



دانشگاه تهران

پردیس دانشکده‌های فنی

دانشکده مهندسی برق و کامپیوتر

---

ترکیب داده / اطلاعات

گزارش پروژه‌ی درس

---

محمدحسین بدیعی

شماره دانشجویی 810199106

گرایش: کنترل – هوش مصنوعی و رباتیک

استاد: جناب آقای دکتر بهزاد مشیری

بهار 1399-1400

## مقدمه

پردازش زبان طبیعی انسان یا NLP ها کاربرد محسوسی در سیستم‌های هوشمند دارند. از جمله کاربردهای آن برقراری یک تعامل دو طرفه بین انسان و ماشین است. در صورت برقراری این ارتباط به نحو مطلوب، ماشین قادر خواهد بود مکالمات و نوشته‌های انسانها را تشخیص داده و پاسخی مناسب را به آنها ارائه کند. نمود اجرا شده از این کاربردها را می‌توان در دستیار هوشمند اپل و گوگل مشاهده کرد. کاربرد دیگری که حائز اهمیت می‌باشد، تحلیل و آنالیز داده‌های گفتاری و نوشتاری انسان‌هاست که هم اکنون شرکت‌هایی مانند کمپانی اپل، ماکروسافت، آمازون و ... بر روی این حوزه تحقیقات وسیعی را انجام می‌دهند. یکی از اهداف این آنالیز کردن‌ها، ارائه‌ی خدماتی بهتر به مشتریان و دستیابی به سوددهی بیشتر برای شرکت است. نمونه‌ی عملی این واقعیت را می‌توان در تحلیل فیدبک‌های گرفته شده از مشتریان آمازون دانست. این فیدبک‌ها می‌توانند یک کالا را از دسته‌ی غیر محبوب به محبوب برده و یا به دسته‌ی پیشنهادهای ویژه‌ی سایت اضافه کند. از دیگر کاربردهای NLP نیز مباحث مربوط به hate speech ها می‌باشد.

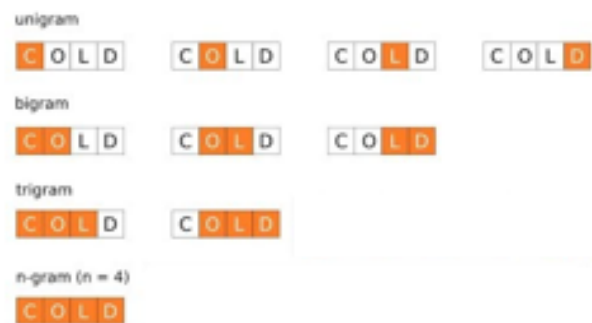
سیستمی که ما در این پروژه موفق به پیاده‌سازی آن شدیم یک سیستم تشخیص شاعر از روی مصراع داده‌شده به آن می‌باشد. این سیستم با آنالیز کلمات موجود در این مصراع قابلیت آن را دارد که با حداکثر دقتی در حدود 82.5% شاعر را تنها از روی یک مصراع تشخیص دهد.

---

## مراحل اجرای پروژه

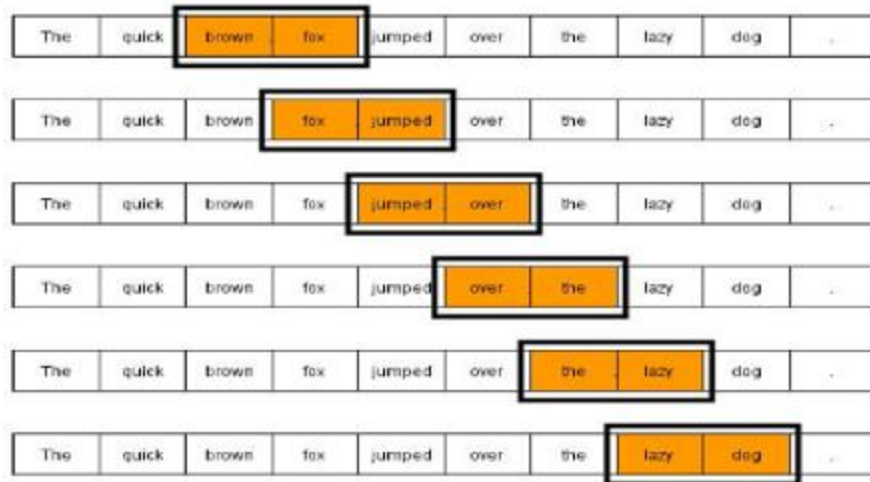
برای پیاده‌سازی این پروژه از مدل‌های unigram, bigram, backoff و همچنین کلاسیفایر naïve Bayes و مالتی کلاسیفایر سیستم soft majority voting و hard majority voting و همچنین از چندین رویکرد ترکیب اطلاعات مدل‌ها استفاده نمودیم که در ادامه به توضیح آن خواهیم پرداخت.

در ابتدا به پیاده سازی از سه مدل زبانی پرداختیم. یک مدل زبانی به صورت یک توزیع احتمالاتی پیاده سازی می شود. این مدل ها شامل تخمین هایی از کلمات با توجه به کلمات ماقبل آن کلمه (به تعدادی مشخص) می شوند. مدل های n-grams در واقع کلمات یا عبارات و .. را به گرام های کوچکتر تقسیم میکنند. این گرام ها برای یک کلمه، حرفهای موجود در کلمه است و برای یک جمله، کلمات موجود در آن و الی آخر. برای نمونه به شکل زیر توجه بفرمایید.



نمونه ای از یک مدل NGRAM که واحدهای آن حروف می باشند

همانطور که در شکل فوق مشاهده می کنید یک کلمه را به گرام های کوچک تقسیم کردیم. این گرام ها ابتدا تکی هستند و در سطر دوم دوتایی و در سطر سوم سه تایی و در سطر آخر نیز چهار تایی می باشند. در بحث مربوط به text classification، اغلب با جملات سر و کار داریم و لذا این جملات را به گرام هایی که در واقع جزئی از جمله یا به عبارت بهتر کلمات موجود در این جملات هستند، تقسیم می نماییم. برای مثال جمله ی زیر را در نظر داشته باشید که به گرام های دوتایی تقسیم شده اند.



نمونه ای از یک مدل BIGRAM که واحدهای آن کلمه می باشند

بحث گرام ها از اینجا اهمیت پیدا می کند که برای تحلیل به جهت دسته بندی متون، تک تک کلمات می توانند اثر گذار باشند و به همین جهت مدل های bigram، trigram و ... نیز حائز اهمیت هستند.

احتمالات شرطی هر یک از کلمات موجود در داده های آموزش به صورت زیر بدست آورده و در یک دیکشنری ذخیره کردیم.

احتمال شرطی در unigram:

$$P(w_i) = \frac{C(w_i)}{N}, \quad C(w_i) = \text{count of occurrence of } w_i,$$

$N$  = total number of words in the training data

سپس احتمالات شرطی bigram را محاسبه نمودیم.

احتمال شرطی در bigram:

$$P(w_i|w_{i-1}) = \frac{C(w_i|w_{i-1})}{C(w_{i-1})}$$

و در نهایت نیز با ترکیب این دو مدل و اضافه نمودن بایاس، یک backoff مدل ایجاد نمودیم.

از این رو برای این بخش یک کلاسِ Documentation در پایتون ایجاد کردیم و تمامی مدل‌های فوق را در آن پیاده نمودیم.

لذا هر یک از این توابع در کلاس به صورت زیر است.

