

1.1

از آن جا که بیش از گرفتن اطلاعات تابع را ارائه کرده ایم پس دنیای ندانیم تا بتوانیم در فضای فرضیه جست و جوی کنیم و  
بنابراین فضای فرضیه مان از ابتدا شامل همان یک فرضیه بوده است  $M=1$

1.2

$$\text{Hoeffding: } P[|E_{\text{out}} - E_{\text{in}}| > \epsilon] \leq 2M e^{-2\epsilon^2 N}$$

$$\rightarrow P[|E_{\text{out}} - E_{\text{in}}| > \frac{2}{100}] \leq 2 \times 1 \times e^{-2 \times 4 \times 10^{-4} \times 10^4} = 2e^{-8}$$

بنابراین طبق هافدینگ به احتمال حداقل  $1 - 2e^{-8}$  خطا کمتر از دو درصد است

1.3

دلیل این اتفاق آن است که دنیایی که بانک در اختیار ما گذاشته نتوانسته به خوبی تابع ما را ارزیابی کند. به این معنا که توزیع داده های  
زمان آموزش و تست متفاوت بوده و distribution shift داریم. به عنوان مثال یکی از عواملی که باعث می شود انتخاب حدسی ما ساده تر  
شود آن است که دنیایی که در اختیار ما قرار گرفته تماماً متعلق به افرادی است که بانک بر اساس روش های پیشین خود آنها را تأیید کرده  
است و بنابراین دنیا به سمتی با بایس است که افراد تأیید شوند.

حال در زمان تست اما داده ها از توزیع دیگری می آیند و با مسئله out of distribution generalization رو به رو می شویم و مدل عملکرد ضعیفی را بروز می دهد

⊗ در حقیقت باند هافدینگ بر فرض iid بودن توزیع های train و test استوار است

⊗ به طور کلی وقتی به این صورت با مسئله generalization رو به رو می شویم یا تعداد سیمپل ها نسبت به مدل استفاده شده کم بوده یا توزیع داده های آموزش و تست  
متفاوت است

1.4

یک راه آن است که سعی کنیم توزیع داده های آموزش و تست را یکی کنیم. هنگامی که متوجه شویم دنیای آموزش به چه سمتی با بایس شده است  
می توان سیمپل هایی که توزیع را به آن سمت برده اند را تا حدی کاهش داد تا این کیفیت از بین نبرد. به عنوان مثال اگر متوجه شدیم که دنیا به سمت  
افراد تأیید شده با بایس شده است، می توان تعدادی از آنها را از دنیا حذف کرد.

البته اضافه کردن دنیایی که در زمان آموزش دیده نشده ولی در زمان تست وجود داشته کار سخت تری است که تنها در شرایط خاص قابل انجام است

پس از اینکه توزیع ها مشابه شوند، باند هافدینگ در خصوص generalization همان مقدار قبلی را می دهد اما این بار عملکرد مدل دیگر perfect نخواهد بود

و در همان زمان evaluation متوجه دقت کم آن می شویم

2.1

در حوزه یادگیری ماشین منظور از درخت تصمیم تفاوت، تفاوت از نظر لیدل زدن سیمپل ها است  
 بنابراین اگر  $m$  ویژگی داشته باشیم،  $2^m$  حالت برای سیمپل ها متصور است و چون دسته بندی binary است  
 پس هر کدام از این سیمپل ها نیز می توانند 2 لیدل بگیرند  $\leftarrow 2^{2^m}$  درخت تصمیم متفاوت داریم

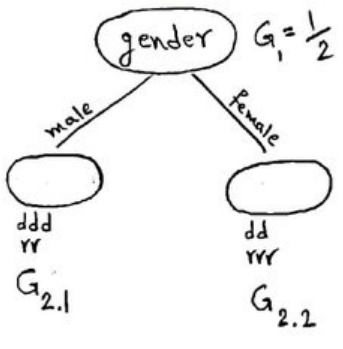
2.2

خیر درخت بهینه نیست  
 انتخاب مریضانه تمامی تواند به ما بگوید که در هر مرحله (به تنهایی)، کدام split بهتر است اما ممکن است یک split در حال حاضر  
 خوب نباشد اما اگر آن را انجام دهیم در ادامه به جواب بهینه ای برسیم. به همین صورت ممکن است یک split که در حال حاضر بهترین gain را دارد  
 در ادامه بسیار بد باشد.  
 با اینکه این الگوریتم ساختار درخت بهینه را به دست نمی دهد اما به دلیل ساده بودن و سریع بودن اجرای آن، همچنان به جواب های مناسبی می رسد

