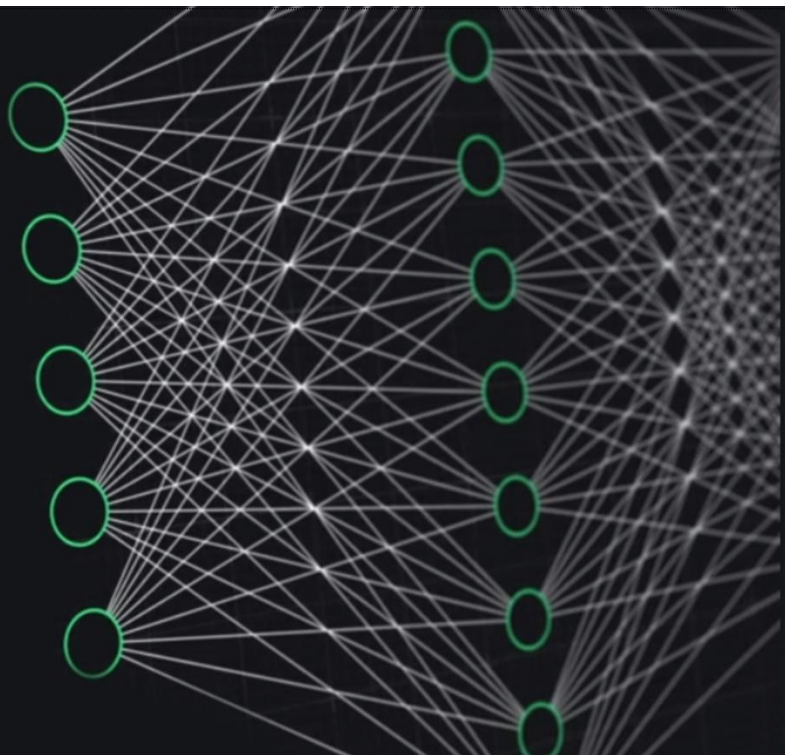


## پروژه نهایی درس یادگیری عمیق

دکتر سید ابوالقاسم میرروشندل

طراح: رضا خان محمدی



### An odd Music Generator

#### داستان پروژه:

امیرحسین که به تازگی به نواختن موسیقی علاقه مند شده، تصمیم به ثبت نام در کلاس پیانو میگیرد. او پس از تلاش فراوان و در کردن اسم و رسمی، از طرف ارکستر سمفونیک تهران دعوت به تک نوازی قطعه های یوهان سباستیان باخ در اجرای تابستان این گروه می شود. اما از آنجا که او شدیداً درگیر پایان نامه ارشد است، از شما کمک خواسته تا با استفاده از دانش یادگیری عمیق و شبکه های عصبی که در طول ترم بدست آورده اید، سیستمی هوشمند طراحی کنید تا به جای او در ارکستر قطعه های پیانو را *باخانه* (همچون یوهان سباستیان باخ) بنوازد.

#### هدف پروژه:

هدف از این پروژه، آشنایی دانشجو با خودرمزگذاران نویزدا (Denoising Autoencoders)، تکنیک های داده افزایی (Data Augmentation)، کار با فایل های صوتی، مدل های رشته-به-رشته (Seq2seq)، و سیستم های هوشمند می باشد.

#### شرح پروژه:

این پروژه شامل چرخه ای متشکل از ۴ بخش اصلی است که شما با تکمیل کردن هرکدام از آن ها می توانید سرآخر با اتصال درست آن ها به یکدیگر، سیستمی هوشمند طراحی کرده باشید که توانایی نواختن موسیقی دارد. به صورت کلی،

این ۴ بخش عبارتند از: (۱) نوپرزدایی ، (۲) تشخیص نت، (۳) پیش‌بینی نت بعدی، و (۴) نوپزگذاری. منطق کلی این چرخه به شکل روبه‌روست:



۱. **نوپرزدایی:** در طول اجرا، نوازنده باید تمام حواسش به صدای موسیقی بوده و نباید اجازه دهد تا سروصدا و آوای تشویق جمع (!) تمرکز او را از بین ببرد.