تعریف کلاس:

```
class DocumentClassification:
    def __init__(self, category):
        self.category = category
        self.documents = []
        self.unary_probabilities = {}
        self.frequency = {}
        self.binary_probabilities = {}
        self.binary_probabilities_backed_off = {}
        self.p = None
```

استخراج فرکانس کلمات در یک عبارت:

```
def get_word_frequency(self):
    words_frequency = {'UNK': 0}
    for i in self.documents:
        for j in i.word_list:
            if j not in words_frequency.keys():
                words_frequency['UNK'] = words_frequency.get('UNK') + 1
                words_frequency[j] = 1
            else:
                words_frequency[j] = words_frequency.get(j) + 1
    total_count = sum([v for v in words_frequency.values()])
    self.frequency = words_frequency
    return [words_frequency, total_count]
```

استخراج احتمالات unary:

```
def calculate_unary_probabilities(self):
    frequency, total_count = self.get_word_frequency()
    self.unary_probabilities = {k: v/total_count for k, v in frequency.items()}
    ...
```

استخراج احتمالاتِ باینری:

```
def calculate_binary_probabilities(self):
    pair_wise_words_frequency = {}

    for i in self.documents:
        for j, k in self.pair_wise(i.word_list):
            # print(j, k)
            if f"{j} {k}" in pair_wise_words_frequency:
                pair_wise_words_frequency[f"{j} {k}"] = pair_wise_words_frequency.get(f"{j} {k}") + 1
            else:
                pair_wise_words_frequency[f"{j} {k}"] = 1
    for key, value in pair_wise_words_frequency.items():
        self.binary_probabilities[key] = value/self.frequency.get(key.split()[0])
```

Pair-wise نمودنِ کلمات در یک عبارت برای بدست آوردنِ احتمالاتِ باینری:

```
@staticmethod
def pair_wise(word_list):
    a, b = tee(word_list, 2)
    next(b, None)
    return zip(a, b)
```

بدست آوردنِ احتمالاتِ backoff برای ساخت مدل:

```
def binary_back_off(self):
    y = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8]
    for key, value in self.binary_probabilities.items():
        self.binary_probabilities_backed_off[key] = \
            max([(i * self.unary_probabilities.get(key.split()[1])) + ((1 - i) * value) for i in y])
```

توجه داشته باشید که یکی از مشکلاتی که در پردازشِ زبانِ طبیعی با آن مواجه هستیم وجودِ داده‌هایی است که در داده‌های آموزش وجود ندارند و مدل‌ها احتمالاتِ مربوط به هر یک را برابر با صفر قرار می‌دهد و لذا کلِ عبارت را از عرصه‌ی تشخیص خارج می‌کند که ما این مشکل را با log likelihood حل نمودیم.

$$\begin{aligned} P_{\text{eval}}(\text{text}) &= \prod_{\text{word}} P_{\text{train}}(\text{word}) \\ \log(P_{\text{eval}}(\text{text})) &= \sum_{\text{word}} \log(P_{\text{train}}(\text{word})) \\ \text{Average log likelihood}_{\text{eval}} &= \frac{\sum_{\text{word}} \log(P_{\text{train}}(\text{word}))}{N_{\text{eval}}} \end{aligned}$$

↑  
total number of words in  
evaluation text

حال که فرکانس‌ها و احتمالات کلمات را به صورت یونری و باینری محاسبه کردیم، می‌توانیم عباراتمان را با توجه به احتمالاتی که برای گرام‌ها بدست آورده‌ایم برای ساخت مدل آماده کنیم و این احتمالات را متناسب با نوع خود در یک دیکشنری ذخیره کنیم.

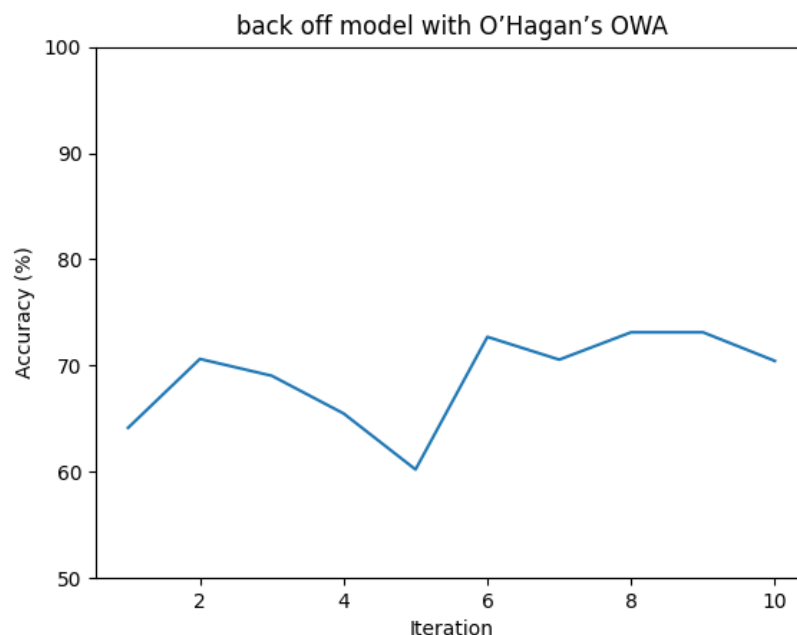
```
def get_unary_probability(self, doc):
    result = 1
    for i in doc.word_list:
        p = self.unary_probabilities.get(i)
        result += math.log10(p) if p is not None else math.log10(0.00000001)
    return result + math.log10(self.p)

def get_binary_probability(self, doc):
    result = 1
    for i, j in self.pair_wise(doc.word_list):
        temp = self.binary_probabilities.get(f"{i} {j}")
        result += math.log10(temp) if temp is not None else math.log10(0.00000001)
    return result + math.log10(self.p)

def get_binary_back_off_probability(self, doc):
    result = 1
    for i, j in self.pair_wise(doc.word_list):
        temp = self.binary_probabilities_backed_off.get(f"{i} {j}")
        result += math.log10(temp) if temp is not None else math.log10(0.00000001)
    return result + math.log10(self.p)
```

حال مدل‌ها را با توجه به توابع اولیه‌ای که در فوق مشاهده می‌کنید تعریف کرده‌ایم که قسمت‌های بیشتر را می‌توانید در فایل‌های موجود در کد مشاهده کنید.

سپس با استفاده از متدهای مختلف دیتا فیوژن ادغام به ترکیب اطلاعات این مدل‌ها نمودیم که نتایج آن به صورت زیر است.



از اولین متدی که برای ترکیب اطلاعات احتمالاتی موجود در مدل‌ها استفاده نمودیم، O'Hagan بود.

صحتی که این متد برای تشخیص متون در اختیار ما قرار داد چیزی در حدود 68.94738214668806 % بود. همانطور که مشاهده می‌نمایید، این نتایج برای ده نمونه‌ی مختلف از داده‌های test و train استخراج شده است. نتایج confusion matrix نیز به صورت زیر است که برای تعدادی از فایل‌های آموزش و تست قرار داده‌ایم.

	ferdowsi	hafez	molavi			ferdowsi	hafez	molavi
ferdowsi	615	29	51	ferdowsi	563	27	47	
hafez	150	291	60	hafez	26	36	16	
molavi	170	52	321	molavi	170	52	321	

	ferdowsi	hafez	molavi			ferdowsi	hafez	molavi
ferdowsi	563	27	47	ferdowsi	516	27	44	
hafez	26	36	16	hafez	99	246	58	
molavi	170	52	321	molavi	152	51	265	

در گام بعدی از متد فیوژن optimistic exponential و البته با اضافه نمودن بایاس به مدل استفاده کردیم که نتایج آن به صورت زیر درآمد.

Orness = 0.8 و لذا وزن‌ها به صورت زیر در آمدند:

$$w_1 = 0.6 \quad w_2 = 0.24 \quad w_3 = 0.16$$

```

C:\Windows\System32\cmd.exe
hafez      23      42      13
molavi     144     70     329

test_file9.txt
Accuracy of program detection : 0.7217806041335453

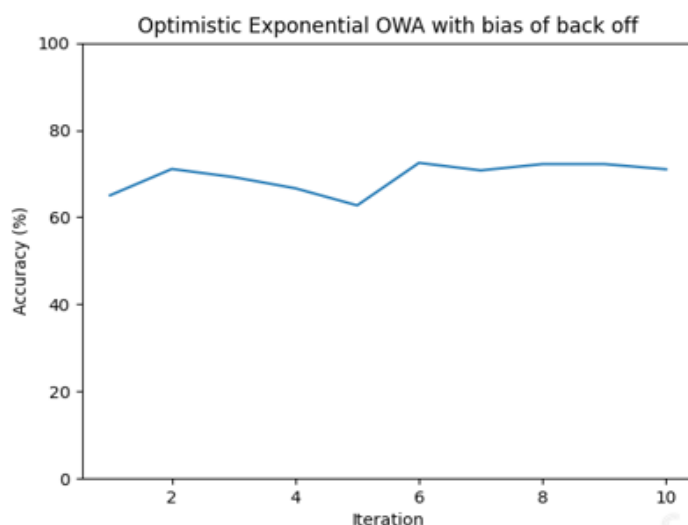
ferdowsi   537     47     53
hafez      23      42     13
molavi     144     70     329

test_file10.txt
Accuracy of program detection : 0.7098765432098766

ferdowsi   495     41     51
hafez      87      264    52
molavi     126     66     276

09.30909386882469
D:\dars\Project\results\owa-optimistic>

