



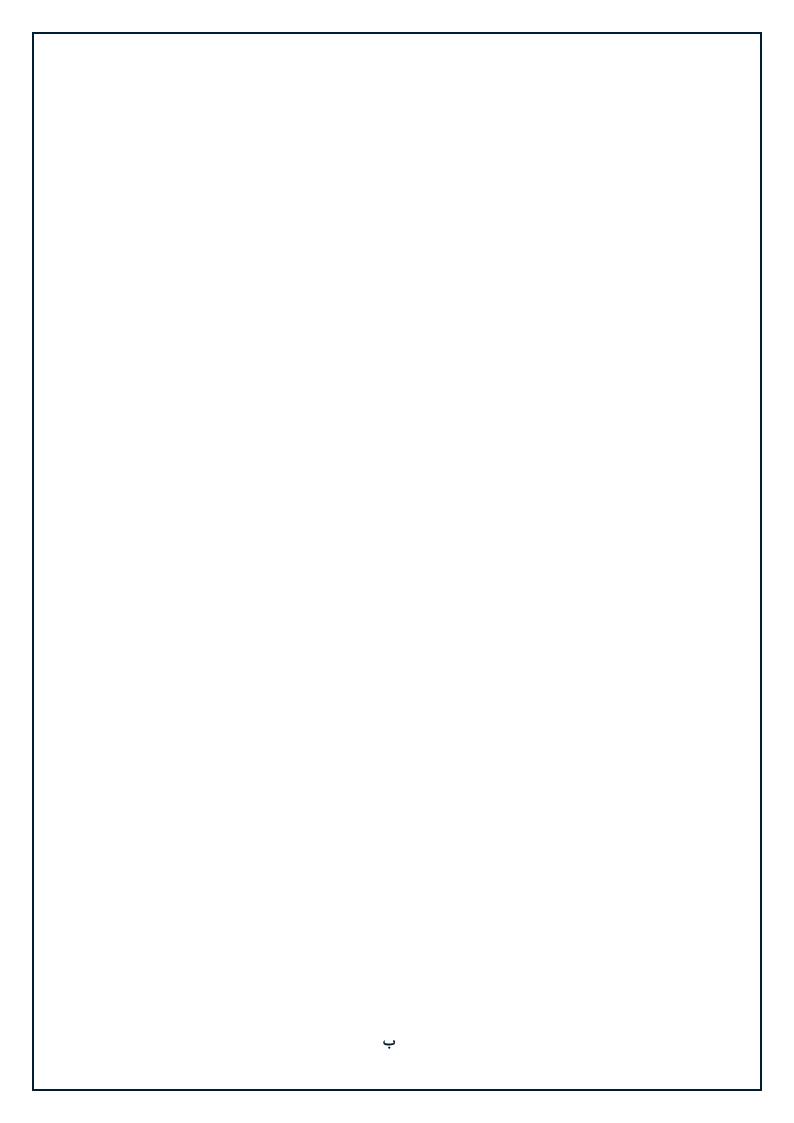
به نام خدا دانشگاه تهران دانشگده مهندسی برق و کامپیوتر