2.3.a

در مرحله اول 3 گزینه برای split داریم که gini gain را برای هر کدام به دست می آوریم:

③ سبیل با لیل جمهوری فدرال ۲ و سبیل بالیل دموکرات را d می نامیم



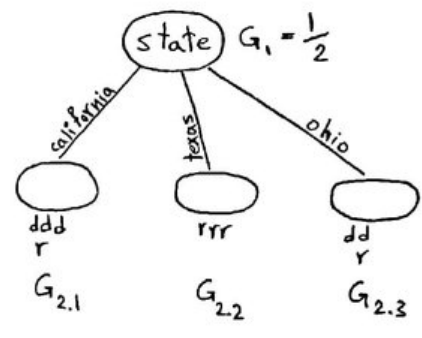
$$G_{2.1} = \frac{3}{5} \times \frac{2}{5} \times 2 = \frac{12}{25}$$

$$G_{2.2} = \frac{2}{5} \times \frac{3}{5} \times 2 = \frac{12}{25}$$

$$\text{avg} \rightarrow \frac{12}{25}$$

$$\text{gini gain} = \frac{1}{2} - \frac{12}{25} = \frac{1}{50}$$

gini gain for "race" is also  $\frac{1}{50}$



$$G_{2.1} = \frac{3}{4} \times \frac{1}{4} \times 2 = \frac{3}{8}$$

$$G_{2.2} = \frac{3}{3} \times 0 \times 2 = 0$$

$$G_{2.3} = \frac{2}{3} \times \frac{1}{3} \times 2 = \frac{4}{9}$$

$$\text{avg} \rightarrow \frac{4}{10} \left( \frac{3}{8} \right) + \frac{3}{10} \left( \frac{4}{9} \right) = \frac{17}{60}$$

$$\text{gini gain} = \frac{1}{2} - \frac{17}{60} = \frac{13}{60} \checkmark$$

پس ابتدا بر اساس state تقسیم می کنیم.

✓  $\text{gain} = \frac{3}{8} - \frac{1}{4} = \frac{1}{8}$  اگر سبیل های کالیفرنیا را بر اساس جنسیت تقسیم کنیم 

male	female
ddd	d
rr	r

 خواص داشت ✓

$\text{gain} = \frac{3}{8} - \frac{1}{3} = \frac{1}{24}$  ~ ~ 

white	black
ddd	d
r	r

 ~ ~ ~ ~ ~

✓  $\text{gain} = \frac{4}{9} - \frac{1}{3} = \frac{1}{9}$  اگر سبیل های اوهایو را بر اساس جنسیت تقسیم کنیم 

male	female
d	d
r	r

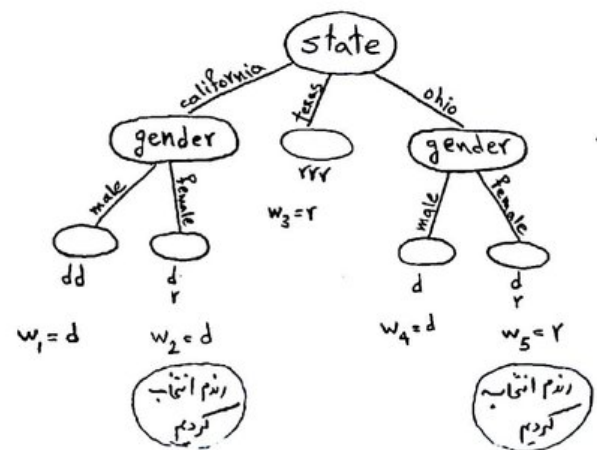
 خواص داشت ✓

$\text{gain} = \frac{4}{9} - \frac{4}{9} = 0$  ~ ~ 

white	black
ddd	dd
r	r

 ~ ~ ~ ~ ~

درفت نهایی با ارتفاع 2



2.3.6

طبق بخش قبل برای اولین انتخاب داریم:

• جنسیت  $gain = \frac{1}{50}$

• نژاد  $gain = \frac{1}{50}$

• مکان کالیفرنیا  $gain = \frac{1}{2} - \left( \frac{4}{10} \left( \frac{3}{8} \right) + \frac{6}{10} \left( \frac{4}{9} \right) \right) = \frac{1}{12}$

california: ddd, r  
others: dd, rrrr

✓ مکان تگزاس  $gain = \frac{1}{2} - \left( \frac{7}{10} \left( \frac{20}{49} \right) \right) = \frac{3}{14}$

texas: rrr  
others: dddd, rr

• مکان اوماها  $gain = \frac{1}{2} - \left( \frac{3}{10} \left( \frac{4}{9} \right) + \frac{7}{10} \left( \frac{24}{49} \right) \right) = \frac{1}{42}$

ohio: dd, r  
others: ddd, rrrr

پس ابتدا بر اساس تگزاس تقسیم می کنیم. سپس برای تقسیم texas = no داریم:

✓ جنسیت  $gain = \frac{20}{49} - \left( \frac{4}{7} \left( \frac{1}{2} \right) \right) = \frac{6}{49}$

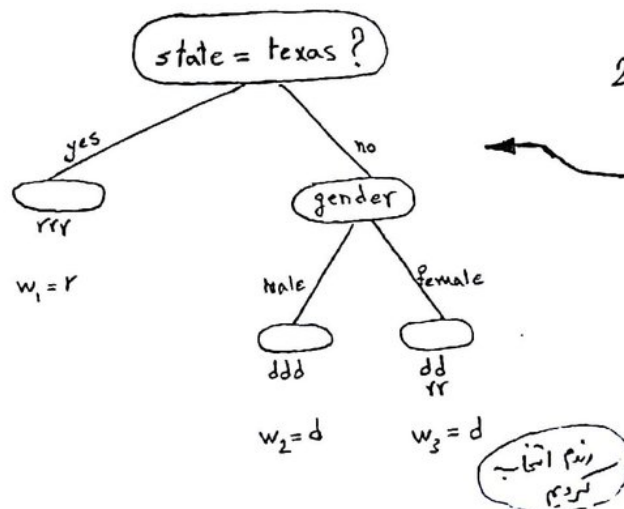
male: ddd  
female: dd, rr

• نژاد  $gain = \frac{20}{49} - \left( \frac{3}{7} \left( \frac{4}{9} \right) + \frac{4}{7} \left( \frac{3}{8} \right) \right) = \frac{1}{294}$

white: dd, r  
black: dd, r

• مکان  $gain = \frac{20}{49} - \left( \frac{4}{7} \left( \frac{3}{8} \right) + \frac{3}{7} \left( \frac{4}{9} \right) \right) = \frac{1}{294}$

california: ddd, r  
ohio: dd, r





## 2.3.c

به طور کلی می توان به تفاوت های زیر اشاره کرد

- در روش multiway split درخت کوچکتری به دست می آید زیرا درجه هر نود بیشتر است

اما در روش binary split عمق درخت بیشتر می شود

اگر عمق درخت را نامحدود فرض کنیم هر دو روش دنیای آموزش را کامل مطابق لیل ها دسته بندی می کنند اما واضح است که با محدودیت حداکثر عمق برای روش multiway split به تعداد کمی بهتر عمل می کند

- روش multiway split درخت های قابل تفسیر تری به دست می دهد زیرا dummy node های کمتری ایجاد می کند

- روش binary split به دلیل ارتفاع بیشتر درخت و dummy node های بیشتر، معایب های بیشتری برای دسته بندی نیاز دارد و از این نظر multiway split می تواند efficient تر باشد

- در روش multiway split ریسک overfitting بیشتر است. مخصوصاً هنگامی که دانش متادیر Feature ها بیشتر باشد

- محاسبه information gain به ازای هر اتامی در روش binary split ساده تر است و بنابراین درایند آموزش سریع تری دارد

به طور کلی انتخاب روش مناسب کاملاً وابسته به نوع دیتا و شرایط مسئله است و نمی توان گفت که یکی بهتر از دیگری است

