

به نام خدا



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

درس پردازش زبان طبیعی  
استاد ممتازی

تمرین دوم

علیرضا مازوچی

۴۰۰۱۳۱۰۷۵

## بخش اول: آشنایی با روش‌های بازنمایی کلمات

گام سوم: یافتن اسناد مشابه

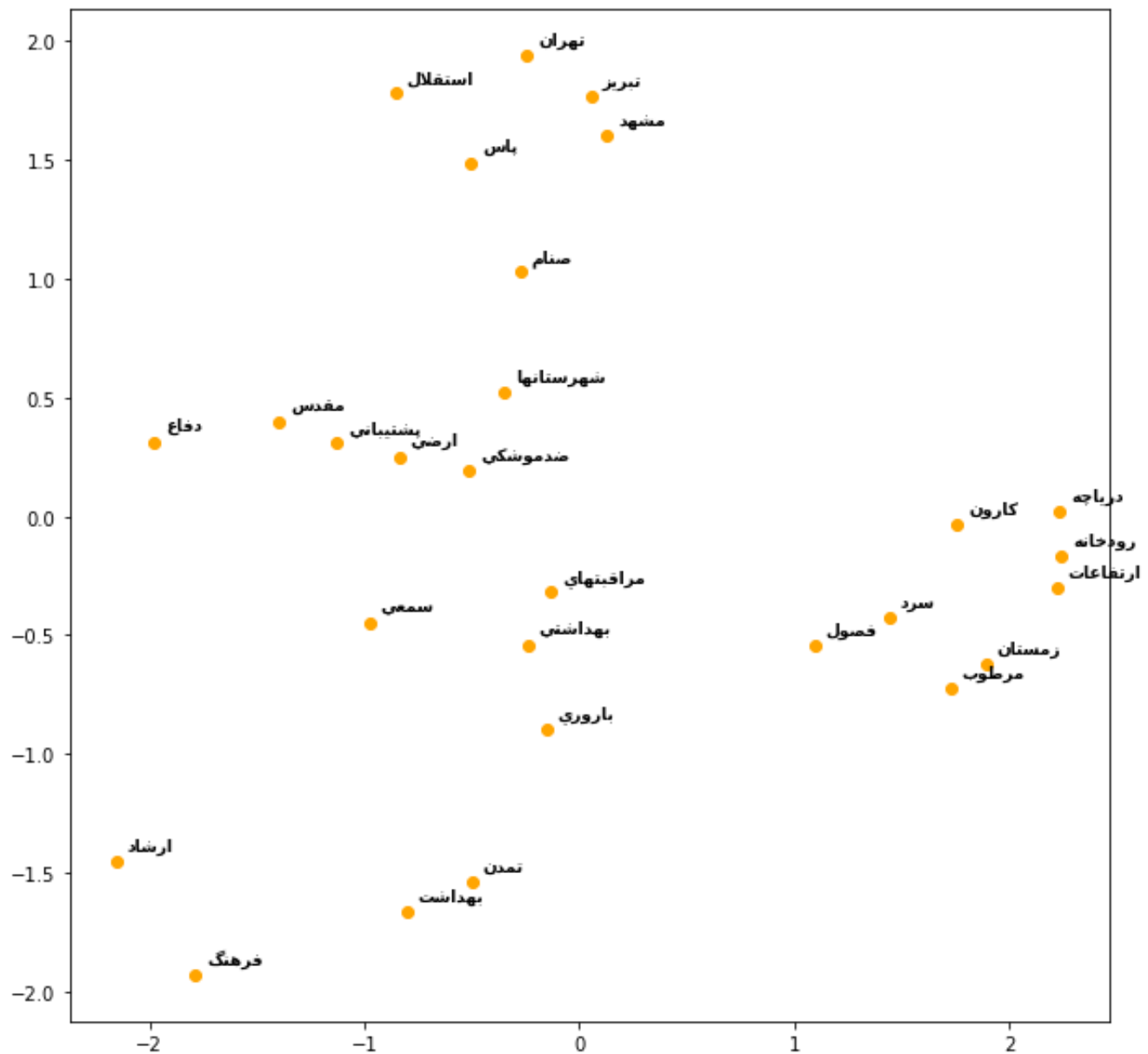
مدل Doc2Vec		مدل Word2Vec و TF-IDF		
امتیاز کسینوسی	شبیه‌ترین سند	امتیاز کسینوسی	شبیه‌ترین سند	سند
79.60%	Doc13	98.76%	Doc13	Doc1
83.68%	Doc20	99.61%	Doc19	Doc3
81.29%	Doc26	98.96%	Doc26	Doc5
95.02%	Doc679	100%	Doc679	Doc25
84.04%	Doc550	98.73%	Doc7	Doc36