```

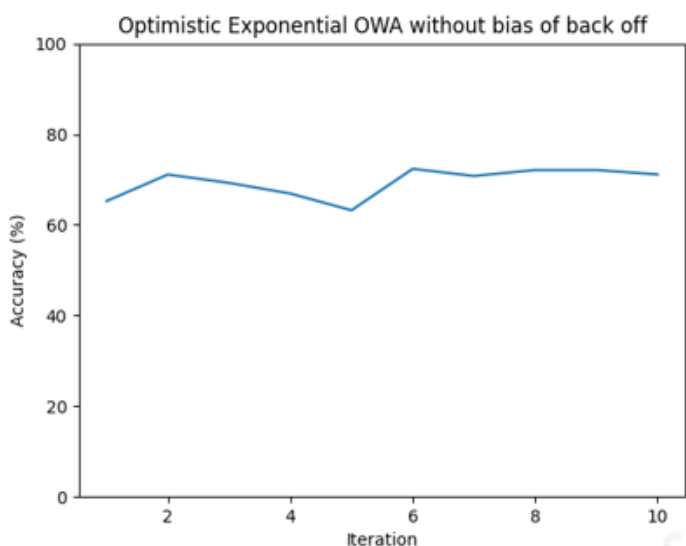




میانگین صحت بدست آمده با  $Orness = 0.8$  برای این مدل نیز برابر با  $69.30909386882469\%$  بدست آمد.

سپس بایاس را حذف نمودیم و تاثیر بایاس در شبکه‌ی گرام ها را بررسی کردیم. نتایج به صورت زیر در آمد

```
C:\Windows\System32\cmd.exe
hafez      21      42      15
molavi     140     69      334
-----
test_file9.txt
accuracy of program detection : 0.7209856915739269
-----
      ferdowsi | hafez | molavi |
-----
ferdowsi      531      52      54
hafez         21      42      15
molavi        140     69      334
-----
test_file10.txt
accuracy of program detection : 0.7112482853223594
-----
      ferdowsi | hafez | molavi |
-----
ferdowsi      489      47      51
hafez         82      266     55
molavi        121     65      282
-----
69.42268641654664
D:\dars\Project\results\owa-optimistic with
```



در این مدت میانگین صحت تا حد بسیار کمی برای دیتاست‌های مشابه افزایش یافت و برابر با  $69.42268641654664\%$  شد.

متد فیوژن بعدی که به backoff اعمال نمودیم Pessimistic Exponential OWA با  $Orness = 0.8$  بود که نتیجه به صورت زیر درآمد.

با توجه به این orness وزن‌ها به صورت زیر درآمدند.

$$w1 = 0.4 \quad w2 = 0.36 \quad w3 = 0.24$$

```

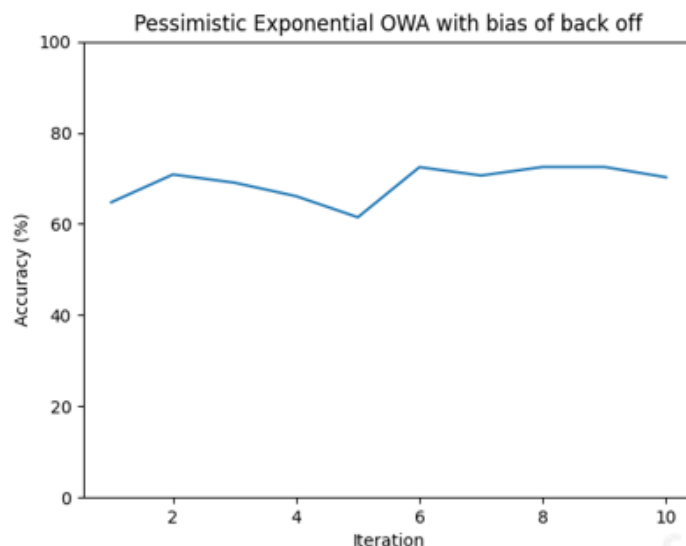
C:\Windows\System32\cmd.exe - python3 main.py
test_file8.txt
Accuracy of program detection : 0.724960254372019
|-----|
| ferdowsi | hafez | molavi | |
|---|---|---|---|
| ferdowsi | 546 | 38 | 53 |
| hafez | 25 | 38 | 15 |
| molavi | 147 | 68 | 328 |
|-----|

test_file9.txt
Accuracy of program detection : 0.724960254372019
|-----|
| ferdowsi | hafez | molavi | |
|---|---|---|---|
| ferdowsi | 546 | 38 | 53 |
| hafez | 25 | 38 | 15 |
| molavi | 147 | 68 | 328 |
|-----|

test_file10.txt
Accuracy of program detection : 0.7023319615912208
|-----|
| ferdowsi | hafez | molavi | |
|---|---|---|---|
| ferdowsi | 501 | 35 | 51 |
| hafez | 94 | 255 | 54 |
| molavi | 134 | 66 | 268 |
|-----|

69.04651213808293

```



میانگین صحت برای این متد در حضور بایاس در شبکه برابر با % 69.04651213808293 بدست آمد. و مجدداً بایاس را از سیستم حذف نمودیم و نتیجه زیر را حاصل کرد که مجدداً با افزایش 0.4 درصد برای صحت همراه بود.

```

C:\Windows\System32\cmd.exe
| molavi | 140 | 69 | 334 |
|-----|

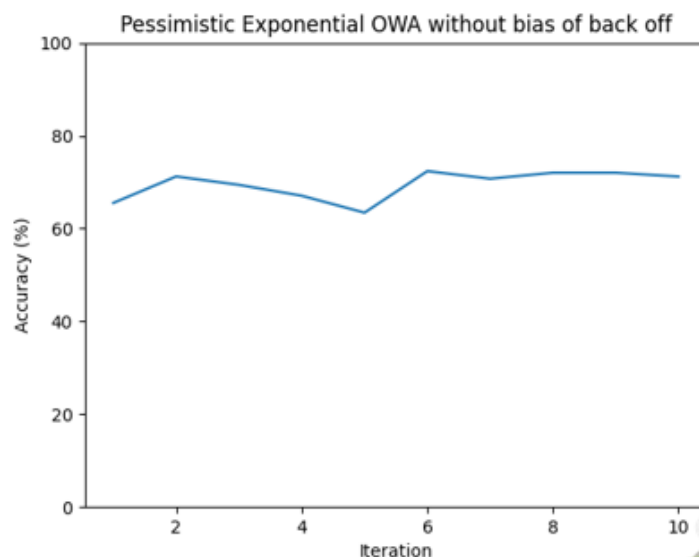
test_file8.txt
Accuracy of program detection : 0.7201907790143084
|-----|
| ferdowsi | hafez | molavi | |
|---|---|---|---|
| ferdowsi | 530 | 53 | 54 |
| hafez | 21 | 42 | 15 |
| molavi | 140 | 69 | 334 |
|-----|

test_file9.txt
Accuracy of program detection : 0.7201907790143084
|-----|
| ferdowsi | hafez | molavi | |
|---|---|---|---|
| ferdowsi | 530 | 53 | 54 |
| hafez | 21 | 42 | 15 |
| molavi | 140 | 69 | 334 |
|-----|

test_file10.txt
Accuracy of program detection : 0.7119341563786008
|-----|
| ferdowsi | hafez | molavi | |
|---|---|---|---|
| ferdowsi | 488 | 48 | 51 |
| hafez | 82 | 266 | 55 |
| molavi | 120 | 64 | 284 |
|-----|

69.49288532940196
D:\dars\... \Project\results\owa-pessimistic without

```



که صحت بدست آمده برای این مدل برابر با % 69.49288532940196 شد.