۲. **تشخیص نت:** در حین اجرا، نوازنده باید در ذهن خود سلسله‌نتهایی را که اجرا کرده را به خاطر سپرده تا از بداند در اجرای یک قطعه تا به هر لحظه چه نتهایی را زده است.

۳. **پیش‌بینی نت بعدی:** در حین اجرا، نوازنده باید در ذهن خود با توجه به سلسله‌نتهایی که تا به لحظه اجرا کرده، نت بعدی را پیش‌بینی کند تا آن را به صدا درآورد.

۴. **نوپزگذاری:** طبیعتاً از آنجایی که جمعیت حاضر در سالن از فرط لذت در پوست خود نمی‌گنجند، در حین اجرای موسیقی کلاسیک با دست و سوت سیستم تولید موسیقی شما را تشویق می‌کنند.

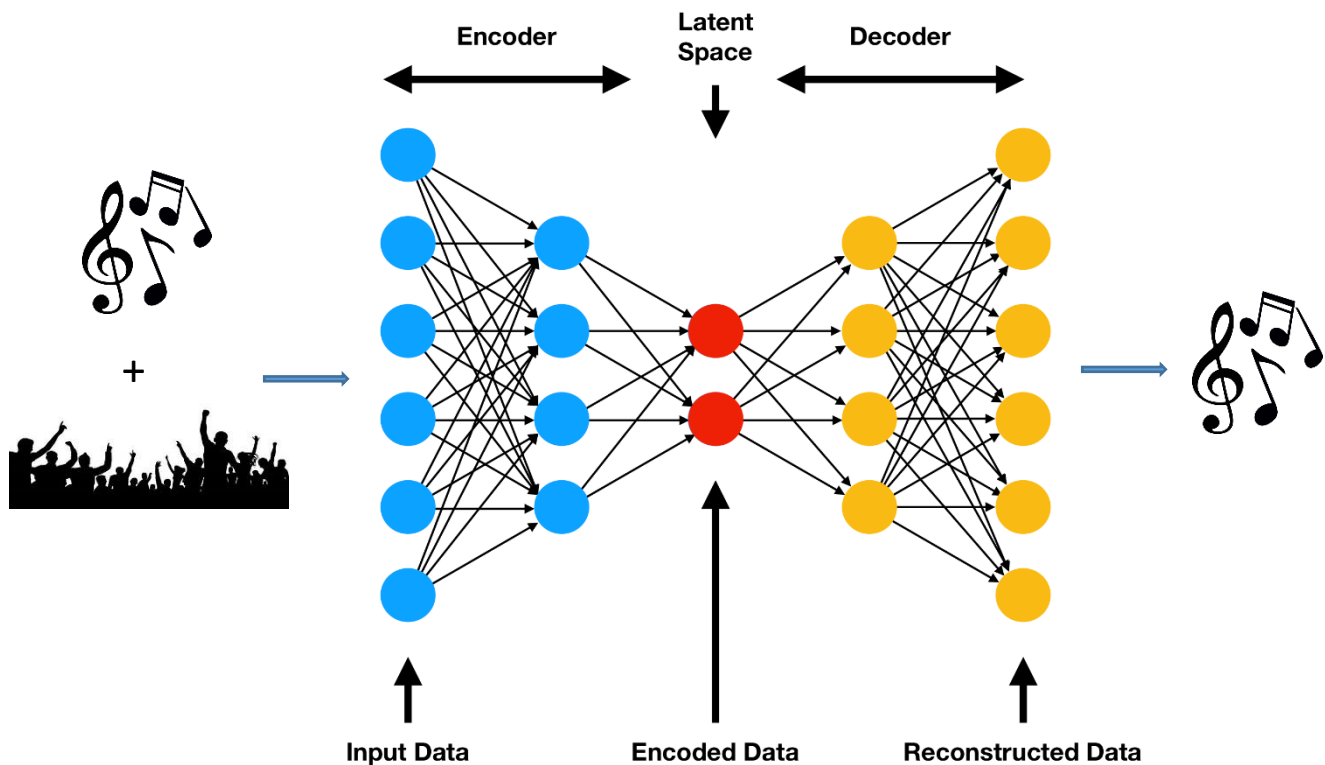
چرخه سیستم شما هم شامل این ۴ بخش اصلی بوده که به نوعی قصد دارد در هر مرحله اجرای زنده را شبیه‌سازی کند. بدیهی است که سیستم هوشمند شما به خودی خود توانایی نوپرزدایی، تشخیص نت، و پیش‌بینی نت بعدی را نداشته و این وظیفه شما است که شبکه‌های عصبی مختلفی را با این اهداف آموزش دهید. بدین منظور، شما نیاز به دیتاهای متفاوتی دارید که تمامی آن‌ها در یک فایل زیپ از آدرس ([shorturl.at/oxBQU](http://shorturl.at/oxBQU)) قابل دریافت است. این فایل زیپ شامل سه فولدر به نام های artifacts، piano\_pieces، و piano\_triads می‌باشد.