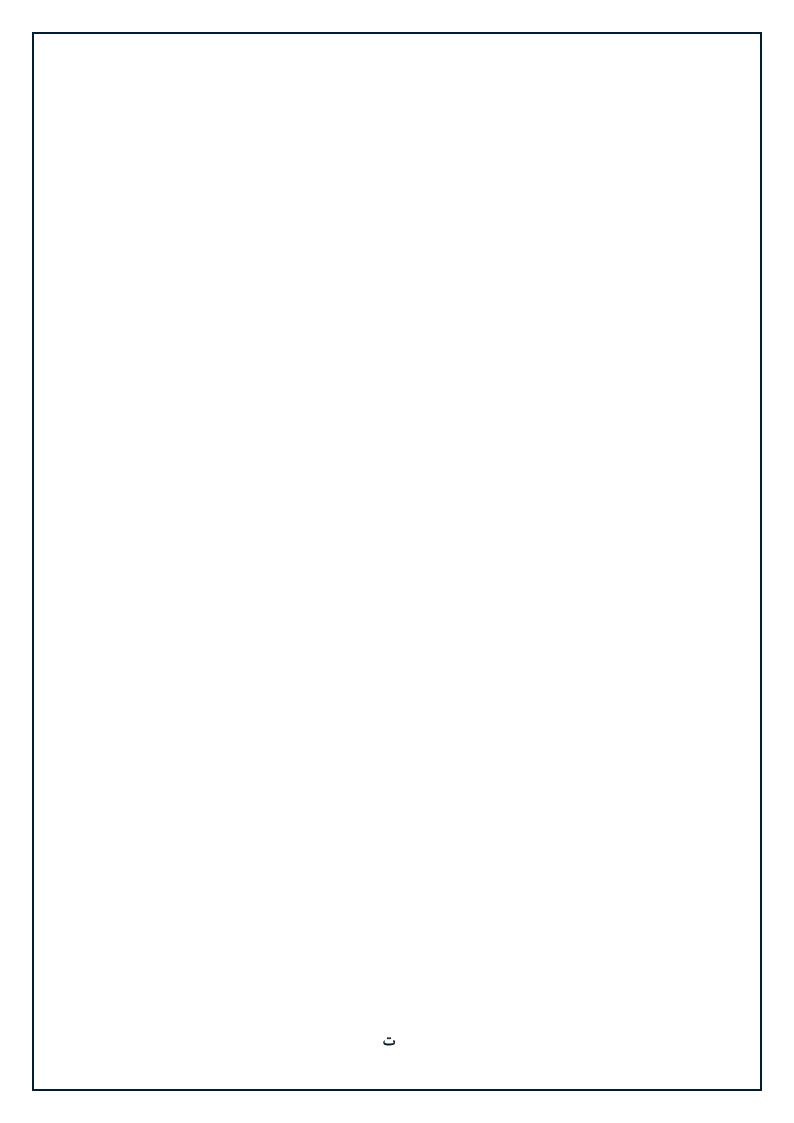
درس شبکههای عصبی و یادگیری عمیق تمرین چهارم

پرسش ۱	نام دستيار طراح	توحید عبدی
	رايانامه	Tohid.abdi@ut.ac.ir
پرسش ۲	نام دستيار طراح	محمد ولی نژاد
	رايانامه	m.valinezhad@ut.ac.ir
	مهلت ارسال پاسخ	14.441

فهرست

١	نا	قواني
١	ش ۱. تحلیل احساسات متن فارسی	پرسن
١	-١. مجموعه داده	1
١	-۲. پیشپردازش دادهها	١
١	-٣. نمايش ويژگى	1
۲	-۴. ساخت مدل	1
۲	–۵. ارزیابی	1
۲	-۶. امتيازي	1
۴	ش ۲ — سامانههای سایبرفیزیکی: نگهداری هوشمند	پرسن
۴	-۱. پیش پردازش داده ها	۲
۵	- ٢. مدل سازى و ارزيابى	۲
۶	– ۳ مقارسه را مدا های را به – ۳ مقارسه را مدا های را به	۲.





قوانين

قبل از پاسخ دادن به پرسشها، موارد زیر را با دقت مطالعه نمایید:

