



چالش سوم: Toxic Comment Classification Challenge

فراہم آورنده: سینا حیدری

زمستان ۹۸

برای مشاهده جزئیات این چالش در گگل به لینک زیر مراجعه کنید.

<https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/overview>

۱ در مورد این چالش

معمولا، به اشتراک گذاردن مسائل شخصی کار دشواری انسان‌ها است؛ تهدید آزار و اذیت آنلاین باعث شده مردم از انجام این کار در پلتفرم‌های آنلاین صرف نظر کنند. پلتفرم‌ها برای آرامش بخشیدن به گفت‌وگوها با مشکلاتی دست‌وپنجه نرم می‌کنند؛ تا جایی که بسیاری از انجمن‌ها، قابلیت نظر دادن کاربران را از دسترس خارج کرده‌اند. تیم Conversation AI، یک ابتکار دانش‌بنیان است که توسط Jigsaw و Google پایه‌گذاری شده است. این تیم در حال کار کردن روی ابزارهایی است که گفت‌وگوی آنلاین را بهبود دهد. یکی از زمینه‌هایی که این تیم روی آن تمرکز کرده، رفتارهای منفی آنلاین همچون نظرات سمی است. اینگونه نظرات بی‌ادبانه‌اند و همچنین، موجب بی‌احترامی به کاربران و در نهایت، ترک آن‌ها از گفت‌وگو می‌شود. تیم Conversation API تا کنون محدودی از مدل‌ها را از طریق Perspective API در اختیار عامه قرار داده است-از جمله مدل سمیت-. اما مدل‌های کنونی هنوز دچار خطا شده و به کاربر اجازه‌ی انتخاب نوع سمیت را نمی‌دهند-برای نمونه؛ بعضی پلتفرم‌ها با ناسزا مشکلی ندارند؛ اما با نوع‌های دیگر سمیت، اینگونه برخورد نمی‌کنند-. در این رقابت، ما به چالش کشیده شده‌ایم تا یک مدل چندسره، با توانایی بهتر از مدل کنونی برای تشخیص انواع سمیت-از جمله؛ تهدید، فحاشی و تنفرهای بر پایه‌ی هویت-سازیم. ما یک مجموعه داده‌ی نظرات که از گفت‌وگوی صفحه‌های ویرایش ویکی‌پدیا جمع‌آوری شده‌اند را استفاده خواهیم کرد. **هشدار:** مجموعه داده‌ی این رقابت شامل متنی است که ناسزا، مبتذل و توهین‌آمیز است.

۲ در مورد این گزارش

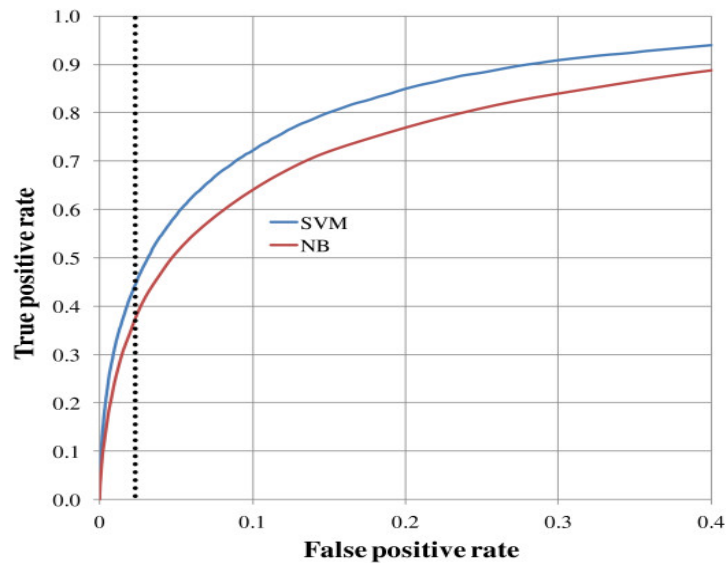
این گزارش چگونگی استفاده از NBSVM^۱ در جهت ساختن یک پایه‌ی قوی برای رقابت را نشان می‌دهد. NBSVM توسط Sida Wang و Chris Manning در مقاله‌ی:

Baselines and Bigrams: Simple, Good Sentiment and Topic Classification

معرفی شد. در این گزارش ما به جای SVM از رگرسیون لجستیک^۲ در کتابخانه‌ی sklearn بهره می‌بریم-از آنجایی که عملا یکسان هستند-.