## 2.4

اگر به اندازه کافی عمق درخت تقسیم را زیاد کنیم می توانیم به ازای هر کدام از سمپل های دنیای آموزش برگ ایجاد کند (در بدترین حالت)

و بنابراین ساختار آن به صورتی است که به راحتی می تواند overfit شود و واریانس بالایی دارد

به طور کلی این مشکل در درخت ها را توسط ensemble learning بهبود می دهند. ایده ی این روش آن است که مدل های مستقل (نسبتاً مستقل) ای

آموزش دهیم و سپس برای prediction جواب مدل ها را aggregate کنیم (مثلاً میانگین بگیریم یا majority vote)

طبق قضیه حد مرکزی می دانیم که اگر  $n$  مقیّر تصادفی iid را میانگین بگیریم واریانس حاصل  $\frac{1}{n}$  می شود و از آنجا که  $E_{out} = bias + variance$  پس عملکرد مدل به طور کلی نیز بهتر می شود

چگونگی تصادفی یک روش ensemble کردن درخت هاست. به این صورت که تصادفی درخت را روی بخش های مختلف داده و با پارامترهای مختلف آموزش می دهد و سعی می کند مدل های نسبتاً مستقلی ایجاد کند. در نهایت نیز بسته به مسئله در زمان prediction جواب ها را aggregate می شوند

⊗ واضح است که بایاس افزایش پیدا نمی کند زیرا  $bias = E_D [E_{in}(h^{(D)}(x))]$  و هیچ کدام از مدل ها  $E_{in}$  شان بهتر از حالت تنگی نمی شود بنابراین با حفظ دقت درخت، واریانس کاهش می یابد

3.1

زیرا لایه ها از هم جدا هستند  $H_T(x^{(i)}) y^{(i)} < 0$

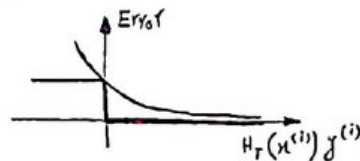
$$E_T = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\text{sign}(H_T(x^{(i)})) \neq y^{(i)})$$

برای  $\text{sign}(H_T(x))$  به صورت زیر 0/1 loss تعریف می شود

با توجه به اینکه به ازای هر سیگنال داریم:  $\mathbb{I}(H_T(x^{(i)}) y^{(i)} < 0) \leq e^{-H_T(x^{(i)}) y^{(i)}}$

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}(H_T(x^{(i)}) y^{(i)} < 0) \leq \frac{1}{N} \sum_{i=1}^N e^{-H_T(x^{(i)}) y^{(i)}}$$

این مقدار را  $E$  می نامیم



و بنابراین تابع هزینه ما یک upper bound است که اگر آن را min کنیم 0/1 loss نیز کمینه می شود. پس در ادامه کار، از آن استفاده می کنیم

3.2

$$D_{t+1}(i) = D_t(i) \frac{e^{-\alpha_t y^{(i)} h_t(x^{(i)})}}{Z_t}$$

رابطه بازگشتی رو به رو و برای  $D_t(i)$  داریم:

$$\begin{aligned} \xrightarrow{\text{جایگزینی}} D_{t+1}(i) &= D_1(i) \times \frac{e^{-\alpha_1 y^{(i)} h_1(x^{(i)})}}{Z_1} \times \frac{e^{-\alpha_2 y^{(i)} h_2(x^{(i)})}}{Z_2} \times \dots \times \frac{e^{-\alpha_t y^{(i)} h_t(x^{(i)})}}{Z_t} \\ &= \frac{1}{N} \times \frac{e^{-y^{(i)} \left( \sum_{t=1}^t \alpha_t h_t(x^{(i)}) \right)}}{\prod_{t=1}^t Z_t} \end{aligned}$$

3.3

طبق تعریف داریم که:

$$H_T(x^{(i)}) = \sum_{t=1}^T \alpha_t h_t(x^{(i)})$$

$$\xrightarrow[\text{از بخش 3.1}]{\text{جایگزینی در E به دست آمده}} E = \frac{1}{N} \sum_{i=1}^N e^{-y^{(i)} \sum_{t=1}^T \alpha_t h_t(x^{(i)})} = \sum_{i=1}^N \underbrace{\frac{1}{N} e^{-y^{(i)} \sum_{t=1}^T \alpha_t h_t(x^{(i)})}}_{(*)}$$

