

به نام خدا



دانشگاه صنعتی امیرکبیر  
دانشکده مهندسی کامپیوتر

استاد درس: دکتر نیک‌آبادی

مرداد ۱۴۰۱

درس مدل‌های گرافیکی احتمالی

گزارش پروژه

سروش مهدی  
شماره دانشجویی: ۹۹۱۳۱۰۵۰

## فهرست مطالب

۱	مقدمه	۱
۲	الگوریتم forward backward	۲
۲	توابع محاسبه ماتریس های $\pi$ , $M$	۱.۲
۲	تابع emission	۲.۲
۲	توابع alpha recursion, beta recursion, gamma	۳.۲
۲	نتایج	۴.۲
۳	الگوریتم EM	۳
۳	pair marginals	۱.۳
۳	تابع EM	۲.۳
۳	نتایج	۳.۳
۴	تحلیل	۴.۳
۵	الگوریتم viterbi	۴
۵	متغیر های زمینه ای	۵
۶	ویژگی پنجم	۶
۶	حالات دیگر	۷
۷	تست	۸

## ۱ مقدمه

در این پروژه فرض‌های مختلفی در نظر گرفته شده که در مورد همه بخش‌ها صادق است مگر آن که صریحا چیز دیگری گفته شود.

اولا به دلیل کمبود داده‌ها نسبت به حالات مختلف ویژگی‌ها (صرفا ویژگی‌های یک تا پنج منظور است) تمامی ویژگی‌ها به شرط متغیر هدف، دو به دو مستقل در نظر گرفته شده‌اند. نتیجه این مستقل بودن این است که احتمال توام ویژگی‌ها به شرط متغیر هدف برابر با ضرب احتمال‌های تکی ویژگی‌ها به شرط متغیر هدف می‌باشد. از این فرض در محاسبه ماتریس emission استفاده شده.

دوما با توجه به مستقل بودن ویژگی‌ها برای حل کردن مشکل مقادیر null ویژگی‌ها از رابطه زیر استفاده شده :

$$P(f_1, f_2, f_3, f_4, f_5=\text{null} \mid \text{target}) = P(f_1, f_2, f_3, f_4 \mid \text{target})$$

یعنی با توجه به مستقل بودن ویژگی‌ها میدانیم احتمال توام برابر با ضرب احتمال‌های تکی هست. از طرفی اگر مقدار یک ویژگی null باشد مانند این است که برای این دسته از ویژگی‌ها فقط احتمال توام بقیه ویژگی‌هایی که مقدار آن‌ها مشخص هست را حساب کنیم. یعنی مقدار ویژگی null را نمیدانیم و از روی بقیه ویژگی‌ها ماتریس emission را محاسبه میکنیم.

در رابطه با ویژگی‌های زمینه‌ای نیز آن‌ها را از ویژگی‌های معمولی مستقل در نظر گرفته ایم که با توجه به داده‌ها فرض منطقی‌ای میباشد چون به ازای هر سمپل، به ازای ویژگی‌های مختلف ویژگی‌های زمینه‌ای یکسان داریم. همچنین به دلیل کمبود داده‌ها ویژگی‌های زمینه‌ای را نیز از یکدیگر مستقل در نظر گرفته ایم و احتمال توام آن‌ها مشابه احتمال توام ویژگی‌های معمولی حساب میشود.

تمامی کد‌ها به همراه توضیحات درباره کد‌ها و کد مورد نیاز برای تست در لینک زیر قرار دارند.

colab-file

این کد تماما به صورتی پیاده‌سازی شده که قابلیت توسعه پذیری داشته باشد. همچنین در اکثر توابع تا جایی که امکان آن بوده به جای دستور حلقه از عملیات ماتریسی برای بهبود عملکرد کد استفاده شده است.

در پیاده‌سازی برای حل مشکل کوچک شدن ضرب احتمال‌ها در تمامی ماتریس‌ها از دقت float128 پیاده‌سازی شده در پکیج نامپای استفاده شده است که برای اجرا شدن بدون مشکل کد باید روی کامپیوتر ۶۴بیتی اجرا شود و البته روی گوگل کلب نیز به درستی اجرا میشود.

نکته دیگر این که با صفر شدن برخی مقادیر ماتریس emission عملکرد الگوریتم افت میکرد برای حل این مشکل مقدار اپسیلون به تمام مقادیر این ماتریس اضافه کردم

پیشنهاد میشود در زمان تست بخش مربوط به تست الگوریتم EM را اجرا نکنید چون اجرا شدن آن زمان زیادی میبرد.

