

به نام خدا
دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

گزارش تمرین دوم درس یادگیری ماشین

پاسخ مسئله‌های تشریحی

استاد درس: دکتر احسان ناظر فرد

دانشجو: فاطمه غلامزاده

۹۹۱۳۱۰۰۳

نیم سال دوم ۱۳۹۹-۱۴۰۰

سوال ۱)

از بیش برآزش شدن جلوگیری می‌کند. زمانی که عمق درخت زیاد می‌شود لازم است هرس کردن انجام شود. با هرس کردن شاخه‌هایی که داده‌های کمی دارند و به یک پدر متعلق هستند را یکی می‌کند و برای هر کلاس احتمال آن‌را نشان می‌دهد.

روش‌های جلوگیری از بیش برآزش در درخت تصمیم:

۱- Pre-pruning : درختی با تعداد شاخه‌های کمتری از درختی که باید تولید می‌شد ساخته می‌شود. یعنی

پیش از آن که درخت تمام داده‌های آموزش را دسته‌بندی کند رشد آن را متوقف می‌کنیم.

۲- Post-pruning: درخت به طور کامل ساخته می‌شود و بعد قسمت‌هایی از آن حذف می‌شود. یعنی درخت

به طور کامل داده‌ها آموزش را دسته‌بندی می‌کند و بعد هرس می‌شود.

۳- فیچرهای غیرمرتبط و نامناسب می‌توانند موجب بیش برآزش در درخت تصمیم شوند بنابراین حذف این

فیچرها می‌تواند از بیش برآزش جلوگیری کند.

۴- افزایش تعداد داده‌های آموزش.

سوال ۲)

الف) ماتریس در هم ریختگی به صورت زیر است:

=== Confusion Matrix ===

```
a  b  <-- classified as
14  6  |  a = bad
 9 28  |  b = good
```

FN = 6 , TP = 14, TN = 28, FP = 9

Precision = $TP / (TP + FP) = 14 / 23 = 0.6$

Recall = $TP / (TP + FN) = 14 / 20 = 0.7$

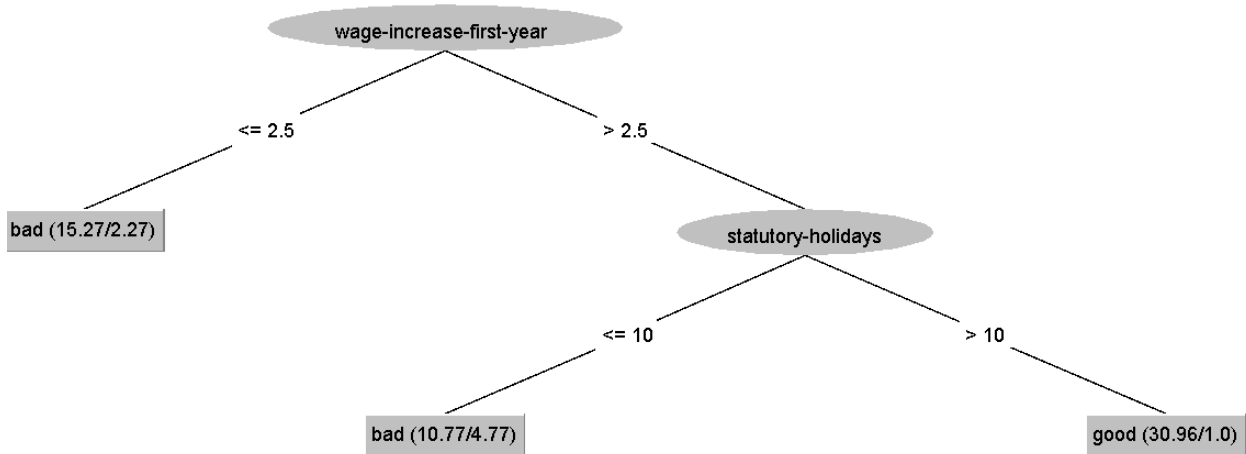
$Accuracy = \frac{\#TP + \#TN}{\#P + \#N} = \frac{\#TP + \#TN}{\#TP + \#FN + \#TN + \#FP}$

⇒ **Accuracy** = $14 + 28 / 57 = 44 / 57 = 0.77$

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$$

$$\Rightarrow \text{F1-measure} = 2 * (0.6) * (0.7) / (1.3) = 0.64$$

درخت تصمیم ساخته شده به صورت زیر است:



داده‌ی مورد سوال، در کلاس b یعنی good قرار می‌گیرد. زیرا ابتدا فیچر wage-increase-first-year مورد بررسی قرار می‌گیرد که مقدار آن ۳ است پس به شاخه سمت راست می‌رویم (بزرگتر از ۲,۵) بعد فیچر statutory-holidays مورد بررسی قرار می‌گیرد که برای داده‌ی مورد سوال مقدار آن ۱۲ است پس به شاخه سمت راست می‌رویم که کلاس این داده مشخص می‌شود و good است.

ب) پارامتر unpruned این مورد را تعیین می‌کند که درخت تصمیم هرس بشود یا خیر، اگر مقدار آن False باشد درخت ساخته شده هرس می‌شود و اگر مقدار آن True باشد درخت هرس نمی‌شود. در واقع در حالتی که عمق درخت زیاد شود و به اصطلاح overfitting رخ دهد، با هرس کردن درخت بعضی از شاخه‌ها حذف شده و در ریشه ی مربوطه احتمال آنها نوشته می‌شود. انتظار می‌رود که با فعال کردن این گزینه عمق درخت بیشتر شود.

