

به نام خدا

پروژه سوم

این شعر از کیست؟

دریا زارع مذهبيه

۹۷۳۱۰۸۶



فهرست

۲ مراحل الگوریتم:
۲ توابع:
۴ تحلیل:

مراحل الگوریتم:

ابتدا با توجه به داده train وردی Unygram و bygram را برای هر یک از شاعر ها می‌شازیم (که یم دیکشنری با کلید "لغت" و مقدار "تعداد تکرار" می‌باشد).

باید داده های unigram را تمیز کنیم بنابراین لغاتی که کمتر از ۲ باز تکرار شده اند را از لغتنامه حذف میکنیم.

سپس با توجه به الگوریتم backoff احتمال وقوع هر کلمه را برای داده تست حساب میکنیم و در نهایت شاعر با احتمال بیشتر که یک شاعر شعری را سروده به عنوان نتیجه برمیگردانیم.

توابع:

```
def backoff(w1,w2,unigram,bygram,lambd1,lambd2,lambd3,beta):
    resault=0
    ww=w1+' '+w2

    if ww in bygram and w1 in unigram:

        resault = lambd1*bygram.get(ww)/unigram.get(w1)
    if w1 in unigram:
        resault+=lambd2*unigram.get(w1)/len(unigram.keys())
    resault+=lambd3*beta
    return resault
```

تابع ای که برای هر کلمه محاسبه میکند احتمال وقوعش را در unigram و bygram ورودی.

```
def howMuchIn(poem,unigram,bygram,lambd1,lambd2,lambd3,beta):
    poem=poem.split()
    prob=0

    for i in range(len(poem)):
        if i==len(poem)-1:
            prob+=backoff(poem[i], '<>',unigram,bygram,lambd1,lambd2,lambd3,beta)
        else:
            prob+=backoff(poem[i],poem[i+1],unigram,bygram,lambd1,lambd2,lambd3,beta)

    return prob
```

محاسبه میکند احتمالی که یک شعر برای شاعری که unygram و bygram آن را ورودی گرفته.

```
def whichOne(poem, lambda1, lambda2, lambda3, beta):
    probF = howMuchIn(poem, ferdowsi, ferdowsi2, lambda1[0], lambda2[0], lambda3[0], beta[0])
    probH = howMuchIn(poem, hafez, hafez2, lambda1[1], lambda2[1], lambda3[1], beta[1])
    probM = howMuchIn(poem, molavi, molavi2, lambda1[2], lambda2[2], lambda3[2], beta[2])
    m = max(probF, probH, probM)
    #print(probF, probH, probM)
    if probF == m:
        return '1'
    if probH == m:
        return '2'
    return '3'
```

با توه به احتمالات حساب شده تصمیم میگیرد کدام شاعر شعر را سروده است.

```
def maxLambda(lambda1, lambda2, lambda3, beta):
    val = [[0, 0, 0], [0, 0, 0], [0, 0, 0]]
    for i in range(len(poems)):
        temp = whichOne(poems[i], lambda1, lambda2, lambda3, beta)

        if poets[i] == temp:
            val[1][int(temp)-1] += 1
        else:
            val[2][int(temp)-1] += 1
        val[0][int(poets[i])-1] += 1

    return val
```

این تابع برای پیدا کردن لامبدا مناسب بکار می‌رود.

به این صورت که برای لامبدا ورودی محاسبه میکند از هر شاعر چند تا را درست تشخیص داده و چند تا را غلط در نتیجه میتوان دقت هر سری لاندا را برای هر شاعر مشاهاذا کرد و مقادیر را بهبود داد.

تحلیل:

با توجه به اندازه لغت نامه ها داده ها بالانس نیستند و دیتا حافظ بیشتر از بقیه است و کمترین را فردوسی دارد. در نتیجه بتا را باید برای حافظ مقدار خیلی کمی در نظر بگیریم و برای فردوسی بیشتر باشد.

در کل مقدار لاندا ای که در بتا ضرب میشود (در صورت پروژه لاندا ۳) مقدار خیلی کمی باید باشد زیرا دقت را کاهش میدهد. اما اگر تعداد غلط هایی که یک شاعر داشت زیاد بود باید ضریب بایگرام افزایش یابد.

ادامه تحلیل ها در کد.