- از پاسخهای خود یک گزارش در قالبی که در صفحهی درس در سامانهی Elearn با نام از پاسخهای خود یک گزارش در قالبی که در صفحه درس در سامانه و REPORTS_TEMPLATE.docx
- \bullet پیشنهاد می شود تمرینها را در قالب گروههای دو نفره انجام دهید. (بیش از دو نفر مجاز نیست و تحویل تک نفره نیز نمره ی اضافی ندارد) توجه نمایید الزامی در یکسان ماندن اعضای گروه تا انتهای ترم وجود ندارد. (یعنی، می توانید تمرین اول را با شخص A و تمرین دوم را با شخص B و ... انجام دهید)
- کیفیت گزارش شما در فرآیند تصحیح از اهمیت ویژهای برخوردار است؛ بنابراین، لطفا تمامی نکات و فرضهایی را که در پیادهسازیها و محاسبات خود در نظر می گیرید در گزارش ذکر کنید.
- در گزارش خود مطابق با آنچه در قالب نمونه قرار داده شده، برای شکلها زیرنویس و برای جدولها بالانویس در نظر بگیرید.
- الزامی به ارائه توضیح جزئیات کد در گزارش نیست، اما باید نتایج بدست آمده از آن را گزارش و تحلیل کنید.
 - تحلیل نتایج الزامی میباشد، حتی اگر در صورت پرسش اشارهای به آن نشده باشد.
- دستیاران آموزشی ملزم به اجرا کردن کدهای شما نیستند؛ بنابراین، هرگونه نتیجه و یا تحلیلی که در صورت پرسش از شما خواسته شده را به طور واضح و کامل در گزارش بیاورید. در صورت عدم رعایت این مورد، بدیهی است که از نمره تمرین کسر می شود.
- کدها حتما باید در قالب نوتبوک با پسوند .ipynb تهیه شوند، در پایان کار، تمامی کد اجرا شود و خروجی خروجی هر سلول حتما در این فایل ارسالی شما ذخیره شده باشد. بنابراین برای مثال اگر خروجی سلولی یک نمودار است که در گزارش آوردهاید، این نمودار باید هم در گزارش هم در نوتبوک کدها وجود داشته باشد.
 - ullet در صورت مشاهده ی تقلب امتیاز تمامی افراد شرکت کننده در آن، 100 لحاظ می شود.
 - تنها زبان برنامه نویسی مجاز **Python** است.
- استفاده از کدهای آماده برای تمرینها به هیچ وجه مجاز نیست. در صورتی که دو گروه از یک منبع مشترک استفاده کنند و کدهای مشابه تحویل دهند، تقلب محسوب میشود.
- نحوه محاسبه تاخیر به این شکل است: پس از پایان رسیدن مهلت ارسال گزارش، حداکثر تا یک هفته امکان ارسال با تاخیر وجود دارد، پس از این یک هفته نمره آن تکلیف برای شما صفر خواهد شد.

- سه روز اول: بدون جريمه
 - o روز چهارم: ۵ درصد
 - ٥ روز پنجم: ١٠ درصد
 - روز ششم: ۱۵ درصد
 - روز هفتم: ۲۰ درصد
- حداکثر نمرهای که برای هر سوال میتوان اخد کرد ۱۰۰ بوده و اگر مجموع بارم یک سوال بیشتر از
 ۱۰۰ باشد، در صورت اخد نمره بیشتر از ۱۰۰، اعمال نخواهد شد.
- برای مثال: اگر نمره اخذ شده از سوال ۱ برابر ۱۰۵ و نمره سوال ۲ برابر ۹۵ باشد، نمره نهایی
 تمرین ۹۷.۵ خواهد بود و نه ۱۰۰.
- لطفا گزارش، کدها و سایر ضمایم را به در یک پوشه با نام زیر قرار داده و آن را فشرده سازید، سپس در سامانهی Elearn بارگذاری نمایید:

HW[Number]_[Lastname]_[StudentNumber]_[Lastname]_[StudentNumber].zip (HW1_Ahmadi_810199101_Bagheri_810199102.zip :مثال)

• برای گروههای دو نفره، بارگذاری تمرین از جانب یکی از اعضا کافی است ولی پیشنهاد میشود هر دو نفر بارگذاری نمایند.

پرسش 1. تحلیل احساسات ۱ متن فارسی

هدف از این سوال آشنایی با وظیفه ۲ تحلیل احساسات بر روی متن فارسی استخراجشده از توییتر می باشد. تحلیل احساسات یا نظرات افراد مرتبط با موضوعات مختلف است. برای آشنایی با روند سوال، این مقاله را مطالعه کنید.

۱-۱. مجموعه داده^۳

(۱۰ نمره)

برای شروع، این مجموعه داده را از Kaggle دریافت کنید. توجه داشته باشید که این مجموعه داده با داده مورد استفاده در مقاله متفاوت است. فقط ستونهای tweet و emotion برای این سوال مورد نیاز هستند. در صورت عدم دسترسی، می توانید از فایل فشرده این مجموعه داده که پیوست شده است استفاده کنید. کلاسهای موجود در ستون emotion و تعداد نمونههای هر کلاس را به کمک یک نمودار میلهای نمایش دهید.