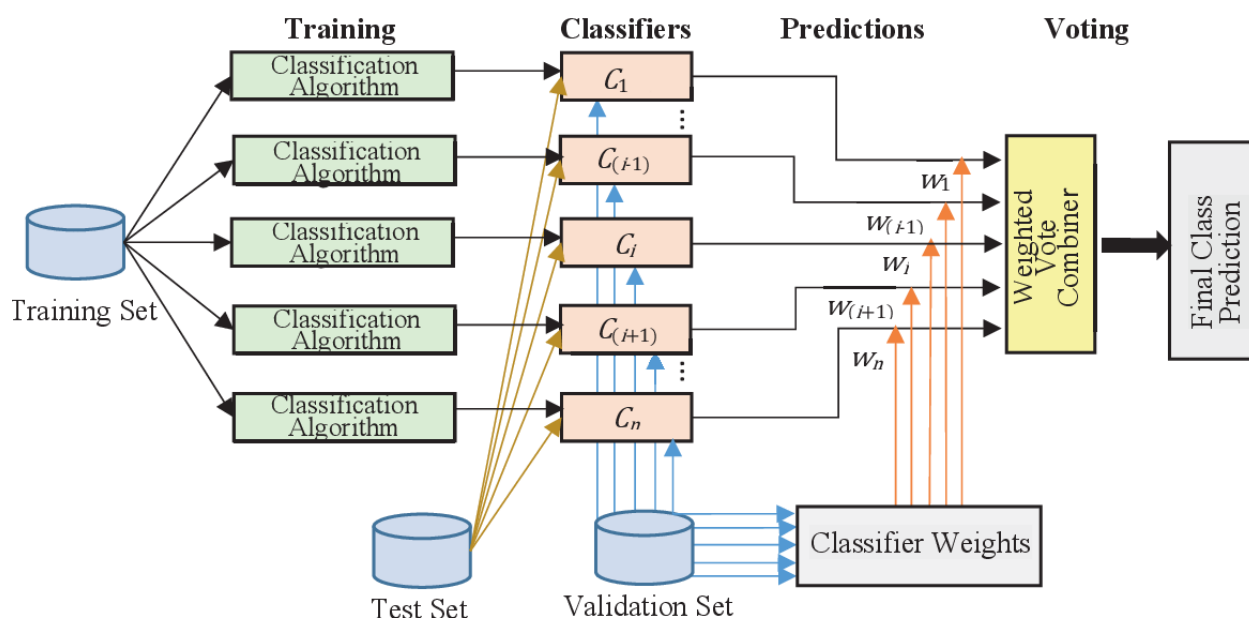
لذا به طور کلی صحت بدست آمده برای متدهای فیوژنِ فوقِ باحضورِ یا عدمِ حضورِ بایاس به شکل زیر شد.

Fusion Method	Average accuracy
O'Hagan's OWA	68.94738214668806 %
Optimistic Exponential OWA with bias of back off	69.30909386882469 %
Optimistic Exponential OWA without bias of back off	69.42268641654664 %
Pessimistic Exponential OWA with bias of back off	69.04651213808293 %
Pessimistic Exponential OWA without bias of back off	69.49288532940196 %

## Multi-Classifer Fusion System

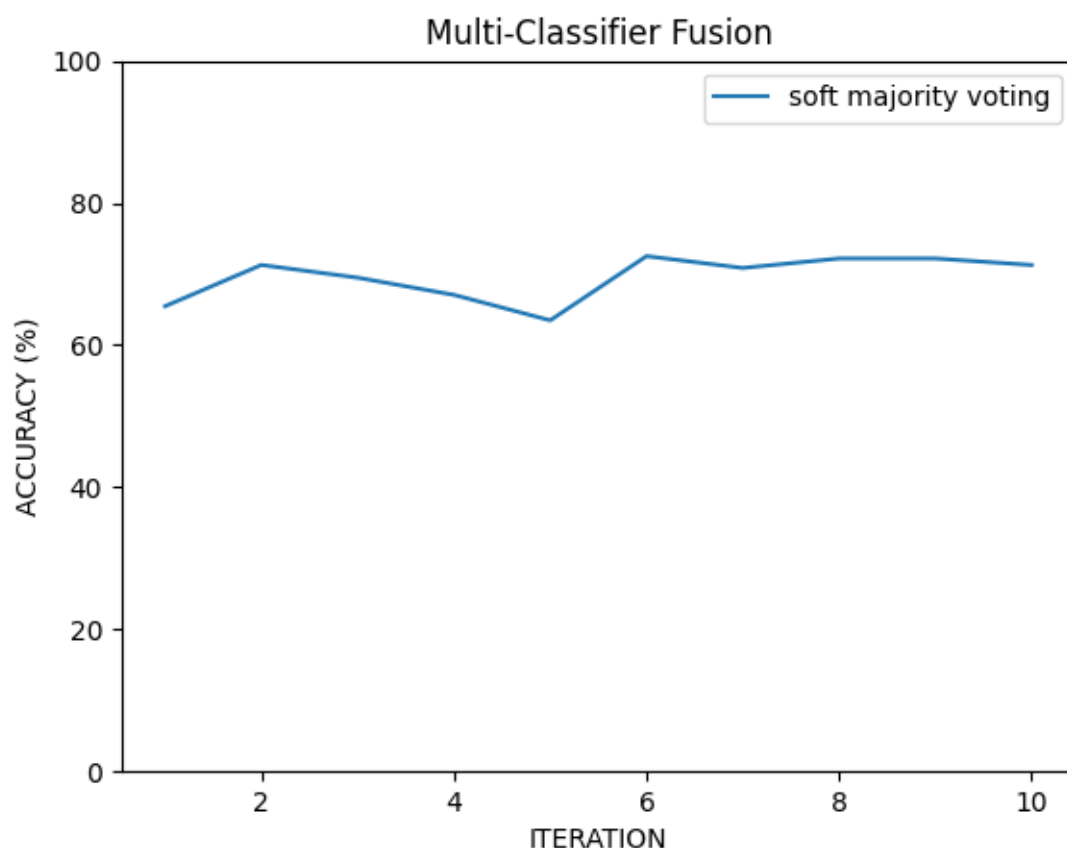
در این گام از soft majority voting و hard majority voting استفاده نمودیم که صحت بالاتری را

تسب به متدهای فیوژنی که در مرحله‌ی قبل بررسی نمودیم، ارائه داد.



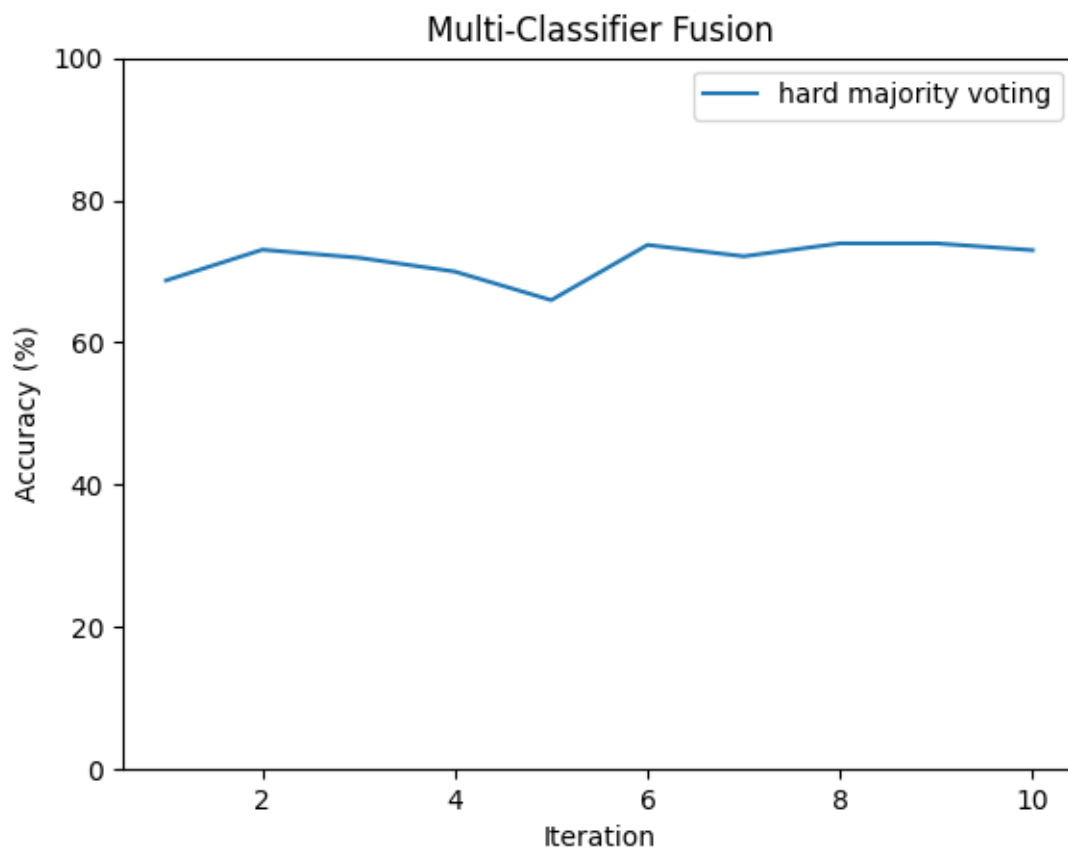
خروجی حاصل برای soft majority voting به صورت زیر است.