گام چهارم: بررسی کلمات مشابه

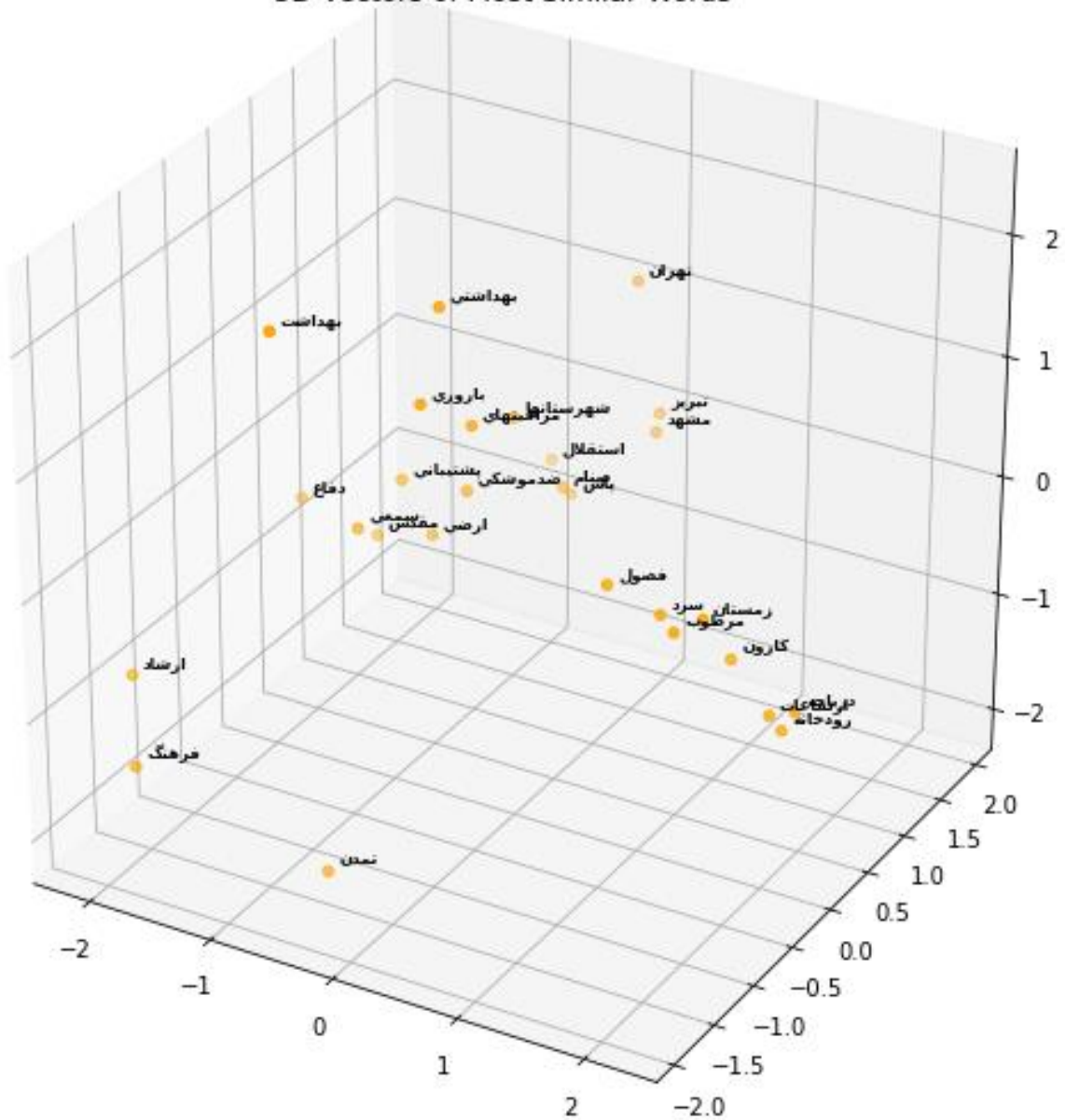
کلمه اول	امتیاز کلمه اول	کلمه دوم	امتیاز کلمه دوم	کلمه سوم	امتیاز کلمه سوم
تهران	۶۱/۲۵%	شهرستانها	۶۰/۱۲%	مشهد	۵۹/۹۵%
بهداشت	۷۸/۱۶%	بهداشتی	۷۴/۱۳%	مراقبت‌های	۷۳/۶۲%
دفاع	۶۹/۹۳%	ضدموشکی	۶۶/۲۴%	پشتیبانی	۶۴/۴۶%
رودخانه	۸۳/۸۸%	کارون	۸۲/۰۹%	ارتفاعات	۸۰/۸۱%
سرد	۷۸/۶۰%	مرطوب	۷۷/۷۳%	فصول	۷۷/۳۳%
فرهنگ	۷۷/۱۹%	تمدن	۶۹/۸۶%	سمعی	۶۷/۳۷%
استقلال	۶۷/۹۸%	ارضی	۶۷/۸۳%	صنام	۶۶/۶۰%