- فولدر artifacts شامل چندین فایل صوتی نوپز است که در مرحله نوپزگذاری استفاده می‌شود.
- فولدر piano\_triads شامل چندین فایل صوتی ۳ ثانیه‌ای از نت‌های (برای درک راحت‌تر، آن‌ها را نت صدا می‌کنیم) پیانو بوده که اسم هر فایل صوتی نشان‌دهنده اسم نت‌ای است که فایل صوتی آن را به صدا درمی‌آورد.
- فولدر piano\_pieces شامل ۴۰۰ فایل متنی آهنگ های پیانو کلاسیک است. در هر کدام از این فایل های متنی، نت‌های هر آهنگ با ' / ' از هم جدا شده است.

```
1 notes_list = select_random_notes(n=3)
2 initial_piece = append_notes(notes_list)
3 noisy_piece = add_random_noise(initial_piece)
4 final_piece = None
5 piece_limit, counter = 30, 0
6 while True:
7     denoised_piece = DAE(noisy_piece)
8     sub_notes = break_down_piece(denoised_piece)
9     sub_notes_classes = [NID(note) for note in sub_notes]
10    sub_notes_vec = note2vec(sub_notes_classes)
11    next_note = S2S(sub_notes_vec)
12    notes_list = sub_notes_classes + [next_note]
13    new_piece = append_notes(notes_list)
14    if next_note == '<END>' or counter == piece_limit:
15        final_piece = new_piece
16        break
17    else:
18        noisy_piece = add_random_noise(new_piece)
19        counter += 1
```

تیکه کد بالا، نمونه‌ای از این چرخه را نشان می‌دهد که سیستم شما بایستی به منظور تولید یک قطعه موسیقی آن را طی کند. در ابتدا، برنامه شما بایستی به کمک تابعی، اسامی  $n$  نت تصادفی را از بین تمامی نت‌های موسیقی که در فولدر piano\_triads موجود است، به عنوان  $n$  نت آغازین قطعه موسیقی در نظر بگیرد (خط ۱). برای مثال برنامه شما سه نت تصادفی A\_dim\_2\_0، B\_min\_3\_1، و E\_dim\_5\_1 را انتخاب می‌کند. سپس، اسامی این نت‌ها به تابع append\_notes فرستاده می‌شود و به عنوان خروجی، یک فایل صوتی به عنوان قطعه آغازین (initial\_piece) برگردانده می‌شود. بدیهی بوده که طول این فایل صوتی که شامل  $n$  نت آغازین ۳ ثانیه‌ای بوده،  $n * 3$  ثانیه است که در طی ۳ ثانیه اول نت A\_dim\_2\_0، در طی ۳ ثانیه دوم نت B\_min\_3\_1، و در طی ۳ ثانیه سوم نت E\_dim\_5\_1 به صدا درمی‌آید. حال، به کمک تابع add\_random\_noise، فایل صوتی آغازین را نوپزی می‌کنیم. شما در این تابع،