^۱ Naive Bayes - Support Vector Machine
^۲ logistic regression

۳ راه حل



ما در این پیاده‌سازی از NB-SVM برای مدل‌سازی استفاده خواهیم کرد.

Bag of Words

در این رویکرد، ابتدا مجموعه داده را که شامل داکيومنت‌ها و لیبل مربوطه است به نمایش bag of words می‌آوریم؛ داده ساختار مورد استفاده یک ماتریس اسپارس می‌باشد که Term-document matrix نامیده می‌شود. متن اولین نظر و لیبل آن را در زیر مشاهده می‌کنید:

```
trn[0]

" A formal orchestra audience is turned into an insane, violent mob
by the crazy chantings of it's singers. Unfortunately it stays absurd
the WHOLE time with no general narrative eventually making it just
too off putting. Even those from the era should be turned off. The
cryptic dialogue would make Shakespeare seem easy to a third grader.
On a technical level it's better than you might think with some good
cinematography by future great Vilmos Zsigmond. Future stars Sally
Kirkland and Frederic Forrest can be seen briefly."

trn_y[0]

0
```

در ادامه، CountVectorizer مجموعه‌ی نظرات را به یک ماتریس که ستون‌های آن نشان دهنده‌ی توکن‌ها و مقادیر

نشان دهنده‌ی تعداد تکرار آن توکن است در می‌آورد. دستور `fit_transform(trn)` مجموعه‌داده‌ی آموزشی را به‌ماتریس `term-document` تبدیل می‌کند.

تعریف بیز ساده

ما در اینجا `log-count ratio` را برای هر کلمه تعریف می‌کنیم.

$$r = \log \frac{\text{ratio of feature } f \text{ in positive documents}}{\text{ratio of feature } f \text{ in negative documents}}$$

در فرمول بالا، **ratio of feature in positive documents** برابر است با: تعداد دفعاتی که یک داکيومنت با لیبل مثبت (یک) شامل کلمه می‌شود، بخش بر تعداد داکيومنت‌هایی که لیبل مثبت دارند. نکته‌ی قابل توجه در اینجا، بهتر عمل کردن بیز ساده‌ی دودویی نسبت به بیز ساده است.

رگرسیون لجستیک

در علوم داده، رگرسیون لجستیک به‌دلیل واقع‌گرایانه بودن بهتر از نسبت به‌دست آمده در مرحله‌ی قبل عمل می‌کند. در زیر، پیاده‌سازی و پیش‌بینی با مدل رگرسیون لجستیک را مشاهده می‌کند. دقت پیش‌بینی برای نسخه‌های `regularize` نشده و `regularize` شده را در زیر مشاهده می‌کنید.

Unregularized version

```
m = LogisticRegression(C=1e8, dual=True)
m.fit(x, y)
preds = m.predict(val_term_doc)
(preds==val_y).mean()

0.8550400000000002

m = LogisticRegression(C=1e8, dual=True)
m.fit(trn_term_doc.sign(), y)
preds = m.predict(val_term_doc.sign())
(preds==val_y).mean()

0.8548799999999997
```

Regularized version

```
m = LogisticRegression(C=0.1, dual=True)
m.fit(x, y)
preds = m.predict(val_term_doc)
(preds==val_y).mean()

0.8827599999999999

m = LogisticRegression(C=0.1, dual=True)
m.fit(trn_term_doc.sign(), y)
preds = m.predict(val_term_doc.sign())
(preds==val_y).mean()

0.88404000000000005
```

Adding bigrams and trigrams

مدل بعدی ما یکی از نسخه‌های رگرسیون لجستیک با ویژگی‌های بیز ساده می‌باشد. حال برای هر داکيومنت ویژگی‌های باینری را-همانگونه که در مراحل قبلی توضیح دادیم-به‌دست می‌آوریم، با این تفاوت که این بار از bigram-ها و trigram-ها نیز بهره می‌بریم. هر ویژگی نسبت log-count می‌باشد. پیش‌بینی وقتی ویژگی‌ها trigram-ها هستند:

```
m = LogisticRegression(C=0.1, dual=True)
m.fit(x, y);

preds = m.predict(val_x)
(preds.T==val_y).mean()

0.90500000000000003

np.exp(r)

matrix([[ 0.94678,  0.85129,  0.78049, ...,  3.      ,  0.5      ,  0.5
]])
```