۱-۲. پیشپردازش دادهها

(۲۰ نمره)

پیشپردازش متن در پردازش زبان طبیعی برای بهبود عملکرد مدل بسیار مهم است. مراحل پیشپردازش ذکر شده در قسمت ۳.۲ مقاله را اعمال کنید و برای هر یک از مراحل، یک مثال که پیشپردازش مورد نظر روی آن اعمال شده است را قبل و بعد از پیشپردازش چاپ کنید. توجه داشته باشید ممکن است برخی از این مراحل برای این مجموعه داده نیاز نباشد و تغییری در هیچیک از سطرها ایجاد نکند. برای انجام پیشپردازشهای این بخش میتوانید از کتابخانه و تعابخانههایی که پیشپردازش زبان فارسی را پشتیبانی میکنند استفاده کنید.

۱-۳. نمایش ویژگی^۶

Bar Plot [†]

Regular Expression ^a

Feature Representation 5

Sentiment Analysis \

Task ^۲

Dataset *

(۲۰ نمره)

در وظایف پردازش زبان طبیعی $^{\prime}$ ، داده هایی که به طور کلی پردازش می شوند، متن خام هستند. با این حال، مدلها فقط می توانند اعداد (Id) را پردازش کنند، بنابراین باید از توکنساز $^{\prime}$ ها برای تبدیل متن خام به اعداد استفاده کنید. دادههای متنی پیش پردازش شده را با توکنساز $^{\prime}$ استفاده کنید و حداکثر طول جملات برای این که تمام سطرها طول یکسانی داشته باشند، از لایه گذاری $^{\prime\prime}$ استفاده کنید و حداکثر طول جملات را برابر با ۳۲ در نظر بگیرید.

تعبیهها کلمات را بهعنوان بردارها در فضایی با ابعاد بالا نشان می دهند که روابط معنایی را به تصویر می کشند. این تعبیهها مدلهای یادگیری ماشین را قادر می سازند الگوها و احساسات را از دادههای متنی بیاموزند. در این مرحله، به کمک مدل از پیش آموزش دیده ParsBERT بردار تعبیه را برای ورودی ها به دست آورید. با تغییر پیکربندی مدل، ابعاد بردار تعبیه را برابر با ۱۲۰ در نظر بگیرید. توجه داشته باشید که برای مدیریت حافظه می توانید از تکنیکهایی مانند تکه تکه کردن و کتابخانه gc^7 استفاده کنید.

ابعاد پیشفرض بردار تعبیه در ParsBERT چقدر است؟ تعداد ابعاد این بردار بیانگر چیست؟ مفهوم بردار تعبیه را توضیح دهید و بیان کنید به نظر شما کدام یک از کلمات موجود در مجموعه داده ممکن است تعبیه نزدیک به هم داشته باشند؟ (نیازی به کد نیست)

١-٢. ساخت مدل

(۳۵ نمره)

داده ها را با نسبت $^{-}$ به دو دسته آموزش و تست تقسیم کنید و $^{-}$ از داده های آموزش را به عنوان اعتبارسنجی $^{-}$ در نظر بگیرید. الگوریتم جستجوی حریصانه $^{+}$ برای یافتن هایپرپارامترهای بهینه برای مدل CNN-LSTM را در فضای جستجو با $^{-}$ حالت زیر اعمال کنید و در نهایت هایپرپارامترهای بهینه که منجر به کمترین خطای اعتبارسنجی می شود را گزارش کنید. هایپرپارامترهای دیگر مدل را مطابق با جدول $^{-}$ مقاله اعمال کنید.

```
batch_sizes = [8, 64]
learning_rates = [0.001, 0.0001]
optimizers = [Adam, SGD]
```

Natural Language Processing '

Tokenizer [†]

Padding *

Embeddings ^f

Pretrained ^a

Configuration 5

Garbage Collector 7

Validation [^]

Greedy Search 9

در مرحله بعد، مدلهای CNN و LSTM ساده را با هایپرپارامترهای بهینه به دست آمده ایجاد کرده و آموزش دهید. نیازی به اعمال الگوریتم جستجوی حریصانه برای این مدلها نیست. به نظر شما، هر یک از این مدلها چه نقاط ضعف و چه نقاط قوتی دارند و ادغام این دو مدل با چه هدفی انجام می شود؟

۱–۵. ارزیابی

(۱۵ نمره)

دادههای تست را به کمک معیارهای ارزیابی ذکر شده در مقاله ارزیابی کنید و یک جدول مشابه جدول ۴ مقاله برای مدلهای CNN ،CNN-LSTM و LSTM چاپ کنید.