برای تست تمامی الگوریتم‌ها از روش leave one out cross validation استفاده شده است. در مورد اکثر توابع در داخل فایل گوگل کلب به صورت کامنت توضیح مختصری داده شده که ورودی و خروجی به چه صورت هست.

تمامی الگوریتم‌ها برای اجرا انتظار دارند که یک پوشه به اسم data در کنار کد باشد که دقیقا مشابه فایل

داده‌ها فرمت بندی شده یعنی در داخل این پوشه فایل مربوط به داده‌های زمینه ای هست و یک پوشه دیگر به نام samples که نمونه‌های مختلف با همان فرمت نام دهی در داخل آن قرار دارند. برای بخش تست مدل‌ها نیز دقیقاً داده‌ها را با همین فرمت باید کنار کد قرار دهیم علاوه بر این که داده‌های تست را نیز باید در پوشه samples قرار دهیم که با فرمت زیر نام گذاری شده‌اند.

test\_id.csv

## ۲ الگوریتم forward backward

توابع مختلفی در این بخش پیاده شده‌اند که در بخش‌های بعدی نیز از آن‌ها استفاده شده.

### ۱.۲ توابع محاسبه ماتریس‌های $\pi$ , $M$

برای محاسبه ماتریس  $\pi$  تابع comp- $\pi$  به ازای تمام سمپل‌ها مقادیر متغیر هدف در قدم اول را شمارش میکند و سپس این مقادیر را با تقسیم بر جمع آن‌ها نرمال میکند. برای محاسبه ماتریس transition نیز منطقی مشابه تابع قبلی در تابع comp-transitions اعمال میکنیم. به این صورت که به ازای سمپل‌های مختلف برای هر گام به جز گام اول مقدار متغیرهای هدف آن گام و گام قبلی را بررسی میکنیم و برای حالت‌های مختلف شمارش میکنیم و در آخر نرمال میکنیم.

### ۲.۲ تابع emission

برای پیاده‌سازی این تابع ابتدا یک تابع کمکی به نام ind-feature-probs پیاده‌سازی شده که برای هر ویژگی مقادیر مختلف آن را در تمام سمپل‌ها به ازای متغیرهای هدف مختلف شمارش میکند و در آخر این مقادیر را نرمال میکند.

سپس از این تابع برای پیاده‌سازی تابع emission استفاده شده به این صورت که برای هر گام در سمپل ورودی مقدار احتمال مربوطه به صورت ضرب احتمال ویژگی‌های مختلف به شرط متغیر هدف محاسبه میشود.

### ۳.۲ توابع alpha recursion, beta recursion, gamma

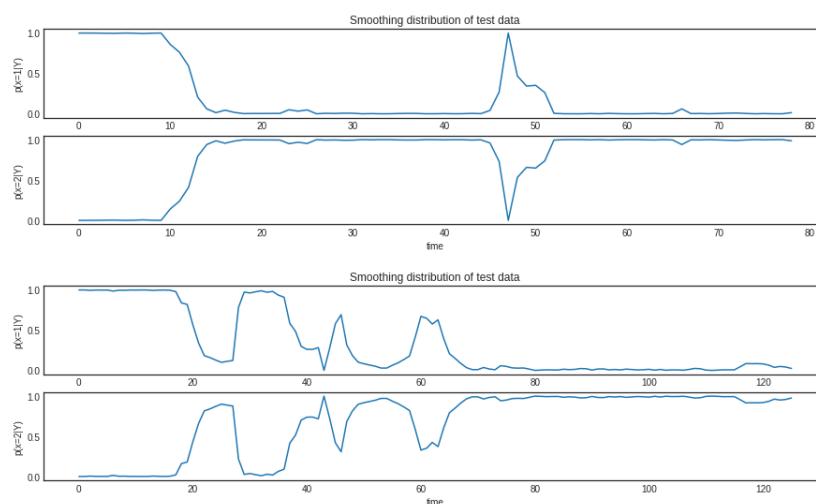
این توابع طبق فرمول‌های منبع زیر و با استفاده از برنامه نویسی پویا پیاده‌سازی شده‌اند که نسبت به حالت بازگشتی سرعت بیشتری دارد