پیش‌بینی وقتی ویژگی‌ها نسبت‌های log-count هستند:

```
x_nb = x.multiply(r)
m = LogisticRegression(dual=True, C=0.1)
m.fit(x_nb, y);

val_x_nb = val_x.multiply(r)
preds = m.predict(val_x_nb)
(preds.T==val_y).mean()
```

0.91768000000000005

مشاهده می‌کنید که دقت پیش‌بینی با این مدل نهایی 0.917680 می‌باشد. ما در پیاده‌سازی از NB-SVM ای که توسط fast.ai ارائه شده استفاده خواهیم کرد، دقت این NB-SVM ، برای همین نمونه‌یی که بررسی شد، به 0.92129 می‌رسد.

۴ پیاده‌سازی

بارگذاری پکیج‌های مورد نیاز

```
1 import pandas as pd, numpy as np
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.feature_extraction.text import CountVectorizer,
  TfidfVectorizer
```

خواندن داده

```
1 train = pd.read_csv('../input/train.csv')
2 test = pd.read_csv('../input/test.csv')
3 subm = pd.read_csv('../input/sample_submission.csv')
4
```

```
1 train.head()
2
```

Out[3]:

	id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
0	0000997932d777bf	Explanation\nWhy the edits made under my usern...	0	0	0	0	0	0
1	000103f0d9c9cfb60f	D'aww! He matches this background colour I'm s...	0	0	0	0	0	0
2	000113f07ec002fd	Hey man, I'm really not trying to edit war. It...	0	0	0	0	0	0
3	0001b41b1c6bb37e	"\nMore\nI can't make any real suggestions on ...	0	0	0	0	0	0
4	0001d958c54c6e35	You, sir, are my hero. Any chance you remember...	0	0	0	0	0	0

دو نمونه از نظرهای این مجموعه داده که اولی سمی و دومی بدون لیبل است را در زیر مشاهده می‌کنید:

```
1 train['comment_text'][0]
```

```
Out[4]:
```

```
"Explanation\n\nWhy the edits made under my username Hardcore Metallica Fan were reverted? The y weren't vandalism, just closure on some GAs after I voted at New York Dolls FAC. And please don't remove the template from the talk page since I'm retired now.89.205.38.27"
```

```
1 train['comment_text'][2]
```

```
Out[5]:
```

```
"Hey man, I'm really not trying to edit war. It's just that this guy is constantly removing relevant information and talking to me through edits instead of my talk page. He seems to care more about the formatting than the actual info."
```

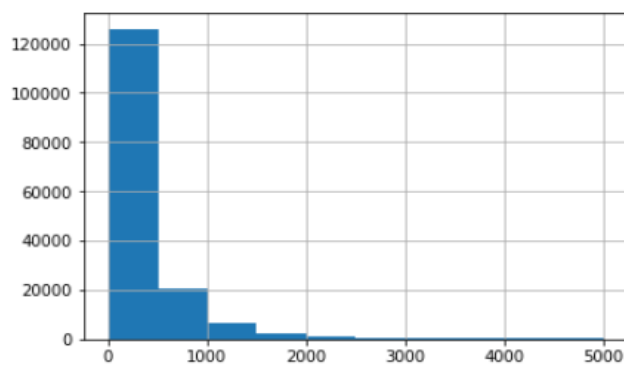
طول کامنت‌ها با هم بسیار متفاوت هستند:

```
1 lens = train.comment_text.str.len()  
2 lens.mean(), lens.std(), lens.max()  
3
```

```
Out[6]:
```

```
(394.0732213246768, 590.7202819048923, 5000)
```

```
1 lens.hist()  
2
```



ما برای پیش‌بینی یک لیست از لیبل‌ها خواهیم ساخت، همچنین یک لیبل 'none' به داده اضافه می‌کنیم تا بدانیم چند نظر هیچ لیبل‌ی نخورده‌اند. سپس می‌توانیم مجموعه داده را خلاصه‌سازی کنیم.

```

1 label_cols = ['toxic', 'severe_toxic', 'obscene', 'threat', 'insult', '
    identity_hate']
2 train['none'] = 1-train[label_cols].max(axis=1)
3 train.describe()
4