3.4

$$E = \sum_{i=1}^N D_{T+1}(i) \prod_{t=1}^T Z_t$$

طبق بخش 3.2  $(*) = D_{T+1}(i) \prod_{t=1}^T Z_t$  پس جایگزینی آن داریم:

$$\frac{\text{عبارت دوم به i وابسته نیست}}{\prod_{t=1}^T Z_t} \times \underbrace{\sum_{i=1}^N D_{T+1}(i)}_1 = \prod_{t=1}^T Z_t$$

تمام  $D_t(i)$  ها توزیع مای هستند که به سیگنال ها (زن می دهند)  
پس طبق تعریف  $\sum_{i=1}^N D_t(i) = 1$

3.5

$$\begin{aligned}
 z_t &= \sum_{i=1}^N D_t(i) e^{-\alpha_t y^{(i)} h_t(x^{(i)})} \\
 &= \sum_{i=1}^N D_t(i) e^{-\alpha_t} \mathbb{I}(h_t(x^{(i)}) = y^{(i)}) + \sum_{i=1}^N D_t(i) e^{\alpha_t} \mathbb{I}(h_t(x^{(i)}) \neq y^{(i)}) \\
 &= e^{-\alpha_t} \underbrace{\sum_{i=1}^N D_t(i) \mathbb{I}(h_t(x^{(i)}) = y^{(i)})}_{1-\varepsilon_t} + e^{\alpha_t} \underbrace{\sum_{i=1}^N D_t(i) \mathbb{I}(h_t(x^{(i)}) \neq y^{(i)})}_{\varepsilon_t} \\
 &= (1-\varepsilon_t) e^{-\alpha_t} + \varepsilon_t e^{\alpha_t}
 \end{aligned}$$

3.6

مقدار  $E$  تا کدام  $t$  را  $E_t$  می نامیم و داریم

$$\left. \begin{aligned}
 \text{3.4} \quad \text{طبق بخش 3.4} \quad E_t &= \prod_{t'=1}^t z_{t'} \\
 \text{3.5} \quad \text{طبق بخش 3.5} \quad z_t &= (1-\varepsilon_t) e^{-\alpha_t} + \varepsilon_t e^{\alpha_t}
 \end{aligned} \right\} \rightarrow E_t = \prod_{t'=1}^t ((1-\varepsilon_{t'}) e^{-\alpha_{t'}} + \varepsilon_{t'} e^{\alpha_{t'}})$$

حال باید مشتق  $E_t$  را بوسیله  $\alpha_t$  حساب کنیم تا مقدار بهینه را برای آن به دست آوریم  
 باید به اینک در هر مرحله به صورت greedy عمل می کنیم و وقتی  $h_t$  و  $\alpha_t$  را به دست می آوریم  $h_j$  و  $\alpha_j$  ثابت خواهند بود پس  $j < t$

$$E_t = \underbrace{\left( \prod_{t'=1}^{t-1} ((1-\varepsilon_{t'}) e^{-\alpha_{t'}} + \varepsilon_{t'} e^{\alpha_{t'}}) \right)}_C \times \left( (1-\varepsilon_t) e^{-\alpha_t} + \varepsilon_t e^{\alpha_t} \right)$$

$$\rightarrow \frac{\partial E_t}{\partial \alpha_t} = C (\varepsilon_t - 1) e^{-\alpha_t} + C \varepsilon_t e^{\alpha_t} = 0 \xrightarrow{\times e^{\alpha_t}} \varepsilon_t - 1 + \varepsilon_t e^{2\alpha_t} = 0$$

$$\rightarrow \alpha_t = \frac{1}{2} \ln \left( \frac{1-\varepsilon_t}{\varepsilon_t} \right)$$

بنابراین به ازای هر  $t$  نشان دادیم که به خاطر انتخاب greedy و ثابت بودن  
 کلام های قبلی مقدار  $\alpha_t$  به این صورت به دست می آید