



Classification of biomedical texts

Data description

Mario Sängner (WBI, HU Berlin)

saengema@informatik.hu-berlin.de

Classification tasks

- You are given a data set of 15301 abstracts of biomedical articles
- Goal: Implement a classification pipeline to perform the following tasks
 - 1) Is the article concerned with cancer? (yes / no)
 - 2) About which specific type of cancer is the article about? (multi-class)

Training data format

- The training data set is given as **tab-separated** (“\t”) file (*train.tsv*)
 - Each line represents one biomedical article (~ training instance)

1	pmid cancer_type doid is_cancer text
2	27581830 no cancer -1 0 The management of benign non-infective pleur
3	1213020 cancer 162 1 Bulky lymphadenopathy with poor clinical outco
4	22374460 endocrine gland cancer 170 1 Multicentre phase II trial of
5	27586118 no cancer -1 0 How safe are children when transported by bi
6	25868852 cancer 162 1 Loss of INPP4B causes a DNA repair defect thro
7	27599572 no cancer -1 0 TAK1 inhibitor NG25 enhances doxorubicin-med
8	27578867 no cancer -1 0 Cd47-Sirpα interaction and IL-10 constrain i
9	27581024 no cancer -1 0 Visualizing the Tumor Microenvironment of Li
10	27576518 no cancer -1 0 Brain damage resembling acute necrotizing en
11	25199829 cancer 162 1 RAF suppression synergizes with MEK inhibition
12	24962318 melanoma 1909 1 Targeting TBK1 inhibits migration and resis
13	22496619 cancer 162 1 Oncogenic KRAS impairs EGFR antibodies' effici
14	27597568 no cancer -1 0 Multiwavelength metasurfaces through spatial
15	18003960 gastrointestinal system cancer 3119 1 Cetuximab for the tre
16	27576507 no cancer -1 0 Expression of Helios in gastric tumor cells

Training data format

- The training data set is given as **tab-separated** ("`\t`") file (*train.tsv*)
 - Each line represents one biomedical article (\sim training instance)
- Data columns:
 - *pmid*: unique (PubMed) identifier of the article
 - *cancer_type*: name of the specific cancer type the article is about or "no cancer", if the article isn't concerned with cancer
 - *doiid*: unique identifier of the specific cancer type or -1, if the article isn't concerned with cancer
 - *is_cancer*: indicates whether the article is about cancer or not
 - *text*: abstract of the article

Evaluation

- Your classification model(s) will be evaluated on a **hold-out evaluation data set**
- For evaluation **only the articles** will be given (in the same tab-separated format)
 - I.e. the columns *is_cancer*, *cancer_type* and *doid* are missing!
 - Your classification model(s) should predict *is_cancer* and *doid*
 - See train_blind.tsv as an example!

1	pmid text
2	27581830 The management of benign non-infective pleural effusions.
3	17213020 Bulky lymphadenopathy with poor clinical outcome is assoc:
4	22374460 Multicentre phase II trial of trastuzumab and capecitabine
5	27586118 How safe are children when transported by bicycle? With tl
6	25868852 Loss of INPP4B causes a DNA repair defect through loss of
7	27599572 TAK1 inhibitor NG25 enhances doxorubicin-mediated apoptos:

Submission format

- For evaluation you have to submit **two files** containing the predictions of your classification model
 - *task1.tsv*: Containing the predictions of your model for task 1 (about cancer?)
 - Two columns: *pmid* and *is_cancer* (0/1)
 - Use `"\t"` as column separator

	<u>pmid</u>	is_cancer
1	27581830	0
2	27599572	1
3	24962318	0
4		

Submission format

- For evaluation you have to submit **two files** containing the predictions of your classification model
 - *task2.tsv*: Containing the predictions of your model for task 2 (cancer-type)
 - Two columns: *pmid* and *doid*
 - Use `"\t"` as column separator
 - Include **all articles from the test set**, i.e. also articles that aren't about cancer at all have to be included!

1	<u>pmid</u>	<u>doid</u>
2	27581830	-1
3	27599572	162
4	24962318	1909
5	24962318	-1
6	...	
7		