### Soft Majority Voting



صحت بدست آمده برای این الگوریتم مالتی کلاسیفایر برابر با  $69.57154727221628\%$  می باشد که به طور میانگین حدود  $0.3$  درصد از میانگین صحت های متدهای فیوژن بیشتر است و نشان از کارایی این روش در هنگام استفاده از چندین مدل را ارائه می دهد. برای soft majority voting احتمالات خروجی از مدل ها را برای تخصیص به یک کلاس خاص با یکدیگر جمع زدیم و هر کدام که بیشترین مقدار تجمیعی را داشت به عنوان برچسب مربوط به آن ذخیره نمودیم.

## Hard Majority Voting

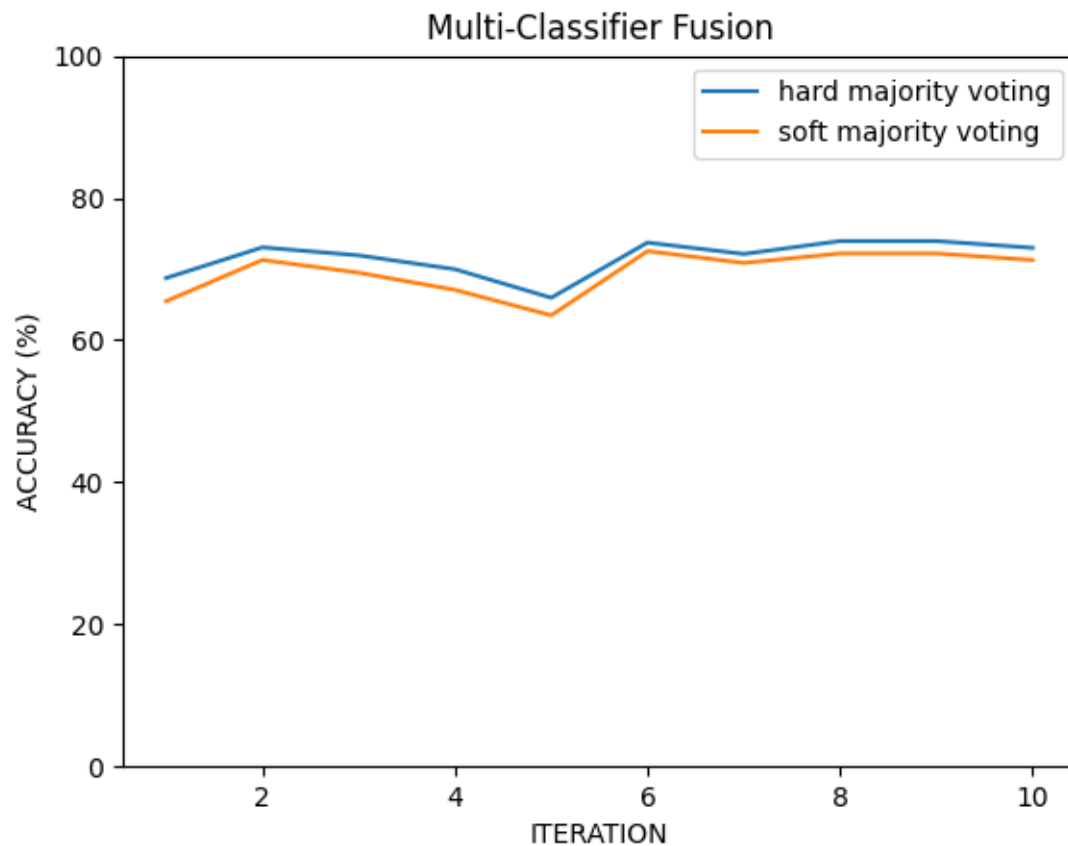


در این روش که خروجی‌های مدل‌ها را نه به صورت احتمالاتی بلکه به صورت تصمیم برای تشخیص برچسب گرفتیم، به بیشترین صحت تا این لحظه رسیدیم.

میانگین صحت این مدل برای مجموعه دیتاست‌های فوق برابر با  $71.62333879424502\%$  و با انحراف معیار  $2.635618776987621$  رسیدیم.

دلیل این امر می‌تواند قدرت بالای majority باشد زمانیکه دو یا تعداد بیشتری مدل تصمیم به برچسب زدن یکسان به یک دیتا به صورت همزمان باشند. طبیعتاً هر چه آرا بیشتر باشد و تصمیم مدل‌های بیشتری برای

اختصاص دادنِ برچسبِ خاصی به داده باشد، لذا آن داده با احتمالِ بالاتری متعلق به همان کلاسِ برچسب زده شده است.



### Naïve Bayes Classifier with Bigram and Unigram model

همانطوری که می‌دانیم قانونِ بیز روشی برای دسته‌بندی بر پایه‌ی احتمال وقوع یا عدم وقوع یک پدیده است که کاربردش را در کنار مدل‌های ترکیب شده با owa و همچنین hard majority voting و soft majority voting مشاهده نمودیم. در واقع در این گام به استفاده از مدل‌های unigram و bigram به صورتِ جداگانه و

با بهره‌گیری از کلاسیفایر naïve Bayes می‌باشد، می‌پردازیم. در واقع ما در این مرحله با استفاده از قانون بیز، برای این مدل‌ها به صورت جداگانه، محتمل‌ترین زبان را انتخاب می‌کنیم.

$$l^* = \operatorname{argmax} P(l|c_{1:N})$$

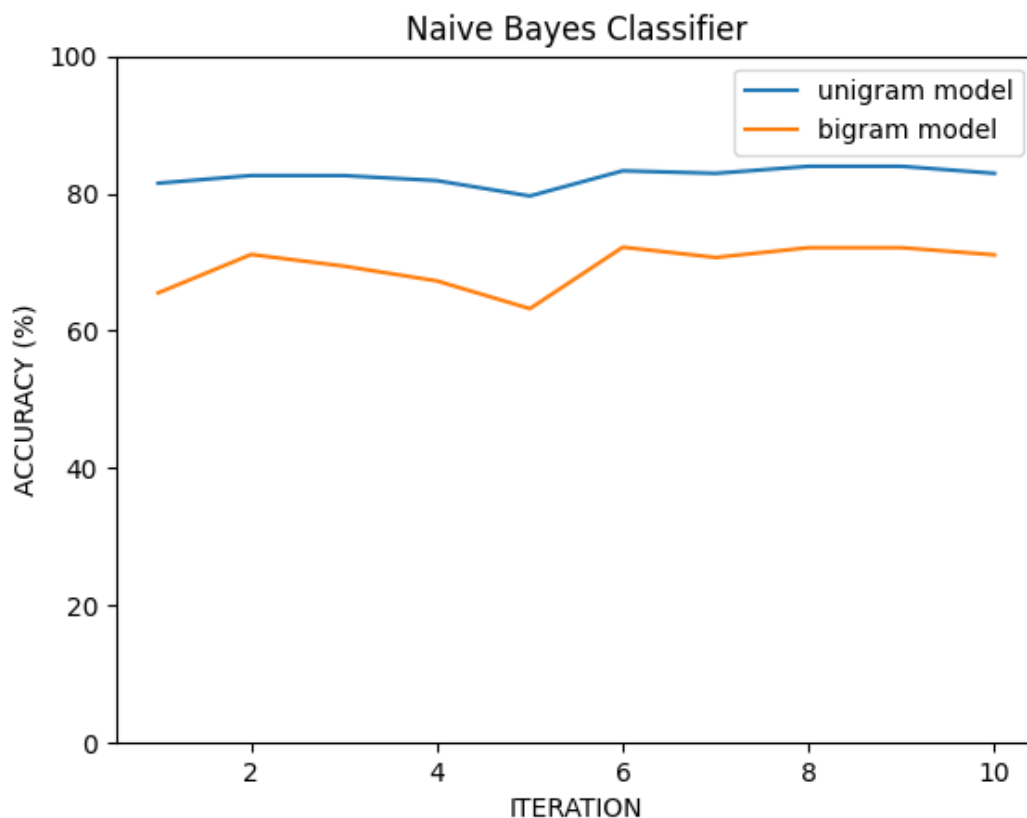
$$= \operatorname{argmax} P(l)P(c_{1:N}|l)$$

$$= \operatorname{argmax} P(l) \prod_{i=1}^N P(c_i|c_{i-2:i-1}, l)$$

$$P(l|c_{1:N}) = \frac{P(l)P(c_{1:N}|l)}{P(c_{1:N})}$$

حال نتایج این دو مدل را به صورت جداگانه در یک شکل ارائه می‌کنیم.

