

فایل های ضمیمه ی گزارش:

test.py

train.py

hmm.py

ex02.ipynb

در پوشه ی output :result1.txt, result2.txt

در پوشه ی models :model_01.txt, model_02.txt, model_03.txt, model_04.txt,

model_05.txt

از روی فایل `model_init` اطلاعات ماتریس `A` و `B` و بردار احتمالات `pi` خوانده می شود. این مقادیر به عنوان مقدار اولیه به هریک از ۵ مدل داده می شوند. کار خواندن این پارامترها از روی فایل را تابع `read_model` انجام می دهد که هم در `test.py` و هم در `train.py` پیاده سازی شده. این تابع یک رشته به عنوان ادرس در ورودی می گیرد و سپس محتویات آن را خوانده و در قالب ماتریس ها و بردار `A, B, pi` خروجی می دهد. با استفاده از این اطلاعات خوانده شده از `model_init` یک مدل اولیه HMM تعریف می شود.

در `train.py` ادرس یک فایل داده های آموزشی مثل `seq_model_0x.txt` داده می شود. خط های این فایل به صورت لیستی از دنباله هایی که هر کدام نظیر یک خط داده ی آموزشی هست را تولید می کند. این لیست به تابع `train_model` از کلاس HMM داده می شود تا آموزش مدل صورت بگیرد. مدت زمانی برای اینکه کل دیتای آموزشی خوانده و آموزش انجام شود طول می کشد. این زمان در پایان نشان داده می شود.

بعد از اینکه مدل آماده شد با استفاده از تابع `write_models_to_files` در پوشه ای که ادرس آن به صورت آرگومان داده شده با همان نام ذکر شده به فرمت پارامترهایی که در `model_init` نوشته شده بودند نوشته و ذخیره می شود.

مدل ها با استفاده از کلاس HMM ساخته و آموزش داده می شوند. این کلاس در فایل `hmm.py` به طور کامل به شرح زیر پیاده سازی شده.

در این کلاس هر مدل علاوه بر ۳ پارامتر ذکر شده برای راحتی محاسبات یک لیست `dic` از حروف استفاده شده در مدل ها (حروف بزرگ ۱ تا ۶ الفبای انگلیسی) و یک عدد `N` که برابر تعداد وضعیت هاست دارد. در این کلاس ۴ تابع پیاده سازی شده.

توابع `forward` و `backward` که هر کدام به ازای دنباله ی ورودی داده شده کار می کنند و با مدلی که از آن فراخوانی می شوند مقدار `p(observation|model)` را محاسبه می کنند. تابع `forward` با محاسبه ی

ماتریس آلفا جلو می رود و در نهایت خود ماتریس آلفا و مقدار مجموع سطر آخر آن را به عنوان مقدار $p(\text{observation}|\text{model})$ برمی گرداند. تابع backward با محاسبه ی ماتریس بتا جلو می رود. و در نهایت حاصل مجموع ضرب درایه به درایه ی سطر اول بتا در π_i را به همراه ماتریس بتا برمی گرداند. هر دو این توابع از برنامه نویسی پویا برای محاسبه ی احتمال یک دنباله مشاهده ی ورودی به ازای مدل مشخص استفاده می کنند.

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

فرمول

$$\beta[T][i] = 1$$

$$\beta[t][i] = \sum_{j=1}^N A[i][j] B[o_{t+1}][j] \beta[t+1][j]$$

پیاده سازی

$$P(o|\lambda) = \sum_{i=1}^N \pi_i \beta_o(i) = \sum_{i=1}^N \pi_i \beta[0][i]$$

تخمین احتمال

$$\alpha_o(i) = \pi_i$$

$$\alpha_{t+1}(j) = \left(\sum_{i=1}^N \alpha_t(i) a_{ij} \right) b_j(o_{t+1})$$

فرمول

$$\alpha[0][i] = \pi_i$$

$$\alpha[t+1][j] = \left(\sum_{i=1}^N \alpha[t][i] A[i][j] \right) \times B[o_{t+1}][j]$$

پیاده سازی

$$P(o|\lambda) = \sum_{i=1}^N \alpha_T(i) = \sum_{i=1}^N \alpha[T][i]$$

تخمین احتمال

تابع دیگری که در HMM پیاده سازی شده تابع `train_model` هست که محتوای فایل آموزشی یا بخشی از آن را به صورت یک لیست ورودی می گیرد و عمل آموزش مدل را به ازای آن روی مدل انجام می دهد. در هر گام از آموزش یک عنصر لیست معادل یک خط فایل داده های `seq_model_0x.txt` خوانده می شود. با استفاده از تابع `baum_welch` مقادیر به دست آمده برای پارامترهای مدل محاسبه می شود. دوتا ماتریس برگشتی مقادیر به صورت دوتایی هایی هستند که اگر بر هم تقسیم کنیم مقادیر جدید `A` و `B` را می دهند. ولی این مقادیر به صورت همین ماتریس دوتایی ها نگهداری می شود تا آخرین خط داده ی آموزشی هم خوانده شود. کل دوتایی های به دست آمده از تمام خط ها باهم جمع می شوند و در نهایت با تقسیم عنصر اول بر دوم این دوتایی ها در هر سطر و ستون که هستند درایه ی نظیر آن سطر و ستون در ماتریس های `A NxN` و `B` به دست می آید. برای بردار احتمالات `pi` هم کافی ست تمام احتمالات به دست آمده در هر خروجی `baum_welch` را باهم جمع زده و در آخر بر تعداد گام های حلقه تقسیم شود (میانگین). در آخر مقدارهای `new_A` و `new_B` و `new_pi` جایگزین مقادیر `A`, `B`, `pi` در مدل می شوند و آموزش روی داده های ورودی به پایان می رسد. به ازای تعداد ایتريشن هایی که به عنوان آرگومان اول از ورودی سیستم می خواند، از روی داده های آموزشی می گذرد و پارامترهای مدل را آپدیت می کند.