در دو صفحه بعد نمودار دو بعدی و سه بعدی درخواست شده ارائه شده است:

2D Vectors of Most Similar Words



3D Vectors of Most Similar Words

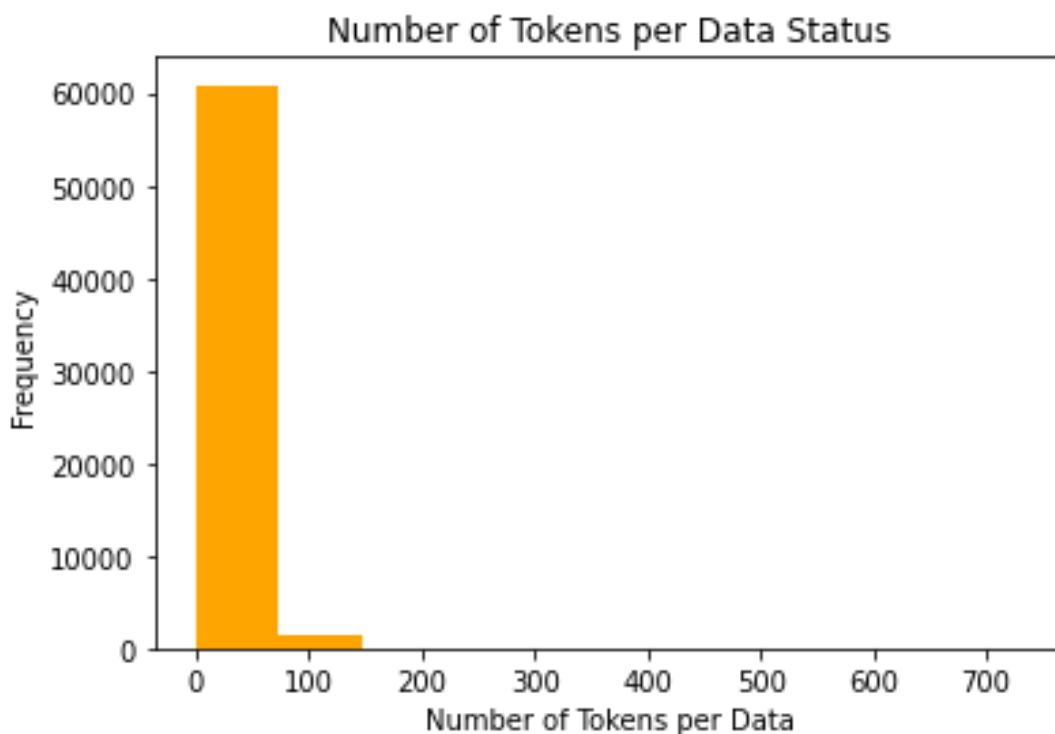


- قبل از بیان هر نکته باید گفت که موارد مطرح شده به طور کلی برقرار است و ممکن است برای برخی از داده‌ها فاصله در نمودار به طرز غیرمنتظره‌ای زیاد شود؛ این امر به دلیل آن است که فضا کاهش بعد شدیدی را تجربه کرده است و طبیعتاً اگر تمام ابعاد قابل حفظ و نمایش بود، موارد استثنا دیده نمی‌شد. باتوجه به کاهش ابعاد حتی برخی از نقاط در نمودار دوبعدی با نمودار سه بعدی فاصله کاملاً متفاوتی دارند؛ مانند «بهداشت» و «تمدن» که در نمودار دوبعدی در کنار هم و در نمودار سه بعدی با فاصله زیاد قرار گرفته‌اند.
- به طور کلی می‌توان دید که هر کلمه با سه کلمه نزدیک به خود در فضای هندسی نزدیک افتاده است. مثلاً «رودخانه»، «دریاچه»، «کارون» و «ارتفاعات» به خوبی در کنار هم قرار گرفته‌اند. این کنار هم بودن تا حد زیادی تابع امتیازهای موجود هم هست. مثلاً برای «رودخانه» و کلمات مشابه آن امتیازها همگی بالای ۸۰٪ است و نزدیکی بیشتر در مورد آن‌ها دیده می‌شود.
- به طور کلی می‌توان دید که از سه کلمه نزدیک به هر کلمه آن کلمه که در رتبه بهتری قرار داشته است، در اینجا هم به آن نزدیک‌تر است. مثلاً برای «بهداشت» به ترتیب کلمات «باروری»، «بهداشتی» و «مراقبتها» به آن نزدیک است. به عنوان مثالی دیگر برای کلمه «فرهنگ»، کلمه «ارشاد» با اختلاف مشابه‌ترین کلمه است و در نمودار هم از دو کلمه مشابه دیگر نزدیک‌تر به «فرهنگ» قرار دارد.
- در نمایش‌های مختلف برخی از کلمات به دسته‌ای دیگر از کلمات نزدیک شده است. مثلاً در نمودار دو بعدی «استقلال» به «تهران» نزدیک است. این مورد می‌تواند به دلیل شهر تیم فوتبال استقلال باشد.
- در دو نمودار یک شکاف و فاصله میان دسته کلمات «رودخانه» و «سرد» و کلمات مشابهشان با بقیه کلمات دیده می‌شود. این نشان می‌دهد که احتمالاً در فضای اصلی هم این دسته از کلمات فاصله زیادی با بقیه دارند که به نظر می‌رسد از نظر مفهومی کاملاً طبیعی باشد. به طور مشابه مفاهیم انتزاعی «فرهنگ» و «ارشاد» و «تمدن» به دور از بقیه کلمات افتاده‌اند.

