



گزارش کار تمرین دوم  
درس یادگیری ماشین

استاد درس: دکتر کمندی

دانشجویان گروه:

محیا معتمدی ۸۱۰۸۹۷۰۵۳

مینو احمدی ۸۱۰۸۹۷۰۳۲

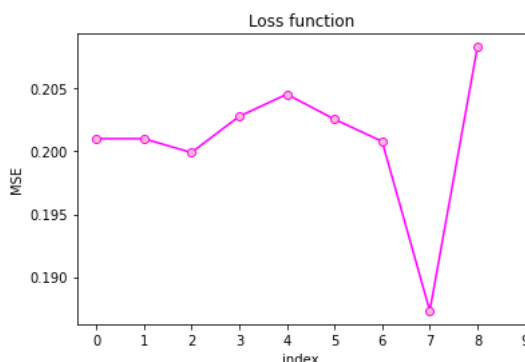
دانشکده علوم مهندسی، دانشگاه تهران

بهار ۱۴۰۱

## سوال ۱)

همان طور که با دیدن دیتا ها مشهود است این دیتا ویژگی های categorical دارد که باید تبدیل به ویژگی های عددی شود به همین دلیل با تخصیص اعداد ویژگی های categorical را به numerical تبدیل کرده ایم (label encoding)، سپس آن ها را به صورت رندوم به ده قسمت تقسیم کردیم و درخت تصمیم را هر بار با استفاده از یکی از ۹ دسته ساخته ایم و بعد با دسته ی دهم تست کرده ایم.

نمودار خطا های تشخیص مدل در زیر آمده است :

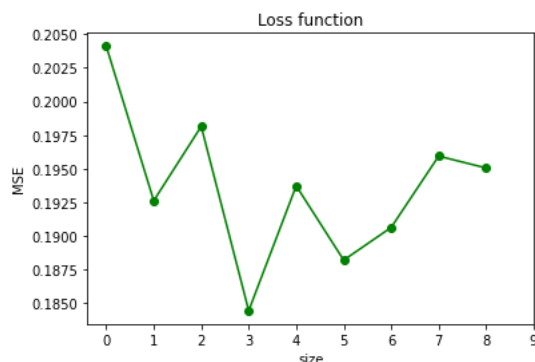


همان طور که در جدول خطا ها میبینیم کمترین خطا را مدلی دارد که در ایندکس هفت ساخته شده (با استفاده از دسته ی هشتم) است .

در کل چون در این بخش دیتا ها شافل شده اند و به صورت رندوم پخش شده اند تقریباً شبیه به هم دسته بندی شده اند هر درخت تصمیم چون با توجه به دیتا های مخصوص به دسته خود ساخته شده است یک مقداری با درخت های تصمیم دیگر تفاوت دارد ولی تعداد دیتا های موجود در هر دسته چون به نسبت زیادند، تقریباً تخمین قابل قبول و با خطای کمی به ما داده است.

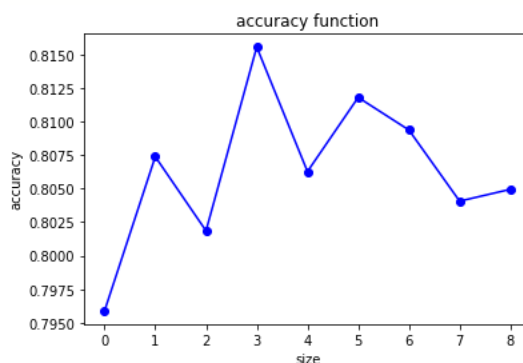
## سوال ۲)

در این سوال ما درخت را ابتدا با یک دسته آموزش داده سپس دو دسته و سپس سه و.. به همین ترتیب بر تعداد دسته های آموزش افزودیم جدول میزان خطای تشخیص درخت های ایجاد شده بر اساس تعداد دسته های آموزش در زیر آورده شده است :