روشهای micro averaging هسمcro averaging برای محاسبه میانگین معیارهای ارزیابی را مقایسه کنید و توضیح دهید هر یک از این روشها چه تاثیری بر مقدار عددی این معیارها در این مسئله دارد.

۱-۶. امتیازی

(۵ نمره)

از روش کیسه کلمات ابرای نمایش ویژگی استفاده کنید و به کمک کتابخانه sklearn روشهای سنتی ماشین لرنینگ که در مقاله ذکر شده اند را آموزش داده و به کمک دادگان تست ارزیابی کنید. نتایج را به جدول نتایج بخش قبل اضافه کنید. برای کاهش استفاده از منابع، در این بخش میتوانید از بخشی از داده ها نمونه گرفته و از این نمونه ها برای آموزش و ارزیابی مدل ها استفاده کنید.

Sample Y Bag of Words Y

پرسش ۲ - سامانههای سایبرفیزیکی۱: نگهداری هوشمند۲

اصطلاح سامانههای سایبرفیزیکی برای اولین بار بین اعضای هیات علمی دانشگاه برکلی در سال ۲۰۰۶ معرفی شد. سامانههای سایبرفیزیکی را می توان به عنوان سامانههایی تعریف کرد که اطلاعات را از محیط فیزیکی با استفاده از حسگر ها و کانالهای ارتباطی جمع آوری می کنند، آنها را با استفاده از کنترل کننده هایی تجزیه و تحلیل می کنند، سپس محیط فیزیکی و فرآیندهای مربوطه را از طریق محرک ها برای دستیابی به یک هدف خاص، در طول عملیاتشان تحت تأثیر قرار می دهند.

گسترش سامانههای سایبرفیزیکی و پیشرفت فناوریهای اینترنت اشیا V منجر به دیجیتالیزه شدن بخش صنعتی شده است. این سامانهها امکاناتی را برای فرآیند تولید با قابلیت اطمینان A ، در دسترس بودن P ، قابلیت نگهداری C و ایمنی C بالا فراهم می کنند، اما از طرف دیگر این سامانه ها ارزیابی سلامت را پیچیده تر و چالش برانگیز می کنند. هر دستگاهی به دلیل فرسودگی و بالا رفتن سن، سلامت آن به تدریج بدتر شده و در نهایت از کار می افتد. به این نقطه از عمر دستگاه یک نقطه خرابی یا failure می گویند. در پیش از این نقطه باید تعمیر و نگهداری را اجرا کنیم تا دستگاه را به وضعیت سالم برگردانیم و سامانه دچار اختلال نشود.

در این تمرین به استفاده از مدل پیشنهادی مقاله که به منظور classification و regression مفید باقی مانده ۱۲ معرفی شده است، میپردازیم. این مقاله یک مدل یادگیری عمیق CNN LSTM را برای این منظور معرفی کرده است و آزمایش هایی را بر روی مجموعه دادههای NASA's C-MAPSS به منظور ارزیابی و مقایسه با مدلهای پایه انجام داد.

۲-۱. پیش پردازش داده ها

(۳۵ نمره)

در این بخش میخواهیم با مجموعه دادههای <u>NASA's C-MAPSS</u> آشنا شویم و پیش پردازشهای مورد نیاز را انجام دهیم. (از طریق این لینک می توانید به مجموعه داده دسترسی پیدا کنید. از آنجا که

Cyber-physical Systems (CPSs) \

Intelligent Maintenance 7

Sensor *

Communication channels *

Controller ^a

Actuator 5

IoT ^v

Reliability [^]

Availability 9

Maintainability '

Safety 11

Remain Useful Life (RUL) 17

داده ها در قالب فایل متنی می باشند، به منظور راحتی برای شما، فایل csv آن ها نیز پیوست شده است. از فایل های test_FD001 ، train_FD001 و RUL_FD001 برای این بخش استفاده کنید.)