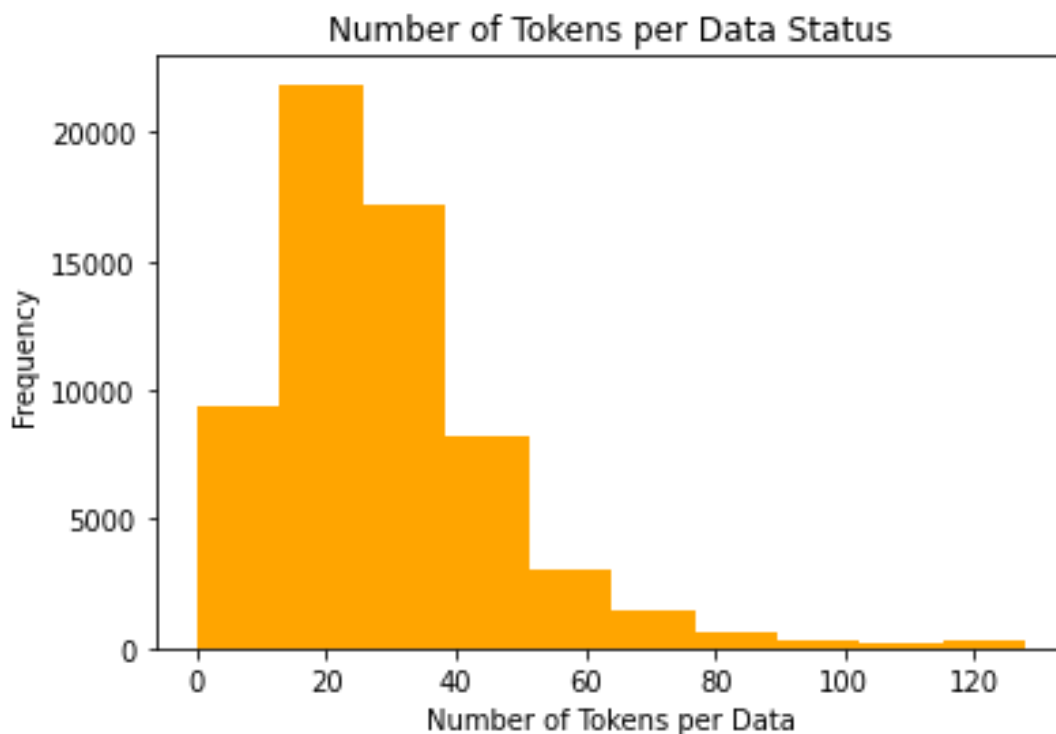
## بخش دوم: تشخیص اجزای سخن (POS)

### گام اول: ایجاد شبکه‌ی عصبی حافظه کوتاه-مدت بلند دوطرفه

ورودی شبکه BiLSTM باید دنباله‌هایی با اندازه‌های یکسان باشد؛ لذا باید طول کوتاه‌ترین دنباله‌ها را با حاشیه‌گذاری (Padding) به طول بلندترین آن برسانیم. بررسی‌های من نشان می‌دهد در سه مجموعه داده، داده‌هایی با ۷۳۴ توکن وجود دارند. اما به طور کلی جملات با تعداد توکن بالا بسیار کم هستند. به بیان دقیق‌تر در هر سه مجموعه میانگین تعداد توکن‌ها به ازای هر جمله کمتر از ۳۰ توکن است! در نمودار زیر که فراوانی داده‌ها با اندازه‌های مختلف توکن را نشان می‌دهد به خوبی مشخص است که بیشتر داده‌ها تعداد نسبتاً کمی توکن دارند.



طبیعی است که در این شرایط حاشیه‌گذاری باعث پیچیده‌شدن شدید مدل و مشکلات حافظه و زمان اجرا می‌شود؛ بنابراین من تصمیم گرفتم که داده‌های بلند را به تعداد کوچک‌تری داده بشکنم. در نمودار زیر مجدداً نمودار فراوانی داده‌ها با اندازه‌های مختلف ترسیم شده است با این تفاوت که این بار تنها در محدوده ۱۲۸ توکن است:



با نگاه به این نمودار من ۶۴ توکن را به عنوان حد بالای تعداد توکن یک داده در نظر گرفتم. بررسی‌های من نشان می‌دهد تنها ۴/۳٪ از داده‌ها بیشتر از ۶۴ توکن دارند بنابراین دقت کلی کاهش چندانی نخواهد داشت ولی آموزش بسیار آسان‌تر خواهد شد.

نهایتاً توجه کنید که برای محاسبه Accuracy خروجی مدل به ازای توکن‌ها حاشیه‌گذاری حذف شده است و تنها دقت به ازای توکن‌ها اصلی محاسبه شده است.

شبکه‌ای که من برای آموزش استفاده کردم متشکل از یک BiLSTM با ۶۴ واحد و یک لایه Dense برای تولید تگ خروجی است. از تابع فعال‌سازی Softmax برای دریافت نتایج بهتر در لایه Dense خروجی کمک گرفته شده است.

برای آموزش هم از پانزده گام استفاده شده است و برای جلوگیری از بیش‌برازش احتمالی از یک EarlyStopping Callback بهره گرفتم.

