

① الف) درخت تصمیم P_{ro} $interpretable$ (مهم رایج های درخت تصمیم آن است)

✓ P_{ro} اکثر نیاز به $preprocessing$ دارد. سرعت بالا و $computation$ کم.

$Cons$ → اگر به درستی $hyperparameter$ ها را ست نکت $overfit$ بسیار محتمل است.

✓ ترانای تک های که داده های آن پیچیده باشد ندارد. (برای کارهای ساده مناسب است)

کجک های عصبی P_{ro} طبق $universal\ classifier$ می توان گفت هر رابطه پیچیده و حتی دیتای پیچیده مانند منظم و عکس را نیز می تواند برداشت کند. (حتی غیر خطی)

$Cons$ $interpretable$ نیست. (قابلیت فهم و تفسیر نتایج را ندارد)

$Cons$ → نیاز به $computation$ بالا و داده $train$ زیاد دارد.

برای داده های ساده مانند عددی از $decision\ tree$ و برای داده های پیچیده مانند متن و عکس از $neural\ networks$ استفاده می کنیم.

ب) اگر از $f(x)=x$ برای $activation\ function$ استفاده کنیم در واقع گنگی نمی کند ما می خواهیم

$non-linearity$ به مدل اضافه کنیم و با این تابع فقط می تواند الگوهای خطی را تشخیص دهد. برای مثال

استفاده از $ReLU$ یا $Gelu$ امکان دارد برای ما انجام می دهد.

ج) این تابع $sub-differentiable$ است یعنی فقط در $x=0$ مشتق پذیر نیست که می توان

به صورت قراردادی مقدار را به مشتق این نقطه نسبت داده همچنین چون به نقطه تقاطع مشتق پذیر اند و

به طور معمول $ReLU$ خیلی خوب $converge$ می کند از این $activation\ func$ مناسب است.

د) در ترانای دیگری ممکن است به بعضی نقاط برسیم که کار ما را سخت کند مانند $local\ minima$ ها و $saddle\ points$ ها.

در بینم محلی ماکزیمم داریم اما به $global\ minima$ و باید به گونه ای عمل کنیم که گریز نکنیم. یا در $saddle\ points$ ها

باید حواسمان باشد که این نقاط از یک طرف گرادیان مثبت و از طرف دیگر گرادیان منفی دارند. یا در نقاطی که

$plateau$ می مانند که یعنی حدودا گرادیان صفر داریم و نمی توانیم از این نقاط محلی نجات بگیریم. بنابراین باید برای جلوگیری

وگرنه فاند در این نقاط روش های optimize معدی مانند Adam, AdaGrad, momentum وجود دارد

$$m_{t+1} = \beta_1 m_t + (1-\beta_1) \nabla_{\theta} L(\theta), m_0 = 0 \quad \text{روش adam}$$

$$v_{t+1} = \beta_2 v_t + (1-\beta_2) \nabla_{\theta} L(\theta)^2$$

$$\theta_j = \theta_j - \frac{\epsilon}{\sqrt{v_{t+1}} + \epsilon'} m_{t+1}$$

هر بار با استفاده از دو مقدار m و v
که انگار سرعت را به ساری کرده اند و واقع
طول step را تغییر می دهند آدیت می کنیم
این الگوریتم ترکیبی از momentum و RMSProp است.

$$g_{t+1} = g_t + \nabla_{\theta} L(\theta)^2 \quad \text{روش AdaGrad}$$

$$\theta_{t+1} = \theta_t - \epsilon \frac{\nabla_{\theta} L(\theta)}{\sqrt{g_t} + \epsilon'}$$

از g که مانند شتاب است استفاده
می شود تا یک step منطقی
تولید کند.