Martin. H. James Jurafsky Daniel Processing. Language and Speech

### ۴.۲ نتایج

نتایج این بخش با اجرای تابع LOOCV بدست می‌آید این تابع برای الگوریتم مورد نظر leave one out cross validation را اجرا میکند و به ازای هر تکرار دقت روی سمپل تست و نمودار gamma مربوط به آن را چاپ میکند. در نهایت نیز دقت میانگین همه تکرارها را چاپ میکند. همه این نتایج در فایل کلب در بخش leave one out cross validation on forward-backward and viterbi algorithms قابل مشاهده هستند. تمامی الگوریتم‌های مربوط به اجرا انتظار دارند فرمت دیتا‌ها دقیقاً مشابه دیتاهای داده شده باشند. یعنی یک پوشه data کنار کد داریم که در آن فایل مربوط به ویژگی‌های زمینه‌های با همان نام و یک پوشه samples قرار دارد که نمونه‌های مختلف با فرمت نام داده شده در این پوشه قرار دارند. در زیر دو نمونه از نمودارهای گاما را می‌بینیم: دقت مربوط به سمپل یک در حالتی که این سمپل به عنوان تست استفاده

شود برابر ۹۱ درصد و دقت مربوط به سمپل دو وقتی که به عنوان تست در نظر گرفته شود ۸۹ درصد میباشد. منظور از دقت، دقت دسته بندی است. در نهایت دقت میانگین نیز برابر ۸۳ درصد میباشد.



شکل ۱: نمودار گاما به ازای سمپل یک در بالا و سمپل دو در پایین

## ۳ الگوریتم EM

برای پیاده سازی این بخش به دلیل این که بیش از یک سمپل داریم از توضیحات صفحه ویکیپدیای الگوریتم baum-welch استفاده شده است

### ۱.۳ pair marginals

برای پیاده سازی این تابع از همان منبع بخش مربوط به alpha-recursion استفاده شده است

### ۲.۳ تابع EM

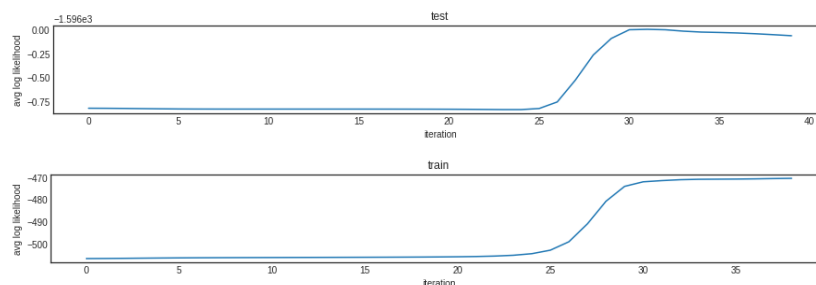
برای پیاده سازی این تابع ابتدا دو تابع کمکی MStepEmission و MStepEmissionTest پیاده سازی شده اند که برای محاسبه ماتریس های emission مربوط به داده های آموزش و آزمون استفاده میشوند. سپس در تابع EM در هر تکرار ابتدا برای هر سمپل E-step را انجام میدهم و سپس با استفاده از گاما و زتای مربوط به همه سمپل ها و با توجه به رابطه صفحه ویکیپدیا پارامترهای HMM را آپدیت میکنیم.

### ۳.۳ نتایج

نتایج مربوط به این بخش در بخش EM algorithm leave one out cross validation آمده اند

برای بررسی این الگوریتم نیز از روش LOOCV استفاده میکنیم که به ازای هر سمپل که به عنوان تست انتخاب میشود دقت مربوط به آن سمپل، توزیع گامای مربوط به آن سمپل و  $\log$ -likelihood روی داده‌های تست و آموزش را چاپ میکند.

در این تابع از پارامترهای بدست آمده توسط EM برای بدست آوردن گاما استفاده میشود این نتایج به ازای تمامی سمپل‌ها در فایل کلب قابل مشاهده اند. مقدار دقت میانگین در نهایت برابر با ۷۳ درصد میشود. برای سمپل شماره یک به عنوان تست مقادیر  $\log$ -likelihood به صورت زیر هست



شکل ۲: مقادیر لاگ لایکلیت به ازای تکرارهای مختلف الگوریتم روی مجموعه تست در بالا و آموزش در پایین

### ۴.۳ تحلیل

همانطور که مشاهده میشود مقدار  $\log$ -likelihood با گذشت زمان روی مجموعه آموزش بیشتر میشود که این همان هدف الگوریتم EM هست که احتمال دیده شده دنباله فیچر‌ها را بیشینه کند

در مورد مجموعه تست نیز ابتدا این مقدار افزایش میابد اما سپس در یک نقطه شروع به کاهش میکند که میتواند نشان دهنده بیش برازش الگوریتم باشد. در واقع این نمودار به ما میگوید باید برای تکرار کمتری الگوریتم را اجرا میکردیم.

