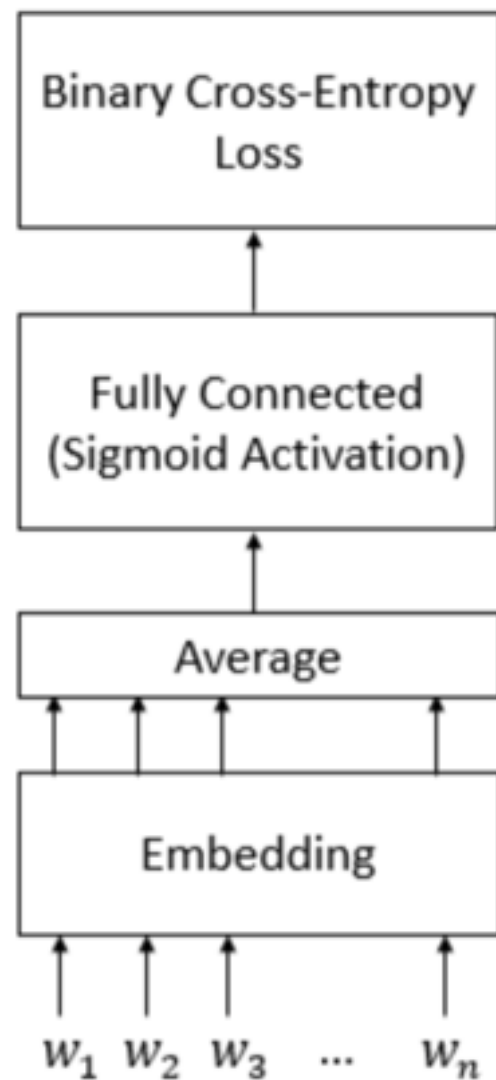
- Reference: *Multi-label classification of patient notes a case study on icd code assignment*
- Problem: assign ICD9 codes to patient's **discharge summaries**.
- Dataset: MIMIC II (22,000+) and MIMIC III (65,000+)
- Labels: 4,000+ and 6,000+ ICD9 codes respectively
- Evaluation: Macro-F

Preprocessing

- Data mainly comes from table **DIAGNOSE_ICD** and **NOTEEVENT**
- Case-sensitive tokenisation
 - How about case-insensitive?
- Digits normalisation
 - Do the numerical values matter? How about non-linear transformation based on numerical values' context?
- Map OOVs to its nearest neighbour by the edit distance
 - How about representing OOV by UNK?
- Hierarchical segmentation: tokens & sentences (LONG DOCUMENTS)
 - How about discourse segmentation

Models

CBOW



Our models