## گام دوم: ارزیابی شبکه‌ی عصبی حافظه کوتاه-مدت بلند دوطرفه

طبیعتاً نمی‌توان معیار Accuracy را به ازای هر برجسب جداگانه حساب کرد. اگر بخواهیم به ازای کل تگ‌ها حساب کنیم، نتایج جدول زیر بدست می‌آید:

صحت (Accuracy)	مجموعه داده
۹۴/۹۷٪	آموزش (Train)
۹۴/۸۱٪	اعتبارسنجی (Valid)
۹۴/۸۰٪	آزمون (Test)

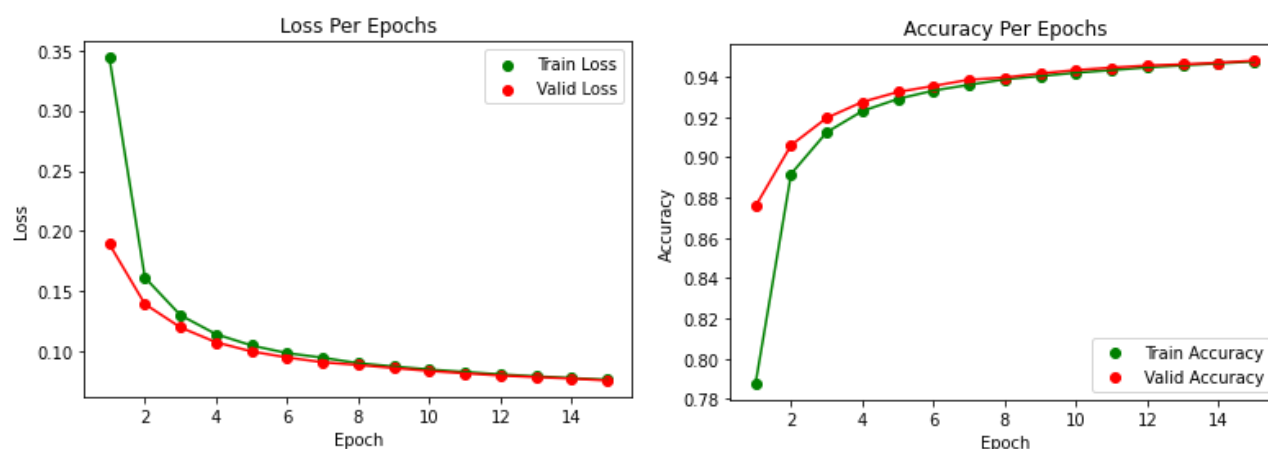
اگر هم بخواهیم برای هر تگ POS دقتی ارائه دهیم، می‌توانیم از معیارهایی نظیر Precision و Recall استفاده بکنیم که نتایج آن در جدول زیر به طور کامل ارائه شده است:

تگ	Precision	Recall
N_PL	93	91
N_VOC	77	53
ADJ_INO	76	42
N_SING	95	95
ADV_COMP	89	75
V_SUB	92	92
DET	96	95
CLITIC	99	99
V_AUX	98	91
P	98	99
ADV_I	90	86
PRO	96	96
ADV_TIME	92	91
V_PRS	96	95
DELM	88	95
FW	80	70



93	96	V_PP
35	68	V_IMP
89	90	ADJ
89	91	NUM
92	92	ADV_LOC
67	81	ADV
95	96	V_PA
35	89	SYM
98	98	CON
92	96	ADJ_SUP
86	92	ADV_NEG
62	78	INT
70	86	PREV
82	90	ADJ_CMPR

در دو نمودار زیر تغییرات خطا و صحت برای مجموعه آموزش و اعتبارسنجی آورده شده است. نکته‌ی عجیب آن است که عملکرد مدل در ابتدا بر روی مجموعه اعتبارسنجی از مجموعه آموزشی بهتر بوده است! نکته دیگری هم که می‌توان در آن دید این است که روند آموزش مدل تقریباً متوقف شده است و ادامه دادن آموزش دیگر موثر نبوده است.