یکی از نویزهای موجود در فولدر artifacts را بر روی فایل صوتی که به عنوان ورودی به تابع داده می‌شود می‌اندازید. خروجی این تابع در اصل همان قطعه آغازین بوده که حال نویزی است. تا به اینجای کار شما مراحل اولیه‌ای که نیاز بوده تا قبل از شروع چرخه انجام شود را پیاده‌سازی کرده‌اید. حال با وارد شدن به حلقه در خط ۵، به توشیح هریک از ۴ بخش اصلی می‌پردازیم. لازم به ذکر است که برای قطعه‌های تولیدی، محدودیت ماکسیمم ۳۰ نت را در نظر می‌گیریم. همانطور که در خط ششم ذکر شده، از یک مدل Denoising Autoencoder (DAE) به منظور نویززدایی استفاده می‌کنیم:



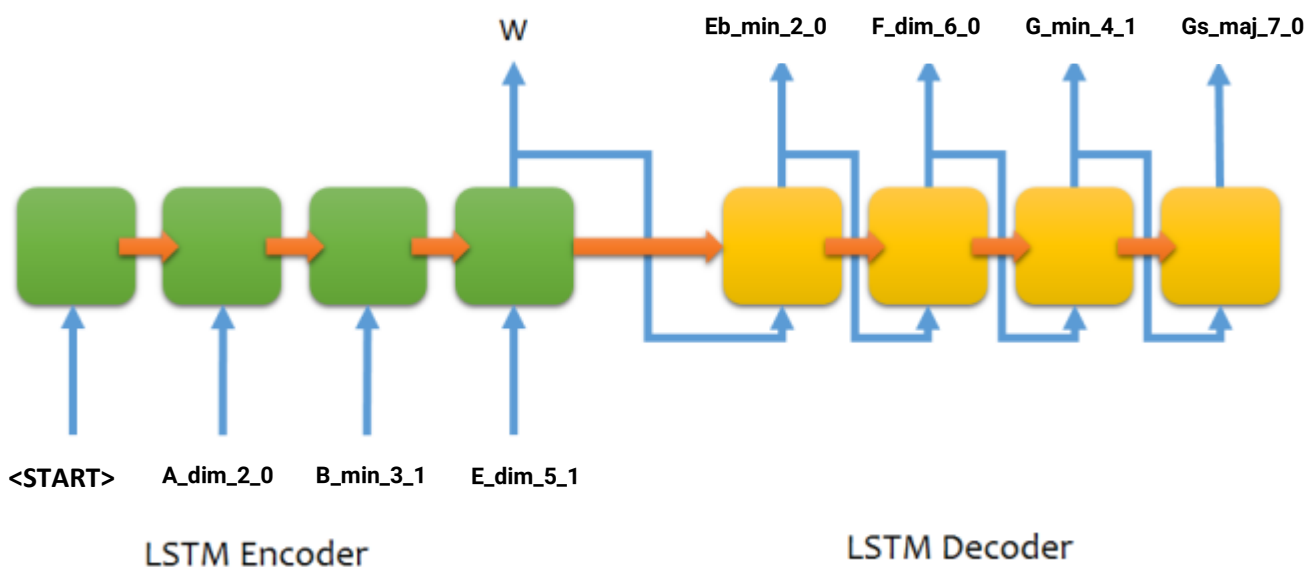
این مدل، بایستی چنان آموزش داده‌شود تا با دریافت فایل صوتی همراه نویز (noisy\_piece)، فایل صوتی عاری از نویز (denoised\_piece) را تولید کند.