## ۴ الگوریتم viterbi

نتایج مربوط به این بخش در بخش forward-backward leave one out cross validation and viterbi algorithm آمده اند

هدف این الگوریتم این هست که با داشتن پارامترهای HMM و با دیدن یک دنباله از مشاهدات دنباله متغیرهای پنهانی را بدست آورد که بیشترین احتمال را دارد. شبه کد این الگوریتم به صورت زیر هست

```
function VITERBI(observations of len T, state-graph of len N) returns best-path, path-prob
create a path probability matrix viterbi[N,T]
for each state s from 1 to N do ; initialization step
    viterbi[s,1] ←  $\pi_s * b_s(o_1)$ 
    backpointer[s,1] ← 0
for each time step t from 2 to T do ; recursion step
    for each state s from 1 to N do
        viterbi[s,t] ←  $\max_{s'=1}^N \text{viterbi}[s',t-1] * a_{s',s} * b_s(o_t)$ 
        backpointer[s,t] ←  $\arg\max_{s'=1}^N \text{viterbi}[s',t-1] * a_{s',s} * b_s(o_t)$ 
bestpathprob ←  $\max_{s=1}^N \text{viterbi}[s,T]$  ; termination step
bestpathpointer ←  $\arg\max_{s=1}^N \text{viterbi}[s,T]$  ; termination step
bestpath ← the path starting at state bestpathpointer, that follows backpointer[] to states back in time
return bestpath, bestpathprob
```

شکل ۳: شبه کد ویتربی

این الگوریتم از گام اول شروع میکند و در هر گام مسیر با بیشترین احتمال برای رسیدن به استیت‌های مخفی مختلف در گام حال حاضر را ذخیره میکند. سپس در گام بعد از این احتمال‌های ذخیره شده برای بدست آوردن محتمل‌ترین مسیر برای رسیدن به آن گام استفاده میکند. برای بررسی این روش نیز از LOOCV استفاده شده که نتایج آن در فایل کلب قابل مشاهده است. میانگین دقت دسته‌بندی برای این روش برابر ۸۲ درصد میباشد. مشاهده میشود که نتایج این الگوریتم نزدیک به الگوریتم forward-backward میباشد. مشاهده میشود هر دو این الگوریتم‌ها در مورد نمونه داده‌های ۱۶ و ۲۱ و ۲۳ عملکرد نامناسبی دارند. در مورد نمونه داده ۲۳ چون نمونه با متغیر پنهان ۲ شروع شده باعث میشود که احتمال اولیه یعنی  $\pi_i$  بسیار کم باشد که در این حالت برابر صفر میشود چون خود این نمونه در نمونه‌های آموزش نیست و در ادامه هر دو الگوریتم را دچار مشکل میکند. در مورد نمونه داده ۲۱ نیز چون متغیر پنهان فرکانس تغییرات بالایی دارد احتمالاً هر دو الگوریتم دچار مشکل میشوند.

## ۵ متغیرهای زمینه‌ای

نتایج مربوط به این بخش در بخش contextual data آمده اند

برای استفاده از متغیرهای زمینه احتمال هر مجموعه از متغیرهای زمینه ای به شرط متغیر هدف با توجه به فرض گفته شده در مقدمه محاسبه میشود و برای هر سمپل احتمال مربوط به متغیرهای زمینه ای آن به ازای متغیرهای هدف مختلف در ماتریس emission آن سمپل ضرب میشوند و سپس از الگوریتم forward-backward استفاده میشود.