## نمودار درهم‌ریختگی برای داده‌های آزمون در ادامه آورده شده است:

N_PL	2.2e+04	0	2	1.2e+03	0	5	0	0	0	7	0	8	1	19	3.2e+02	2	2	2	3.2e+02	60	0	37	9	0	10	8	0	0	0	4
N_VOC	0	7	0	2	0	1	0	0	0	0	0	0	0	2	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
ADJ_INO	20	0	2.4e+02	51	0	0	0	0	0	2	0	0	0	1	36	0	63	0	1.4e+02	4	0	2	5	0	0	0	0	0	0	1
N_SING	8.6e+02	0	11	1.3e+05	31	34	15	1	0	4.2e+02	10	38	24	54	1.7e+03	55	23	10	1.9e+03	2.6e+02	13	1.2e+02	87	2	55	12	12	10	4	10
ADV_COMP	6	0	0	24	3.7e+02	0	0	0	0	0	0	0	0	1	27	0	0	0	3	4	0	9	0	0	44	0	0	0	0	0
V_SUB	10	0	0	91	0	4.4e+03	0	0	10	1	2	0	0	61	48	1	12	7	26	1	0	4	66	0	1	0	0	0	0	1
DET	4	0	0	46	0	0	7.8e+03	0	0	0	1	2.4e+02	0	0	34	0	0	1	36	0	0	14	0	0	0	0	0	0	0	0
CLITIC	0	0	0	0	0	0	0	5.4e+03	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
V_AUX	0	0	0	8	0	8	0	0	1.8e+03	0	0	0	0	1.4e+02	2	0	1	0	1	0	0	0	12	0	0	0	0	0	0	0
P	8	0	1	1.9e+02	0	1	0	0	0	4.6e+04	0	10	8	10	56	0	1	2	54	3	8	19	2	0	7	0	0	0	6	0
ADV_I	1	0	0	9	0	0	0	0	0	0	4.8e+02	1	0	2	3	0	0	0	2	0	0	0	0	0	57	0	0	0	0	0
PRO	11	0	0	44	0	0	2e+02	0	0	2	0	9.8e+03	0	14	4	3	0	0	37	7	0	9	1	0	62	0	0	0	0	1
ADV_TIME	5	0	1	55	1	2	0	0	0	6	0	0	2e+03	0	32	0	2	1	34	5	4	37	1	0	0	1	0	0	0	0
V_PRS	36	0	0	1.3e+02	0	98	2	0	15	2	0	2	0	1.4e+04	1.6e+02	2	14	1	41	8	1	13	67	0	3	0	0	0	0	1
DELM	38	2	0	4.7e+02	0	4	0	0	1	10	0	0	1	18	3.7e+04	35	6	1	93	5.8e+02	0	27	0	0	4	0	0	2	1	0
