

Multi-label Classification on Networked Data

Han Xiao

Department of Computer Science
Aalto University

September 2017

overview

- ▶ **goal:** multi-label classification of posts on QA-sites
- ▶ **ingredient 1:** deep learning for supervised learning
- ▶ **ingredient 2:** graph embedding learning for under supervised learning
- ▶ both requires convex optimization techniques

motivation 1

- ▶ goal: learn a automatic question tagger
- ▶ example: stackoverflow questions
- ▶ tags: python, java, android, ios, ...

Correlating activity between entities using Python

python

statistics

modified 9 mins ago [xboard](#) 16

How to know the model has started overfitting?

neural-network

overfitting

answered 16 mins ago [David Makovoz](#) 11

What are the tools to plot cluster results?

machine-learning

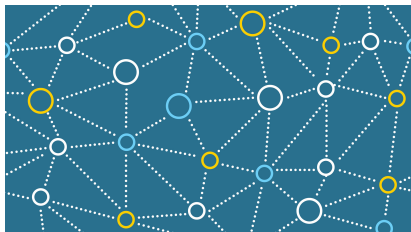
python

clustering

answered 42 mins ago [Anony-Mousse](#) 3,185

motivation 2

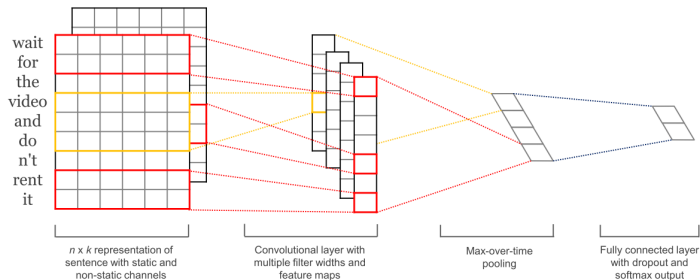
- ▶ posts and users are linked in a network
 - ▶ user *post/likes/comment* posts
- ▶ can we leverage *network information* for tagging purpose?



multi-label classification

- ▶ given n pairs of feature vectors and label vectors, $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^D$ and $\mathbf{y}_i \in \{0, 1\}^L$
- ▶ learn a function $f : \mathbb{R}^D \rightarrow \{0, 1\}^L$ that minimizes some error function
- ▶ our problem can be modeled as such

multi-label classification: CNN-based approach



1

- loss function: softmax is replaced by sigmoid

¹Kim, Yoon. "Convolutional neural networks for sentence classification." EMNLP (2014).

graph embedding

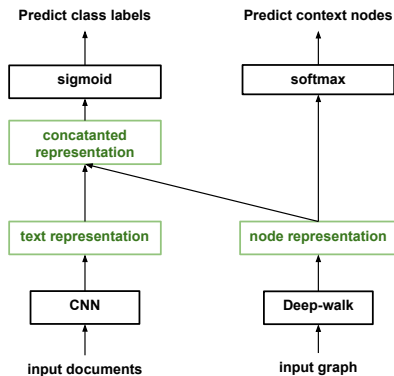
- ▶ given $G = (V, E)$, $V = \{1, \dots, n\}$, $E \subseteq V \times V$
- ▶ learn a function $f : V \rightarrow \mathbb{R}^d$
- ▶ $f(v)$ is the embedding for node v
- ▶ demonstrated to improve clustering, classification performance in general

example: DeepWalk ²

- ▶ idea borrowed from word2vec (a word predicts its context/surrounding word)
- ▶ here, think of a node as a word
- ▶ a node predicts its “context” nodes
- ▶ context nodes are collected by random walk

²Perozzi, Bryan, Rami Al-Rfou, and Steven Skiena. “Deepwalk: Online learning of social representations.” KDD, 2014.

model & architecture



training objective: $\mathcal{L} = \mathcal{L}_s + \mathcal{L}_u$

model: unsupervised term \mathcal{L}_u

- ▶ essentially deepwalk
- ▶ given $\{(i, c)\}$ pairs, learn a model that minimizes:

$$\begin{aligned}\mathcal{L}_u &= -\mathbb{E}_{(i,c)} \log \Pr(\mathbf{e}_c \mid \mathbf{e}_i) \\ &= -\mathbb{E}_{(i,c)} \log \frac{\mathbf{e}_c^T \mathbf{e}_i}{\sum_{c'} \mathbf{e}_{c'}^T \mathbf{e}_i}\end{aligned}$$

- ▶ (i, c) pairs collected from random walk

model: supervised term \mathcal{L}_s

- ▶ given $\{(\mathbf{x}_i, \mathbf{e}_i, \mathbf{y}_i)\}_{i=1}^L$ triplets, learn a model that minimizes:

$$\mathcal{L}_s = \frac{1}{L} \sum_{i=1}^K \log p(\mathbf{y}_i \mid \mathbf{x}_i, \mathbf{e}_i)$$

- ▶ model similar to KimCNN (using sigmoid instead of softmax)

training

- ▶ alternating update of \mathcal{L}_u and \mathcal{L}_s
- ▶ optimizer: ADAM

experiment: datasets & evaluation

- ▶ 3 sites from stackexchange.com:

	data science	emacs	software engineering
#instances	5145	4536	10336
#labels	327	618	1344
avg #labels per	2.7	2.0	2.6

- ▶ question network construction: two questions are linked they are linked to at least one common user
- ▶ train/dev/test ratio: 80%/10%/10%
- ▶ evaluation: precision@k, $k = \{1, 3, 5\}$

methods

- ▶ *fastxml*: tf-idf features
- ▶ *cnn*: sigmoid loss function
- ▶ *cnn* + *deepwalk*

results: data science

	<i>fastxml</i>	<i>cnn</i>	<i>cnn+deepwalk</i>
<i>p@1</i>	0.26	0.53	0.57
<i>p@3</i>	0.17	0.37	0.38
<i>p@5</i>	0.15	0.29	0.28

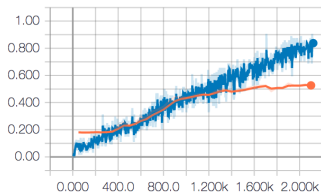
results: emacs

	<i>fastxml</i>	<i>cnn</i>	<i>cnn+deepwalk</i>
<i>p@1</i>	0.21	0.58	0.56
<i>p@3</i>	0.10	0.30	0.30
<i>p@5</i>	0.07	0.20	0.20

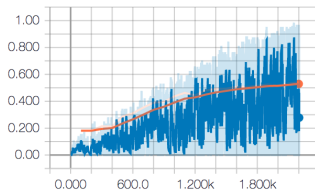
results: software engineering

	<i>fastxml</i>	<i>cnn</i>	<i>cnn+deepwalk</i>
<i>p@1</i>	0.06	0.41	0.43
<i>p@3</i>	0.06	0.26	0.25
<i>p@5</i>	0.05	0.18	0.18

train/dev curve



(a) cnn



(b) cnn+deepwalk

Figure: precision@1 for *emacs*

conclusion

- ▶ I'm surprised at CNN's performance
- ▶ also surprised that network embedding does not help much
- ▶ possible reasons: improper network construction, formulation, training, etc

future plan

- ▶ try different ways of learning network embedding
- ▶ understand how graph embedding can help
- ▶ design scalable methods for classification (especially when label space is large)