در واقع میتوان اینگونه فرض کرد که برای هر سمپل متغیرهای زمینه ای آن را نیز به بقیه ویژگی‌ها اضافه میکنیم به این صورت که این متغیرها را برای یک سمپل به ازای تمامی گام‌ها کنار بقیه ویژگی‌ها تکرار میکنیم و سپس الگوریتم forward-backward را مانند بخش دوم اجرا میکنیم  
برای بررسی این بخش نیز از LOOCV استفاده میشود که نتایج مربوط به آن در فایل کلب قابل مشاهده است. در نهایت میانگین دقت برای این روش ۷۸ درصد میباشد. مشاهده میشود که استفاده از این ویژگی دقت را نسبت به دقت بدست آمده در بخش دوم یعنی ۸۳ درصد کاهش داده است

## ۶ ویژگی پنجم

نتایج مربوط به این بخش در بخش LOOCV on forward-backward algorithm without fea-ture 5 آمده اند

با حذف ویژگی پنجم یک مشکل پیش می‌آمد که برخی از گام‌ها همه ویژگی‌هایشان null میشد برای رفع این مشکل در این بخش و بخش بعد همه مقادیر null را با صفر جایگزین کردم که با توجه به ناپیوسته بودن ویژگی‌ها این حالت نیز میتواند به عنوان یک حالت دیگر از ویژگی‌ها یادگیری شود. یعنی انتظار داریم الگوریتم یاد بگیرد که احتمال null بودن یک ویژگی به شرط متغیر پنهان چیست.

در ادامه از ۴ ویژگی دیگر برای اجرای الگوریتم forward-backward استفاده کردم. برای بررسی این روش نیز از LOOCV استفاده کردم که نتایج آن در فایل کلب قابل مشاهده است. میانگین دقت برای این حالت برابر ۸۲ درصد هست که نشان میدهد وجود متغیر پنجم در دسته بندی به ما کمک میکرد چون دقت کمتر شده.

## ۷ حالات دیگر

نتایج مربوط به این بخش در بخش different combination of features آمده اند

برای مدل‌های بیشتر من با استفاده از ایده بخش قبل یعنی صفر کردن متغیرهای null از تمامی ترکیب‌های ممکن ویژگی‌ها برای آموزش مدل استفاده کردم.  
یعنی دقیقاً طبق بخش قبل عمل کردم با این تفاوت که در هر مدل ترکیب ویژگی‌های استفاده شده متفاوت

هست

f1 , f2 , f3 , f4 , f5

f1 , f2 , f3 , f5

f1 , f2 , f4 , f5

در این حالت نیز با پنج ویژگی به دقت ۸۳ درصد رسیدم که مشابه دقت بدست آمده در بخش دوم هست. این موضوع نشان میدهد که دو روش برخورد با متغیرهای null که در این بخش (صفر کردن این متغیرها) و بخش دوم (در نظر نگرفتن این ویژگی‌ها) استفاده شده عملکرد مشابهی در نتیجه نهایی دارند. البته نتیجه این حالت حدود ۰/۰۳ درصد از روش قبلی بهتر هست



نتیجه اجرای LOOCV برای تمامی ترکیب‌های ویژگی‌ها در فایل کلب قابل مشاهده هست. مشاهده مشود ترکیب ویژگی‌های یک و دو چهار و پنج عملکرد بهتری از حالتی دارد که همه ویژگی‌ها را استفاده کنیم.

### ۸ تست

برای بررسی بهترین ۳ مدل به دست آمده در بخش قبل می‌توانید از تابع TEST استفاده کنید که ورودی آن لیستی از اعداد int هست که نشان دهنده ایدی فایل‌های تست می‌باشد انتظار داریم این فایل‌ها نیز در مسیر همان سمپل‌های آموزشی باشند. توضیحات بیشتر در مورد نحوه قرارگیری نمونه‌های تست در مقدمه آمده است. همچنین انتظار می‌رود ستون‌ها دقیقاً به همان شکل ستون‌های داده‌های آموزش ناگذاری شده باشند و ستون target را نداشته باشیم. نتایج نیز در مسیر کد با فرمت اسم زیر ذخیره می‌شوند.

output\_id\_model model id.csv  
که در آن id همان id فایل تست هست که همان شماره ای هست که بعد از test\_ در نام فایل‌های تست قرار می‌گیرد.  
همچنین model id نیز نشان دهنده مدل استفاده شده می‌باشد. که به صورت زیر مشخص می‌شود.

1 : (f1 , f2 , f3 , f4 , f5)

2 : (f1 , f2 , f4 , f5)

3 : (f1 , f2 , f3 , f5)

فایل‌های نتیجه ذخیره شده دارای یک ستون به نام pred هستند که پیش‌بینی برای هر گام زمانی در این ستون ذخیره شده است.