برای مثال این تابع خطا در سائز سه ( استفاده از چهار دسته برای آموزش) کمترین خطا را داشته است و بعد از آن در سائز ۵ ( استفاده از ۶ دسته برای آموزش) خطای کمی داریم.  
تحلیل :

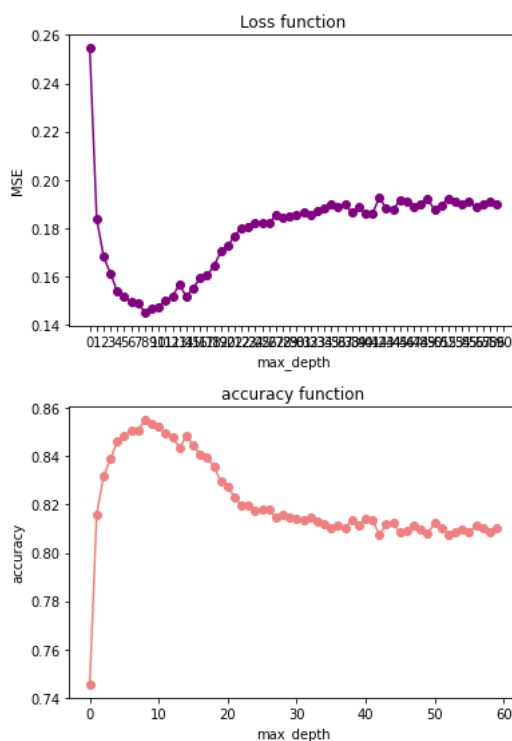
در این نمودار میبینیم که ابتدا با افزایش تعداد دسته ها برای آموزش، خطا کاهش یافته است و از یک جایی به بعد، به جای کاهش، اتفاقا افزایش خطا را داریم که نشان میدهد افزایش میزان دیتای آموزش تا حدی مناسب است ولی بیش از آن مقدار، باعث ایجاد **over fitting** میشود و مدل در داده های آموزش نتیجه ی خوبی دارد ولی روی داده های تست نتیجه ی مطلوبی نمیگیریم که این موضوع باعث افزایش خطا میشود.



در اینجا نمودار accuracy یا صحت تشخیص مدل را میبینیم که کاملاً با نمودار خطا هم خوانی دارد. بدیهی است که هرچه خطا تشخیص کمتر باشد صحت و درستی تشخیص مدل بالاتر است که میبینیم در سائز سه که کمترین خطا بوده بالا ترین صحت تشخیص را داریم.