```

Out[8]:

	toxic	severe_toxic	obscene	threat	insult	identity_hate	
count	159571.000000	159571.000000	159571.000000	159571.000000	159571.000000	159571.000000	159571.000000
mean	0.095844	0.009996	0.052948	0.002996	0.049364	0.008805	0.898321
std	0.294379	0.099477	0.223931	0.054650	0.216627	0.093420	0.302226
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.000000
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.000000
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.000000
max	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

```

1 len(train), len(test)
2

```

Out[9]:

```
(159571, 153164)
```

تعدادی نظر خالی در مجموعه داده وجود دارد که باید حذف شود؛ در غیر این صورت، sklearn خطا خواهد داد.

```

1 COMMENT = 'comment_text'
2 train[COMMENT].fillna("unknown", inplace=True)
3 test[COMMENT].fillna("unknown", inplace=True)
4

```

Out[9]:

```
(159571, 153164)
```

ساختن مدل

کار را با پدیدآوردن یک نمایش bag of words، به عنوان term document matrix آغاز می کنیم. همانگونه که در مقاله‌ی NBSVM پیشنهاد شده، از ngrams استفاده خواهیم کرد.


```

1 import re, string
2 re_tok = re.compile(f'([{{string.punctuation"}}&','"«»@´°%&_i$£])')
3 def tokenize(s):
4     return re_tok.sub(r' \1 ', s).split()
5

```

```

Out[9]:
(159571, 153164)

```

ظاهراً استفاده از TF-IDF، prior بهتری نسبت به ویژگی‌های دودویی استفاده شده در مقاله نتیجه می‌دهد. در این چالش، امتیاز در جدول با استفاده از TF-IDF، از 0.59 به 0.55 کاهش می‌یابد.

```

1 n = train.shape[0]
2 vec = TfidfVectorizer(ngram_range=(1,2), tokenizer=tokenize,
3 min_df=3, max_df=0.9, strip_accents='unicode', use_idf=1,
4 smooth_idf=1, sublinear_tf=1 )
5 trn_term_doc = vec.fit_transform(train[COMMENT])
6 test_term_doc = vec.transform(test[COMMENT])
7

```

```

Out[9]:
(159571, 153164)

```

این یک ماتریس اسپارس با تعداد کمی از مقادیر غیر صفر می‌سازد-ذخیره شده در نمایش زیر-.

```

1 trn_term_doc, test_term_doc
2

```

```
Out[13]:
(<159571x426005 sparse matrix of type '<class 'numpy.float64'>'
  with 17775104 stored elements in Compressed Sparse Row format>,
 <153164x426005 sparse matrix of type '<class 'numpy.float64'>'
  with 14765755 stored elements in Compressed Sparse Row format>)
```

تابعی برای محاسبه‌ی پیش‌آمدها توسط فرمول بیز:

```
1 def pr(y_i, y):
2     p = x[y==y_i].sum(0)
3     return (p+1) / ((y==y_i).sum()+1)
4
```

```
1 x = trn_term_doc
2 test_x = test_term_doc
3
```

تابعی برای فیت کردن مدل‌ها-برای هر لیبل یک مدل فیت می‌کنیم-:

```
1 def get_mdl(y):
2     y = y.values
3     r = np.log(pr(1,y) / pr(0,y))
4     m = LogisticRegression(C=4, dual=True)
5     x_nb = x.multiply(r)
6     return m.fit(x_nb, y), r
7
```

فیت کردن مدل‌ها و اضافه کردن پیش‌بینی برای هر سطر به preds:

```
1 preds = np.zeros((len(test), len(label_cols)))
2
3 for i, j in enumerate(label_cols):
4     print('fit', j)
5     m,r = get_mdl(train[j])
6     preds[:,i] = m.predict_proba(test_x.multiply(r))[:,1]
7
```

```
fit toxic
fit severe_toxic
fit obscene
fit threat
fit insult
fit identity_hate
```

References

- [1] Naive Bayes(NB)-Support Vector Machine(SVM): Art Of State Result Hands-On Guide using Fast.ai
- [2] NB-SVM strong linear baseline
- [3] NER - training using spacy (Ensemble)
- [4] Baselines and Bigrams: Simple, Good Sentiment and Topic Classification (Also saved in the ZIP file!)