**خروجی مدل‌های unigram و bigram:**



بررسی میانگین صحت و انحراف معیار:

	Average accuracy of accuracy curve	Standard deviation of accuracy curve
<b>unigram</b>	<b>82.51816581708587</b>	<b>1.2860424633208134</b>
<b>bigram</b>	<b>69.4594504770292</b>	<b>3.11665339550396</b>

همانطور که مشاهده می‌کنید بهترین نتیجه‌ی ممکن را unigram با صحتی نزدیک به 82.51% در اختیار ما قرار داده است.

حال تمامی نتایج را در کنار یکدیگر قرار می‌دهیم و یک ارزیابی کلی انجام داده و سپس به سراغ الگوریتم MapReduce در سیستم HDFS برای پردازش‌های زبان طبیعی خواهیم رفت.

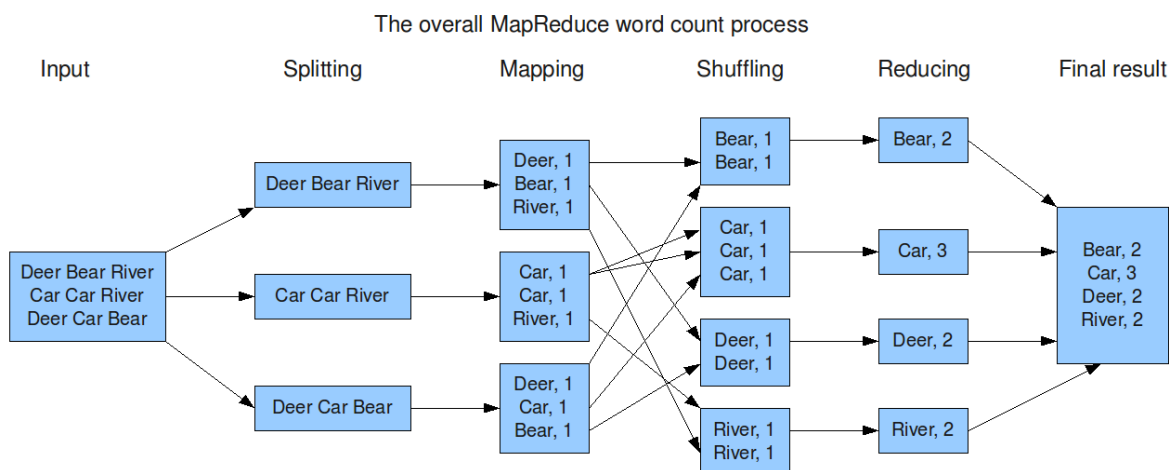
	Average accuracy of accuracy curve	Standard deviation of accuracy curve
<b>Backoff model with O'Hagan's OWA</b>	<b>68.947 %</b>	<b>4.320</b>
<b>Backoff model with optimistic exponential OWA</b>	<b>69.422 %</b>	<b>3.210</b>
<b>Backoff model with optimistic exponential OWA and bias</b>	<b>69.309 %</b>	<b>3.396</b>
<b>Backoff model with pessimistic exponential OWA</b>	<b>69.492 %</b>	<b>3.111</b>
<b>Backoff model with pessimistic exponential OWA and bias</b>	<b>69.0465 %</b>	<b>3.762</b>
<b>Hard majority voting</b>	<b>71.623 %</b>	<b>2.635</b>
<b>Soft majority voting</b>	<b>69.571 %</b>	<b>3.160</b>
<b>Unigram model with naïve Bayes</b>	<b>82.518 %</b>	<b>1.286</b>
<b>Bigram model with naïve Bayes</b>	<b>69.459 %</b>	<b>3.116</b>



همانطور که از نمودارِ فوق مشاهده می‌شود بهترین عملکرد را unigram و سپس hard majority voting و سپس soft majority voting، سپس Backoff model with pessimistic exponential OWA و سپس bigram و بعد از آن نیز متدهای دیگر قرار می‌گیرند. دلیل این امر را می‌توان در دیتاست جستجو کرد که به صورت شعرهایی هستند که ممکن است هر شاعر از کلماتی منحصر بفرد استفاده نماید که شاعر دیگری از این کلمات بهره نگیرد. همچنین تعداد تکرارِ کلماتی خاص توسطِ شاعران نیز از این قاعده مستثنی نیست. دلایلِ مربوط به majority voting و عموماً خوب آن را نیز در بخشِ خود بررسی نمودیم. و همچنین در اینجا بین متدهای مختلفِ OWA، نمایی بدبینانه بهترین عملکرد را از خود نشان داده که از آن می‌توان استنباط کرد که از آنجا که در این روش، محتاطانه وزن‌ها انتخاب می‌شوند، لذا نتایجِ مدل‌ها هر یک برای تشخیص دارای اطلاعاتِ کافی برای تصمیم‌گیری می‌باشند و به همین دلیل است که وزن‌ها می‌توانند در یک محدوده‌ی نزدیک به هم برای ضریبِ هر یک از مدل‌ها قرار گیرند.

### استفاده از سیستم HDFS برای انجام پردازش‌های زبانِ طبیعی

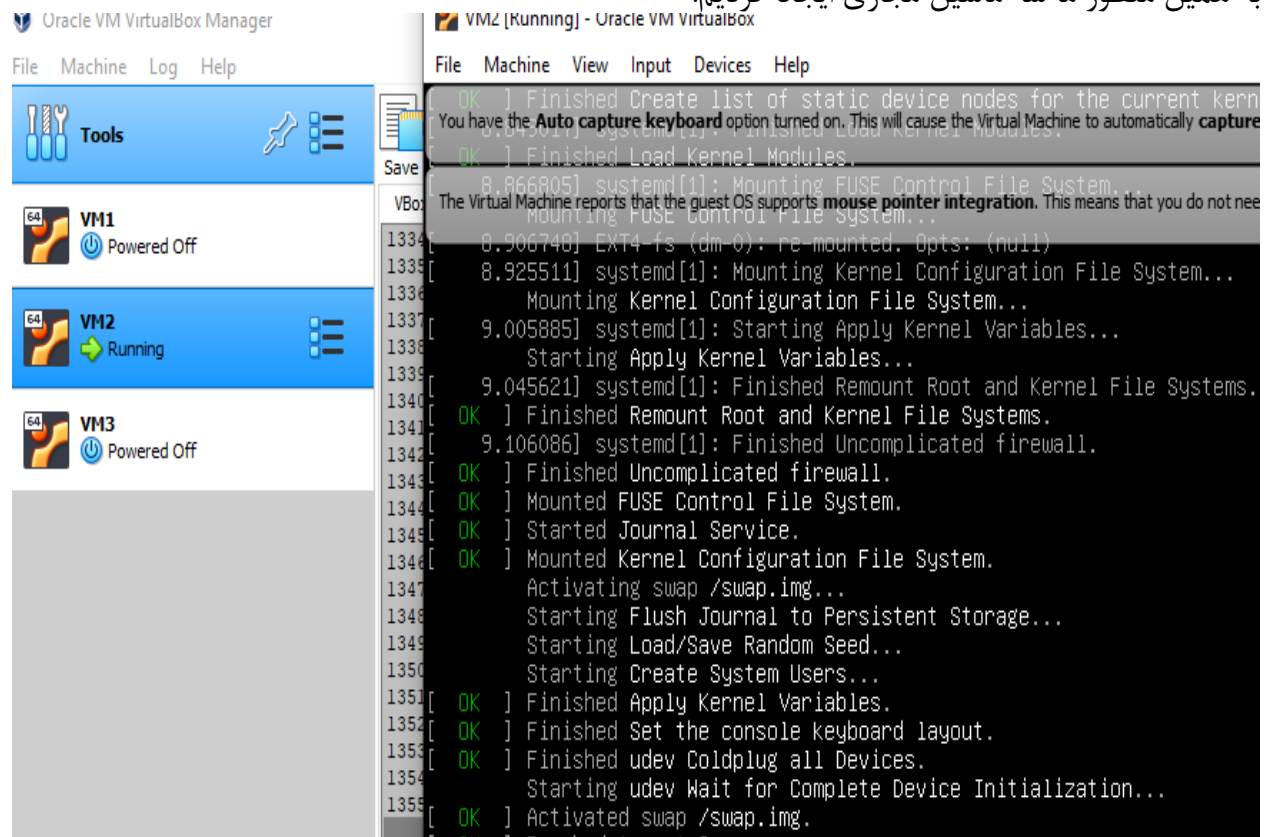
برای انجام این پردازش‌ها از ابزارِ هادوپ استفاده نمودیم. الگوریتمِ اجرا شده بر روی این بخش الگوریتمِ MapReduce است. شمای کلی این الگوریتم را در زیر مشاهده می‌کنید.