FW	2	0	0	76	0	0	0	0	0	4	0	5	2	4	63	4.7e+02	0	0	14	7	0	2	0	0	12	0	0	0	0	0
V_PP	12	0	13	62	0	38	0	0	0	3	0	1	0	62	22	0	5.3e+03	0	62	1	1	8	54	0	2	0	0	0	0	0
V_IMP	1	0	0	44	0	36	0	0	0	1	0	0	0	18	6	1	1	66	4	0	0	1	9	0	0	0	0	0	0	0
ADJ	3.4e+02	0	40	2.5e+03	0	15	26	0	0	48	0	13	75	30	5.6e+02	2	38	4	3.2e+04	61	5	1.6e+02	21	0	3	6	2	2	1	25
NUM	25	0	0	2.1e+02	0	0	0	0	0	2	0	3	0	0	9.2e+02	4	0	0	25	1.1e+04	0	7	0	0	1	5	0	0	0	0
ADV_LOC	1	0	0	14	0	0	0	0	0	8	0	0	2	0	3	0	1	0	1	0	3.9e+02	1	1	0	0	0	0	0	0	0
ADV	46	0	1	2.6e+02	4	7	3	0	1	13	1	2	35	26	4.5e+02	0	3	1	2.7e+02	24	0	2.7e+03	14	0	1e+02	0	0	2	0	36
V_PA	22	0	3	1.5e+02	0	99	0	0	1	2	0	1	3	1e+02	53	0	36	0	29	3	0	9	1e+04	0	3	0	0	0	0	0
SYM	2	0	0	27	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	17	1	0	0	0	0	
CON	15	0	0	42	7	1	0	0	2	70	38	16	1	5	1e+02	6	0	0	18	11	0	98	2	0	3.1e+04	0	19	0	0	1
ADJ_SUP	35	0	0	28	0	0	0	0	0	0	0	0	1	1	9	0	0	0	12	1	0	1	0	0	0	1.1e+03	0	0	0	0
ADV_NEG	0	0	0	11	0	0	6	0	0	0	0	0	6	0	10	0	0	0	29	0	0	6	0	0	8	0	5.1e+02	0	0	0
INT	1	0	0	11	0	2	1	0	0	0	1	0	1	0	4	0	0	0	2	0	0	1	0	0	4	0	6	57	0	0
PREV	0	0	0	5	0	0	0	0	0	19	0	0	0	0	0	0	0	0	4	0	0	4	0	0	0	0	0	0	76	0
ADJ_CMPR	13	0	1	44	0	1	0	0	0	1	0	0	0	1	17	0	1	0	70	1	0	21	0	0	0	1	0	0	0	8e+02
N_PL	N_VOC	ADJ_INO	N_SING	ADV_COMP	V_SUB	DET	CLITIC	V_AUX	P	ADV_I	PRO	ADV_TIME	V_PRS	DELM	FW	V_PP	V_IMP	ADJ	NUM	ADV_LOC	ADV	V_PA	SYM	CON	ADJ_SUP	ADV_NEG	INT	PREV	ADJ_CMPR	