پس از اتمام مرحله نویززدایی، نوبت به تشخیص نت می‌رسد. در این مرحله قصد داریم لیست نت‌هایی که تا به این مرحله در denoised\_piece نواخته شده را تشخیص دهیم. بدیهی است که می‌توان لیستی از اسامی نت‌ها در حافظه ذخیره کرد، ولی هدف ما در این بخش استفاده از مدلی جدید بدین منظور است. لذا، همانطور که در خط هفتم قابل مشاهده است، از تابعی با نام break\_down\_piece استفاده می‌کنیم تا فایل صوتی denoised\_piece را به تکه‌های ۳ ثانیه‌ای تشکیل دهنده‌اش تقسیم کنیم. حال، به کمک مدل Note Identification که بایستی قبل از ورود به چرخه آموزش دیده باشد، اسم نت‌ای که هر کدام از این فایل‌های صوتی ۳ ثانیه‌ای به صدا درمی‌آوردند را تشخیص می‌دهیم.

و اسامی نت‌ها را در sub\_notes\_classes ذخیره می‌کنیم (همانطور که احتمالاً متوجه شده‌اید، محتویات این لیست در اولین پیمایش حلقه برابر با محتویات لیست notes\_list است). مدل Note Identification صرفاً یک طبقه‌بند (Classifier) است که می‌تواند با دریافت یک فایل صوتی ۳ ثانیه‌ای، اسم نت‌ای که در این فایل به صدا درآمده را تشخیص دهد (تشخیص نت).

تا به اینجا کار، سیستم هوشمند شما می‌تواند درمیان شلوغی، صدای نت‌ها را تمیز دهد و نوع آن‌ها را شناسایی کند. حال، بایستی با توجه به توالی نت‌هایی که تا به لحظه نواخته، نت بعدی که می‌خواهد *باخانه* بنوازد را پیش‌بینی کند. بدین منظور، شما بایستی به کمک داده‌هایی که در اختیارتان قرار داده شده، یک مدل Sequence-to-sequence (S2S) را آموزش داده تا بتواند احتمال شرطی زیر را در ازای هر نت ممکن محاسبه کند:

$$\operatorname{argmax}_{w_{t+1}, \dots, w_{t+T}} \prod_{j=1}^T P(w_{t+j} | w_1, \dots, w_{t+j-1})$$



در رابطه بالا،  $T$  برابر است با تعداد تمام نت‌های موجود و توالی  $w_1, \dots, w_{t+j-1}$  به عنوان یک توالی نت شناخته می‌شود. طبق رابطه بالا، این وظیفه شبکه S2S شما بوده تا بتواند احتمال شرطی نت بعدی را با توجه به توالی نت داده شده حساب کند. نتیجتاً، خروجی خط دهم (next\_note)، برابر با آن نت‌ای بوده که احتمال شرطی محاسبه شده آن ماکسیمم باشد. نکته قابل توجه این است که (برای مثال) توالی  $[A\_dim\_2\_0, B\_min\_3\_1, E\_dim\_5\_1]$  نمی‌تواند مستقیماً به عنوان ورودی به شبکه عصبی داده‌شود و ما نیاز داریم تا این توالی را به صورت عددی که برای شبکه قابل فهم باشد درآوریم. لذا، قبل از اعمال S2S، در خط نهم تابع note2vec را فراخوانی می‌کنیم.

پس از اتمام مرحله پیش‌بینی نت بعدی، نام نت جدید را به لیست نت‌ها اضافه می‌کنیم و فایل صوتی جدید را (new\_piece)، که حال ۳ ثانیه طولانی‌تر از فایل اولیه است را به کمک تابع append\_notes تشکیل می‌دهیم. منتها، از آنجایی که شبکه‌های S2S نیازمند دو نشانگر آغاز (<START>) و پایان (<END>) هستند، شبکه می‌تواند به صلاح



خودش یادگرفته تا در زمان مناسب و زودتر از محدودیت ۳۰ نت، قطعه را به پایان برساند. بدین منظور پس از پیش‌بینی `next_note`، چک می‌کنیم که ۱) آیا محدودیت ۳۰ نت به پایان رسیده و ۲) آیا این متغیر حاوی نشانگر پایان است یا خیر. در صورت درست بودن هرکدام از این دو شرط، متغیر `new_piece` حاوی فایل صوتی نهایی سیستم است. در غیر این صورت، مرحله نویزگذاری اعمال شده که در طی آن باری دیگر یک نویز تصادفی بر قطعه جدید (`new_piece`) اضافه می‌شود.