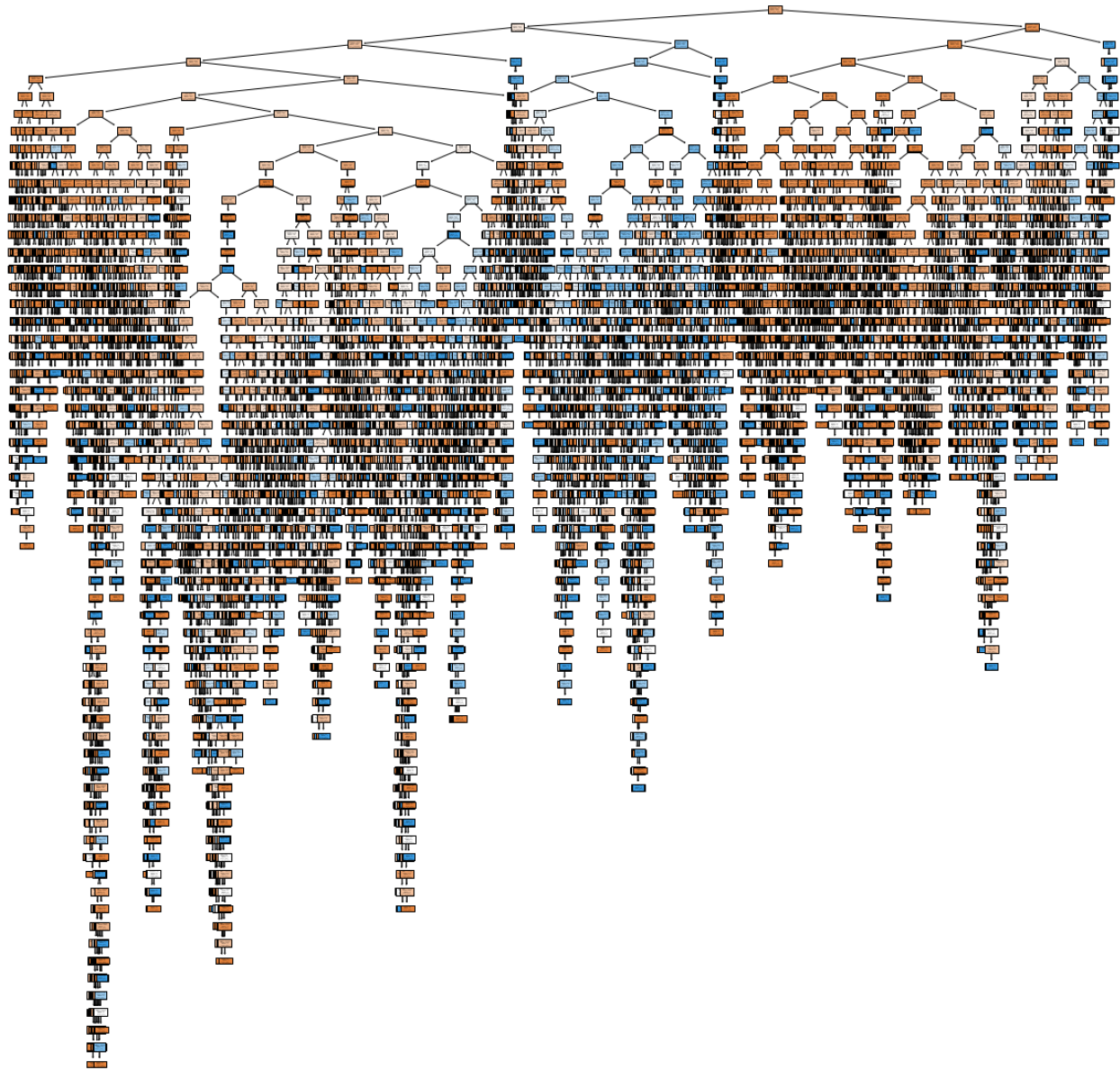
### سوال ۳)

درخت تصمیم با عمق بالا دچار over fitting میشود در نتیجه از روش هرس کردن یا به عبارتی کاهش عمق استفاده میکنیم تا از این مشکل جلوگیری شود. بیشترین عمق درخت ساخته شده ۶۱ است در نتیجه تمامی عمق ها از ۱ تا ۶۱ بررسی شده و همان طور که در نمودار مشخص است تا عمق ۸ خطا کاهش پیدا کرده ولی از این عمق به بعد از انجایی که over fitting رخ میدهد خطا افزایش پیدا میکند و accuracy کاهش میابد.