چهار جفت تگ هستند که بالای هزار مرتبه باهم دیگر اشتباه گرفته شده‌اند که اسامی آن‌ها در جدول زیر آورده شده است:

تعداد اشتباه	تگ درست	تگ پیش‌بینی‌شده
۱۱۶۰	N_PL	N_SING
۱۶۷۳	N_SING	DELM
۱۹۱۱	N_SING	ADJ
۲۴۸۰	ADJ	N_SING

طبیعتاً مدل برای داشتن Accuracy بهتر باید بتواند تفاوت میان این جفت‌ها را بهتر تشخیص دهد. البته عاملی که شاید از آن غفلت شده باشد آن است که برخی از این تگ‌ها بسیار پرتکرار هستند و یک درصد خطای کوچک روی آن هم تواند تعداد خطای بالایی را ایجاد کند. به عنوان واضح‌ترین مثال، N\_SING در هر چهار سطر جدول فوق آمده است ولی دارای Precision و Recall ای معادل ۹۵٪ است که کاملاً منطبق بر عملکرد مدل است اما تعداد تگ‌های آن به طور عمومی زیاد است به گونه‌ای که حدود ۱۳۰ هزار پیش‌بینی درست برای آن ثبت شده است.

باتوجه به آنچه گفته شد می‌توانیم ماتریس درهم‌ریختگی را سطری یا ستونی نرمال کنیم و بر مبنای درصد خطا نظر بدهیم. در جدول زیر چهار جفت از بزرگترین درصدهای خطا آورده شده است:

نوع نرمال‌سازی	تگ پیش‌بینی‌شده	تگ درست	تعداد اشتباه	درصد اشتباه
درست	N_SING	SYM	۲۷	۵۶/۲۵٪
درست	ADJ	ADJ_INO	۱۳۹	۲۴/۸۶٪
درست	N_SING	V_IMP	۴۴	۲۳/۴۰٪
پیش‌بینی‌شده	N_VOC	DELM	۲	۲۲/۲۲٪

در این جدول به عنوان مثال ۲۷ تا از تگ‌های SYM به اشتباه N\_SING تشخیص داده شده‌اند. این مقدار اگرچه زیاد نیست ولی برای تگ بسیار کم تکرار SYM زیاد است

و این مسئله باعث شده است تا ۵۶/۲۵٪ از تگ‌های درست SYM اشتباهی N\_SING تشخیص داده شوند! مقدار Recall پایین ۳۵٪ ای SYM هم این موضوع را تایید می‌کند.