- در ابتدا توضیح مختصری در مورد مجموعه داده ذکر شده دهید.
- سپس با توجه به عملیاتهای صورت گرفته در مقاله برای پیش پردازش داده ها اقدام کنید. پیش پردازشهای صورت گرفته شامل Data labeling ، Data Normalization ، Data Selection و Timing window
- این عملیات را برای مدل regression و classification انجام دهید، به صورت کامل عملیاتی که انجام داده اید را شرح داده و تحلیلهای خود را با دلیل بیان کنید.

راهنمایی: برای محاسبه ستون مربوط به Remain Useful Life برای هر سطر دادههای test، ابتدا بیشترین cycle مربوط به هر موتور در دادههای test را پیدا کرده و با مقادیر متناظر در فایل RUL_FD001 بیشترین خده و با مقادیر متناظر در فایل cycle هر سطر تست تفریق کنید.

۲-۲. مدل سازی و ارزیابی

(۳۵ نمره)

حال در این بخش قصد داریم مدل ارائه شده توسط مقاله برای regression و بیاده سازی و ارزیابی کنیم. (نکته: دستیابی به عددهای ذکر شده در مقاله ضرورتی ندارد. همچنین معیار های ارزیابی مدل ها در مقاله توضیح داده شد، باتوجه به آن توضیحات اقدام به ارزیابی مدل ها کنید.)

الف) Classification: در این قسمت مدل پیشنهادی مقاله را پیاده سازی کنید. سپس باتوجه به پارامترهای بیان شده در مقاله و دادههای تهیه شده در بخش قبل اقدام به آموزش مدل کنید.

- ابتدا مدل را بدون استفاده از رویکرد Early stopping آموزش دهید.(,coptimizer=adam
- در گام بعد اقدام به ارزیابی مدل با استفاده از دادههای test تهیه شده در بخش قبل کنید و معیارهای در گام بعد اقدام به ارزیابی مدل با استفاده از دادههای accuracy را برای هر کلاس گزارش کنید و نتایج را تحلیل کنید.
- همچنین ماتریس درهم ریختگی ٔ و نمودار منحنی ROC را ترسیم کنید و نمودارها را تحلیل کنید.
- حال با استفاده از رویکرد Early stopping این مراحل را تکرار کنید. نتایج حاصل از دو بخش را مقایسه و تحلیل کنید.

Receiver Operating Characteristic curve ^r Confusion matrix ^l

ب) Regression: در این قسمت مدل پیشنهادی مقاله را پیاده سازی کنید. سپس باتوجه به پارامترهای بیان شده در مقاله و دادههای تهیه شده در بخش قبل اقدام به آموزش مدل کنید.

- ابتدا مدل را بدون استفاده از رویکرد Early stopping آموزش دهید. (optimizer=rmsprop
- در گام بعد مدل را در دوحالت استفاده از تمام پنجرههای زمانی و حالت استفاده از آخرین پنجره زمانی دادههای test تولید شده در بخش قبل ارزیابی کنید و معیارهای MAE ،MSE ،RMSE و MAP و MAP را گزارش کنید و نتایج آن ها را تحلیل کنید.
- همچنین نمودار مقادیر واقعی و تخمین زده شده برای RUL^1 را ترسیم کنید. نمودار ها را تحلیل کنید.
- حال با استفاده از رویکرد Early stopping این مراحل را تکرار کنید. نتایج حاصل از دو بخش را مقایسه و تحلیل کنید.

۲-۳. مقایسه با مدلهای یایه

(۳۰ نمره)

در این بخش با استفاده از معماری CNN LSTM پیشنهادی توسط مقاله، معماری یک مدل CNN را برای دو عملیات classification و classification طراحی کنید. حال با استفاده از دادههای تهیه شده در بخش ۱، مدلهای طراحی شده را با استفاده از رویکرد Early stopping آموزش داده و مانند بخش ۲ آن ها را ارزیابی کرده و نتایج و تحلیلهای خود را بیان کنید. (نکته: دستیابی به عددهای ذکر شده در مقاله ضرورتی ندارد.)

Remain Useful Life 1