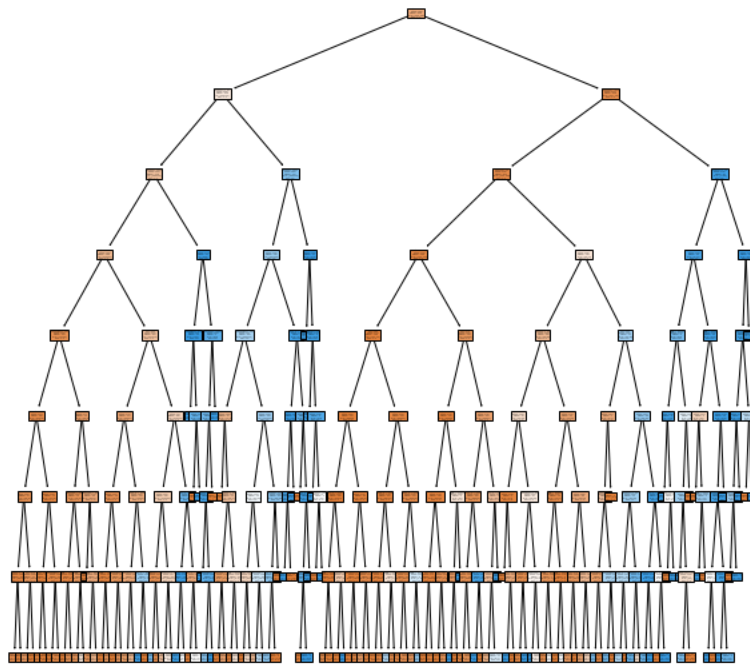


در ادامه تصویر درخت هرس نشده و هرس شده را برای مقایسه بهتر تاثیر هرس شونده آورده ایم.

در زیر درخت با عمق ۶۱ را میبینید:

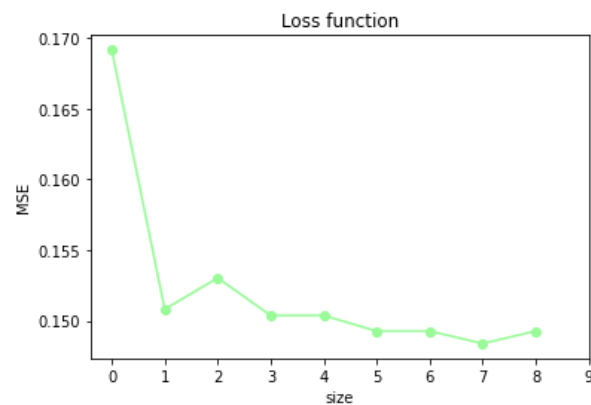


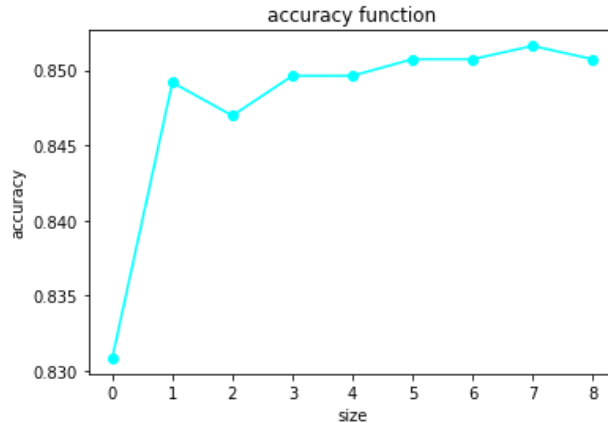
و در این بخش درخت هرس شده با عمق ۸ را میبینیم :



به وضوح با مقایسه دو درخت تاثیر هرس کردن را در تسریع رسیدن به درخت مطلوب را میبینیم.

در ادامه این سوال ما منحنی آموزش را برای درخت هرس شده رسم کردیم که در زیر آمده است :





حال که درخت تصمیم هرس شده است میبینیم که این منحنی اکنون روند بسیار منطقی تری را طی میکند.

#### سوال ۴)

در این قسمت ما برای انتخاب ویژگی مناسب برای هر گره از روش رندوم استفاده کردیم که میزان خطا و صحت آن در زیر آمده است :

accuracy: 0.8246351172047767

loss: 0.17536488279522336

این موضوع نشان میدهد که روش رندوم دقت کمتری نسبت به روش information gain دارد ( در روش رندوم accuracy حدود هشتاد و دو درصد بود که از هشتاد و شش درصد کمتر است). روش دیگر gain ratio میباشد که میزان خطا و صحت آن در زیر آمده است :

accuracy: 0.8538257408226448

loss: 0.14617425917735516

که نسبت به رندوم بهتر بود ولی نسبت به information gain دقت کمتری داشت.