پس از آموزش مدل‌های ذکر شده و پیاده سازی چرخه، سیستم شما بایستی در مرحله تست (`inference`)، با لود کردن مدل‌های آموزش دیده، با هربار اجرا، قطعه‌ای جدید تولید کند که طبیعتاً طول آن ۳۰ نت یا کمتر است.

## نکات تکمیلی:

- قطعه تولید شده نهایی باید قابلیت پخش شدن را داشته باشد. دقت کنید که در حین تحویل، کد نهایی بایستی اجرا شده تا در حضور هم از اجرای زنده قطعات موسیقی که سیستم هوشمندتان تولید می‌کند لذت ببریم.
- محدودیتی در انتخاب هایپرپارامترها و ساختار دقیق شبکه‌ها وجود ندارد. هرچند دانشجو موظف است تا با انتخاب مقادیر درست و آموزش بهتر مدل، کیفیت قطعه موسیقی پایانی را افزایش دهد.
- محدودیتی در انتخاب نرخ نمونه‌برداری (`sampling rate`) وجود ندارد.
- در آموزش شبکه‌های عصبی، نحوه استفاده درست از داده‌هایی که در اختیار دانشجو قرار داده شده دارای اهمیت بالایی بوده و صحت کار مدل‌ها را تعیین می‌کند.
- از آنجایی که نت‌ها متنوع بوده و به طبع تعداد کلاس‌های قابل پیش‌بینی در آموزش مدل `Note Identification` بالاست، دانشجو موظف است به مقدار نیاز از تکنیک‌های داده‌افزایی (`Data Augmentation`) استفاده کند.
- بخش قابل توجهی از نمره دهی این پروژه بر مبنای کیفیت قطعه‌های موسیقی خروجی است.
- دانشجو موظف است به منظور اثبات عملکرد مدل `Note Identification`، لیستی از نت‌های اولیه (`notes_list`) و تولیدشده (`next_note`) را در لیست `NOTES_GROUND_TRUTH` ذخیره کند و میزان شباهت اعضای آن‌را با اعضای `sub_notes_classes` مقایسه کند.
- مراحل اصلی چرخه بایستی حتماً به همین ترتیب ذکر شده انجام شود. دانشجو مختار است به غیر از ترتیب چرخه، جزییات کدها را بنا به خواسته‌اش تغییر دهد.
- گزارش پروژه بایستی کامل بوده و دقیق به سوالات و بخش‌های مربوطه پاسخ داده‌باشد.
- پروژه دانشجو بایستی نمودارهای تغییر `Loss` و `Accuracy` هر دو مجموعه آموزش و تست را داشته باشد.
- پیاده سازی به صورت گروهی (حداکثر ۲ نفره) است و هیچ محدودیتی برای زبان برنامه نویسی و فریم‌ورک یادگیری عمیق مورد استفاده وجود ندارد.
- بحث و بررسی میان دانشجویان آزاد است اما هر دانشجو موظف است پروژه را به تنهایی انجام دهد و در هنگام تحویل حضوری، به تمام جزئیات کد کاملاً مسلط باشد. با موارد **تقلب** و **کپی کردن**، طبق تشخیص دوستان حل تمرین، برخورد جدی خواهد شد.
- توجه کنید که کدهای شما باید خوانا و دارای کامنت گذاری مناسب باشد.
- زمان بندی و چگونگی تحویل حضوری پروژه، متعاقباً اعلام خواهد شد.