این سیستم عملیات ها را بدین صورت که مشاهده می کنید انجام می دهد. مثلا برای آموزش مدل از طریق دیتاها باید جملات را به گرام‌هایی که کلمات هستند تقسیم کرده و سپس کلمات را شمارش کرده و فرکانس آنها را بیابیم. این الگوریتم در ابتدا Splitting را انجام می دهد و جملات موجود در Dataset را متناسب با خواسته‌ی ما تفکیک کرده و به گروه‌هایی تقسیم می کند و سپس جملات را با استفاده از عملیات splitting به گرام‌ها تفکیک می نماید. حال در عملیات Mapping به ازای حضور هر کلمه، یک واحد به آن اختصاص می دهد. سپس در عملیات shuffling کلمات مشابه را در کنار یکدیگر قرار داده و نهایتا در عملیات Reducing آنها را با یکدیگر جمع زده و تعداد تکرار کلمات را با سرعت بالایی می یابد. حال که تعداد تکرار هر کلمه را یافت، آنها را به عنوان خروجی بر می گرداند. این نحوه‌ی عمل سیستم فوق برای هر تسکی از جمله پردازش های NLP است که مثال آن را ذکر کردیم.

لازم است که ماشین های دیگری در خدمت این سیستم باشند و سیستم بتواند تسک ها را بین آنها تقسیم کرده و سپس خروجی های را گرفته و ترکیب کند.

به همین منظور ما سه ماشین مجازی ایجاد کردیم.



و پس از نصب این سیستم و قرار دادن الگوریتم لازم برای ترکیب در دایرکتوری آن، اقدام به start کردن yarn میکنیم تا نودها شروع به کار کنند.

```
Select hadoop@vm1: ~
hadoop@vm1:~$ sudo chown -R 755 hdfs
hadoop@vm1:~$ sudo mkdir -p /home/hadoop/hdfs/{namenode,datanode}
hadoop@vm1:~$ ls hdfs/*
hdfs/datanode:

hdfs/namenode:
hadoop@vm1:~$ sudo chown -R hadoop:hadoop /home/hadoop/hdfs
hadoop@vm1:~$ start-dfs.sh
Starting namenodes on [vm1]
vm1: namenode is running as process 19177. Stop it first.
Starting datanodes
Starting secondary namenodes [vm1]
vm1: secondarynamenode is running as process 19516. Stop it first.
hadoop@vm1:~$ jps
20357 Jps
19177 NameNode
20057 DataNode
19516 SecondaryNameNode
hadoop@vm1:~$ start-yarn.sh
Starting resourcemanager
Starting nodemanagers
hadoop@vm1:~$ jps
20736 Jps
19177 NameNode
20057 DataNode
20588 NodeManager
19516 SecondaryNameNode
20447 ResourceManager
hadoop@vm1:~$ stop-dfs.sh
Stopping namenodes on [vm1]
Stopping datanodes
Stopping secondary namenodes [vm1]
hadoop@vm1:~$ stop-yarn.sh
Stopping nodemanagers
Stopping resourcemanager
hadoop@vm1:~$
```

همانطور که مشاهده می کنید تمامی نودها در حال فعالیت و در خدمت ماشین پردازنده ی NLP هستند. این ماشین ها در هر لحظه می توانند وضعیت خود را به اشتراک گذارند و سیستم نیز در این صورت هیچ اضافه باری به ماشینی نخواهد داد.

همچنین در صورت fail شدن یک تسک مثل یافتن تعداد تکرار کلمه در دیتاست، ماشین اصلی این تسک را مجدداً به یک ماشین دیگر می دهد.

در شکل زیر می‌توانید دسترسی به هدوپ را از طریق ویندوز مشاهده کنید. این دسترسی پس از **start** خوردن و شروع فعالیت‌ها برقرار می‌شود و همانطور که مشاهده می‌کنید سیستم به درستی پردازش‌ها را انجام داده و خروجی‌های خود با فراخوانی **url** متناظر مربوط به ماشین اصلی نمایش می‌دهد.

**Hadoop Overview** 'vm1:9000' (active)

Started:	Tue Jan 05 18:40:21 +0330 2021
Version:	3.2.1, rb3cbb0467e22ea829b3808f4b7b01d07e0bf3842
Compiled:	Tue Sep 10 20:26:00 +0430 2019 by rohitsharmaks from branch-3.2.1
Cluster ID:	C10-406326d5-9657-409e-b00e-109261438c0e
Block Pool ID:	BP-10167620-192.168.1.4-1609859349803

**Summary**

Security is off

Safe mode is ON. Resources are low on NN. Please add or free up more resources then turn off safe mode manually. NOTE: If you turn off safe mode before adding resources, the NN will immediately return to safe mode. Use "hdfs dfsadmin -safemode leave" to turn safe mode off.

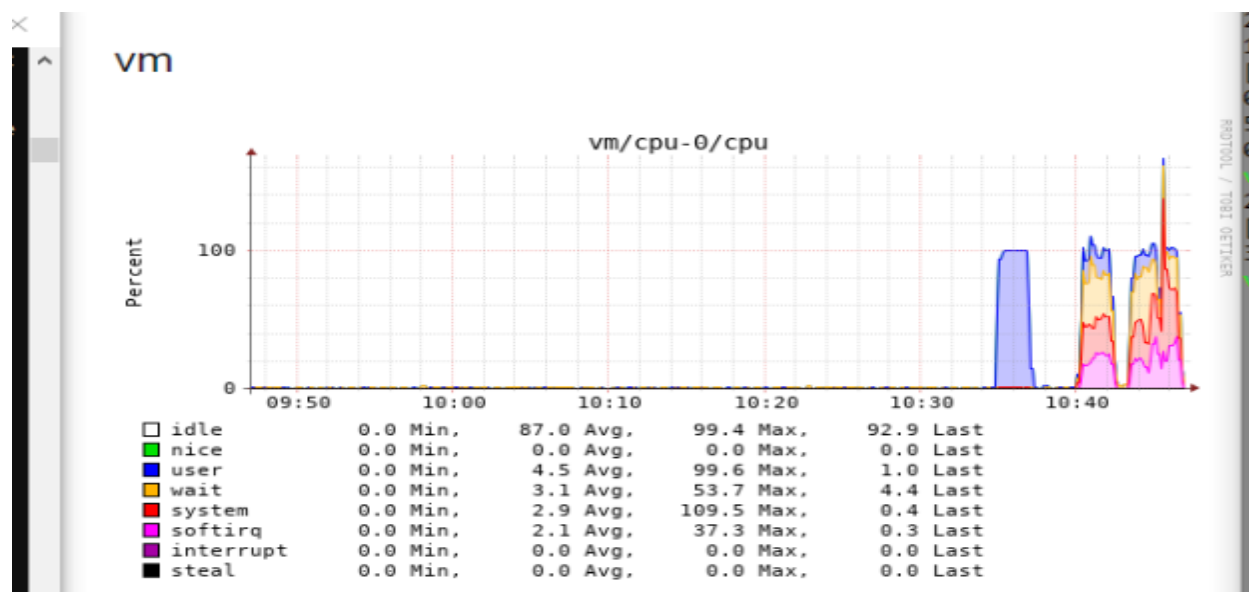
1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

Heap Memory used 20.06 MB of 29.32 MB Heap Memory. Max Heap Memory is 237.81 MB.