تابع پیاده شده ی `baum_welch` به ازای یک خط ورودی مشخص الگوریتم باوم-ولش را پیاده سازی می کند. از برنامه نویسی پویا استفاده شده و با کمک ماتریس های `alpha` و `beta` که از الگوریتم های `forward` و `backward` به ازای همان ورودی آموزشی گرفته و با پارامترهای فعلی مدل مقادیر ماتریس های `gamma` و `xi` را محاسبه می کند. به ازای ورودی داده شده، سطر اول گاما به عنوان `pi` جدید و به ازای دو ماتریس `A` و `B` جدید هم دو ماتریس `NxN` که هر درایه اش یک دوتایی هست که بالا توضیح دادم نتیجه می دهد. برای `A` جدید به ازای کل سطرهای `xi` و `gamma` مقادیر را باید باهم جمع زد و به ترتیب در درایه های اول و دوم دوتایی ها قرار داد. برای ماتریس `B` جدید باید برای تمام سطرها مقادیر `gamma` نظیر ستون هایشان جمع شوند تا درایه دوم به دست بیاید. در صورتی که سطر مورد نظر نظیر مشاهده ای در ورودی باشد که مقدار مشخص نظیر سطر `B` را داشته باشد، این ها مجموع `gamma` های نظیرشان به عنوان درایه اول دوتایی ماتریس `B` محاسبه می شود.

$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)}$$

فرمول

$$\gamma[t][i] = \frac{\alpha[t][i] \beta[t][i]}{\sum_{j=1}^N \alpha[t][j] \beta[t][j]}$$

پیاده‌سازی

$$\kappa_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{k=1}^N \alpha_t(k) \beta_t(k)}$$

فرمول

$$\kappa[t][i][j] = \frac{\alpha[t][i] A[i][j] \beta[o_{t+1}][j] \beta[t+1][j]}{\sum_{k=1}^N \alpha[t][k] \beta[t][k]}$$

پیاده‌سازی

$$p_i^* = \gamma_1$$

$$A[i][j]^* = \frac{\sum_{t=1}^T \kappa_t[t][i][j]}{\sum_{t=1}^T \gamma[t][i]} = \frac{\text{از } i \text{ به } j \text{ بر } \square}{\text{از } i \text{ به } \square \text{ بر } \square}$$

$$\beta[i][j]^* = \frac{\sum_{t=1}^T \gamma[t][j]}{\sum_{t=1}^T \gamma[t][i]} : \text{if } o_t = v_i$$

تخمین پارامتر

در فایل `test.py` مدلی که ادرس آن داده شده از فایل متنی با تابع `read_model` خوانده می‌شود. فایل داده های تست هم از ادرس داده شده خط به خط خوانده و به صورت لیستی از دنباله ها آماده میشود. مقدار $p(o|\text{model})$ برای هر خط به ازای هر ۵ مدل محاسبه می‌شود. این محاسبه را `measure_p_o_lambda` انجام می‌دهد. در `write_results_to_file` بهترین نتیجه یعنی بیشترین احتمالی که محاسبه شده انتخاب می‌شود و مدلی که آن نتیجه را به دست داده و همراه اندازه ی احتمال محاسبه شده روی فایلی به نام و ادرس تعیین شده نوشته می‌شود.

فایلهایی که با اجرای این فایلها به دست آمده در پوشه ی `output` هست. این فایلها با دستورات زیر تولید شده اند. نتیجه ها با ۵ ایتريشن حاصل شده :

```
python train.py 5 ./hmm_data/model_init.txt ./hmm_data/seq_model_01.txt ./models/model_01.txt
python train.py 5 ./hmm_data/model_init.txt ./hmm_data/seq_model_02.txt ./models/model_02.txt
python train.py 5 ./hmm_data/model_init.txt ./hmm_data/seq_model_03.txt ./models/model_03.txt
python train.py 5 ./hmm_data/model_init.txt ./hmm_data/seq_model_04.txt ./models/model_04.txt
python train.py 5 ./hmm_data/model_init.txt ./hmm_data/seq_model_05.txt ./models/model_05.txt
```

و برای تولید نتایج و ارزیابی مدل روی دو مجموعه ی تست :

```
python test.py ./hmm_data/modellist.txt ./hmm_data/testing_data1.txt ./output/result1.txt
python test.py ./hmm_data/modellist.txt ./hmm_data/testing_data2.txt ./output/result2.txt
```

که نتیجه ی این اجراها به ترتیب فایل های پوشه ی `models` و `output` هستند.

قبل از این که این فایل ها تهیه شوند تمام این کد ها برای هر ۵ مدل و ساخت و ارزیابی شان در `ex02.ipynb` نوشته شده بود. در این فایل ۵ مدل به صورت `dictionary` با نام `models` در دسترس بود. تمامی توابع توضیح داده شده در کلاس `HMM` روی این عینا موجود هستند. توابع فایل `train.py` و `test.py` با اندکی تغییر اینجا وجود دارند و تغیر هم مربوط به این است که ادرس فایل ها را از ورودی سیستم دریافت شود به علاوه اینکه در هر اجرا فقط یه مدل آموزش داده شود.

در این نوت بوک یک تابع به نام `measure_precision_on_testing_set1` هست که دقت نهایی `result1.txt` را بر حسب جواب های فایل `testing_answer.txt` شمارش می کند. هر ۵ مدل را با ۵ ایتريشن آموزش دادم و دقت را محاسبه کردم. این دقت برابر هست با ۰,۷۵۶