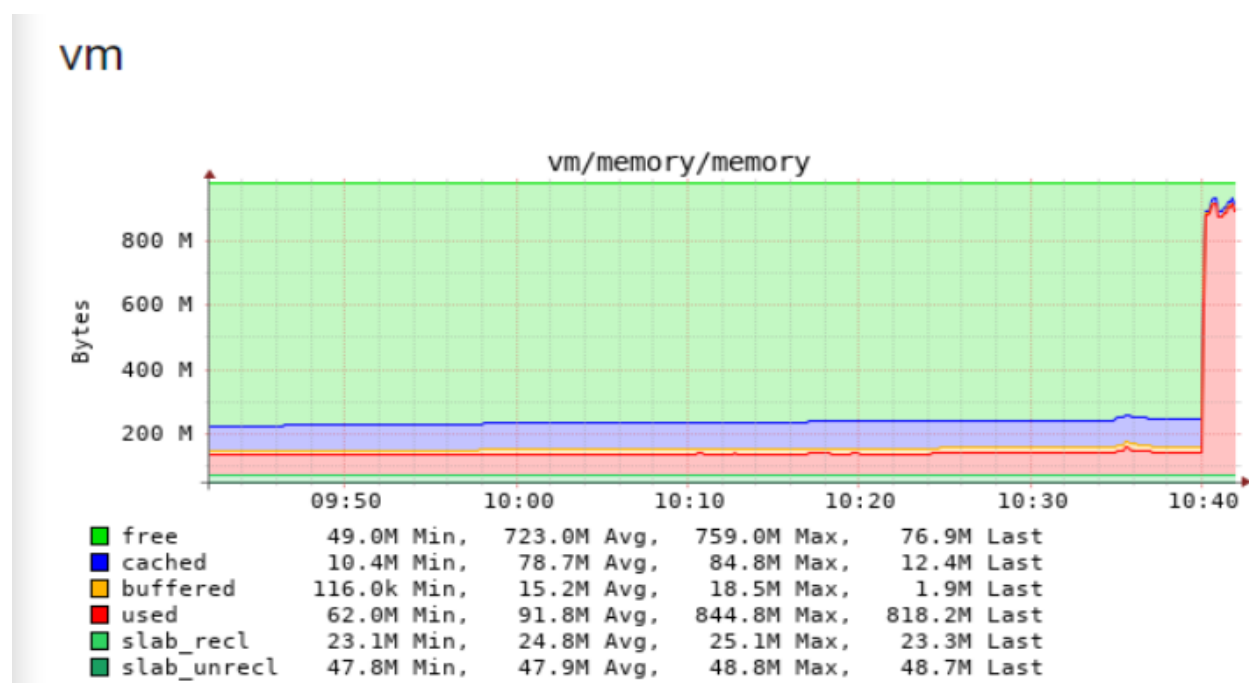
Non Heap Memory used 47.02 MB of 48.19 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

در نهایت طبق مقاله‌ای که ضمیمه نمودیم و از این سیستم نیز برای پردازش‌های **NLP** استفاده نموده، ما نیز عملیات‌ها را بر روی آن انجام دادیم و نتیجه‌ی زیر برای پردازش‌های **cpu** حاصل شد.

وضعیتِ cpu :



همچنین وضعیتِ memory نیز به صورت زیر در آمده است.



با توجه به میزان فضای کم و cpu کمی که vm ها در اختیار دارند، وضعیت خوبی را مشاهده می کنیم و به نظر می رسد این سیستم با الگوریتم های ذکر شده می تواند یک سیستم کارا برای پردازش های سنگین NLP باشد.

در آخر نیز نمای کلی سیستم تشخیص شاعر از روی شعر را قرار داده‌ایم. همانطور که مشاهده می‌کنید این سیستم برای ID های مختلف که در واقع سطرهای فایل ورودی هستند، شاعر را تشخیص می‌دهد. در واقع ستون وسط، شاعر اصلی این شعر و ستون آخر شاعر پیش‌بینی شده نمایش داده می‌شود که با مراجعه به ستون اول و تطابق دادن ID با سطر موجود در فایل، مصراعی را که پیش‌بینی بر روی آن انجام گرفته است را با حداکثر صحت 82.5% (دقت مربوط به unigram model with naïve Bayes) مشاهده خواهید کرد.

The screenshot shows a Windows command prompt window titled 'C:\Windows\System32\cmd.exe' displaying the output of a program. The output lists 20 lines of text, each with a line number and a poet's name. The last line shows the accuracy of the program detection as 0.8259515570934256. Below this, a table shows the counts for each poet across the three columns.

	ferdowsi	hafez	molavi
ferdowsi	910	79	61
hafez	25	585	109
molavi	71	158	892

Overlaid on the command prompt is a Notepad window titled 'test\_file.txt - Notepad' containing a Persian poem by Hafez:

```

3  هین گمان بد برای ذولباب *
3  با خود ، والله اعلم بالصواب
3  جهد ببقمیر به فتح مکه هم
3  کد بود در حب دنیا متهم ؟
3  آنکه او از مخزن هفت آسمان
3  چشم و دل بر بست روز امتحان
3  از بی نظاره اش حور چنان
3  کرده بر آفاق هر هفت آسمان
3  قدسیان افتاده بر خاک رهش *
3  صد جو یوسف اوقاده در جهش
3  خویشتن آراسته از بهر او
3  خود و را پروای غیر دوست کو ؟
3  آنچنان پر گشته از اجلال حق
3  کاندرا او هم ره نیابد آل حق
3  لا یسج فینا نبی مرسل
3  و الملک و الروح ایضا فاعقلوا
3  گفت ما زاعیم ، همچون زاغ نی
3  مست صباغیم ، مست باغ نی
3  حدنکه مک نماه ، افلاک ه عتفا .
  
```

### پیشنهادهات:

در این نرم‌افزار حداکثر دقتی که در بین تمامی الگوریتم‌ها پیاده‌سازی شده داستیم چیزی در حدود 82.5% بود. پیشنهاد من این است که از الگوریتم‌های learning استفاده کنیم و فیوژن را با وزن‌هایی به مدل‌ها اعمال کنیم، که جمله‌ی ورودی را به بهترین نحو تشخیص دهد. این وزن‌ها را باید از روی احتمال کلمات و نیز گرام‌های موجود در جمله‌ای را که به عنوان ورودی به سیستم داده می‌شود، بیابیم. همچنین از جمله برنامه‌های من بررسی

SVM نیز علاوه بر متدهای فوق بود و پیشنهاد می‌کند الگوریتم SVM نیز در کنار مقایساتی که انجام شده، مورد بررسی قرار گیرد.

## مقالات

نرم‌افزار فوق در واقع پیاده‌سازی چهار رفرنس اصلی به همراه اضافه نمودن متدهای فیوژن به الگوریتم‌های این رفرنس‌ها می‌باشد.

- ✓ **A Comparative Study of Bing Web N-gram Language Models for Web Search and Natural Language Processing, Jianfeng Gao, Patrick Nguyen, Xiaolong Li, Chris Thrasher, Mu Li, Kuansan Wang**

این مقاله به پیاده‌سازی مدل‌های پایه‌ای ngram برای کاربردهای NLP پرداخته بود که ما دو مدل از مدل‌هایش را پیاده‌سازی نمودیم. و البته کلاسفرهای و مالتی کلاسفرهای جدیدی را اضافه کردیم.

- ✓ **Comparing Approaches to Convert Recurrent Neural Networks into Backoff Language Models For Efficient Decoding, Heike Adel, Katrin Kirchhoff, Ngoc Thang Vu, Dominic Telaar, Tanja Schultz**

این مقاله به پیاده‌سازی مدل backoff پرداخته بود و ما نیز این مدل را پیاده کردیم و الگوریتم‌ها ترکیب اطلاعات را روی آن اعمال نمودیم.

- ✓ **Managed N-gram Language Model Based on Hadoop Framework and a Hbase Tables, Tahani Mahmoud Allam, Hatem M. Abdulkader, Alsayed Abdelhameed Sallam**

این مقاله پردازش‌های NLP را بر روی بستر هادوپ انجام داده بود و کارایی هادوپ را به عنوان محیطی برای مدیریت big data مورد ارزیابی قرار داد که در پیاده‌سازی ما نیز این بخش مشهود است.

علاوه بر مقالات فوق از مقالات مربوط به پیاده‌سازی HDFS و الگوریتم MapReduce نیز استفاده کردیم که در مراجع موجود در درفت مقاله ذکر نموده‌ایم.