در این قسمت پارامتر unpruned به true مقداردهی شد. ماتریس درهم‌ریختگی به صورت زیر در آمد:

=== Confusion Matrix ===

```

a  b  <-- classified as
14  6 |  a = bad
 6 31 |  b = good
  
```

FN = 6 , TP = 14, TN = 31, FP = 6

$$\text{Precision} = TP/(TP+FP) = 14/20 = 0.7$$

$$\text{Recall} = TP/(TP+FN) = 14/20 = 0.7$$

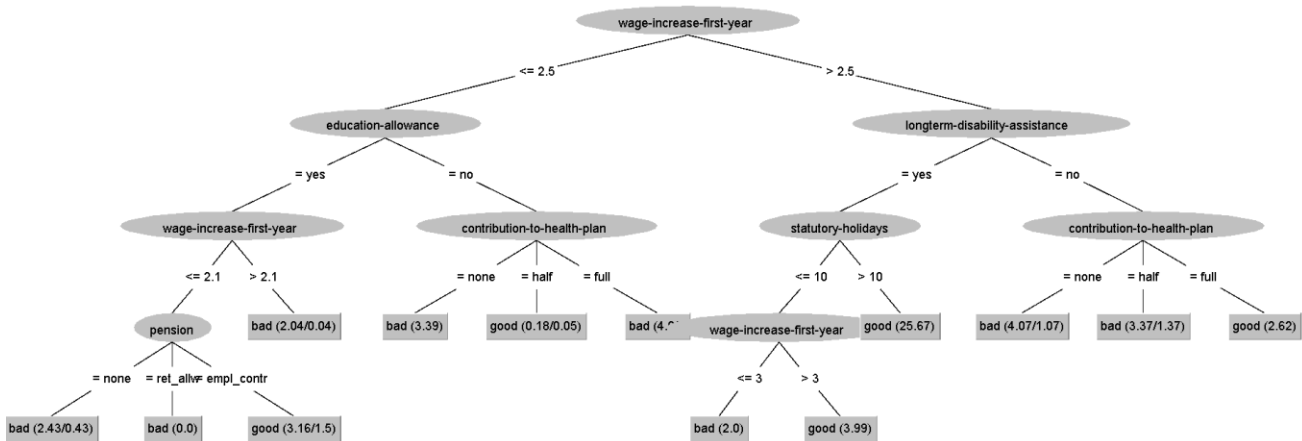
$$\text{Accuracy} = \frac{\#TP + \#TN}{\#P + \#N} = \frac{\#TP + \#TN}{\#TP + \#FN + \#TN + \#FP}$$

$$\Rightarrow \text{Accuracy} = 14 + 28 / 57 = 45 / 57 = 0.78$$

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2TP}{2TP + FP + FN}$$

$$\text{F1-measure} = 2 * (0.7) * (0.7) / (1.4) = 0.7$$

درخت تصمیم ساخته شده به این صورت است:



برای داده مورد سوال ابتدا فیچر wage-increase-first-year بررسی می‌شود و چون $3 < 2.5$ است به شاخه سمت راست می‌رویم. سپس فیچر longterm-disability-assistant بررسی می‌شود و چون yes است به شاخه سمت چپ می‌رویم. برای فیچر statutory-holidays چون مقدارش ۱۲ است و از ۱۰ بیشتر است به شاخه سمت راست می‌رویم و کلاس داده برابر b می‌شود.

تفاوت درخت‌ها: در این قسمت چون پارامتر unpruned را true کردیم درخت رسم شده هرس نشده و عمق و عرض آن بیشتر است و فیچرهای بیشتری در آن مورد بررسی قرار می‌گیرد. با توجه به دقت‌های به دست آمده مشاهده می‌شود که با هرس نکردن درخت، نسبت به حالت هرس شده دقت‌ها افزایش یافته‌اند.

(سوال ۳)

بله افزایش می‌یابد. هر چقدر که ابعاد داده‌ها بیشتر باشد این فضای با ابعاد بالا با چگالی و فشردگی کمتری از نقاط داده‌های آموزش روبه روست و نیاز داریم تا حجم زیادی از فضا را جستجو کنیم تا همسایه‌های داده تست را پیدا کنیم. هم‌چنین با افزایش ابعاد داده‌ها، فاصله دوبره دوی میان نقاط نیز افزایش پیدا می‌کند. در این صورت همسایه‌های یک داده تست ممکن است به حدی از آن دور باشند که در واقع هیچ وجه اشتراکی با این داده نداشته باشند و این امر موجب خطا در دسته بندی می‌شود. به این مسئله، **Curse of dimensionality** گفته می‌شود.

فرض کنیم که داده‌های آموزش در یک فضای D بعدی به طور یکنواخت در مکعب واحد پراکنده شده‌اند و هم‌چنین فرض کنیم که می‌خواهیم چگالی برچسب‌های کلاس‌ها را با رشد یک مکعب واحد در اطراف داده تست تخمین بزنیم. یعنی مکعب آن قدر رشد می‌کند تا تعداد داده مورد نظر ما در آن قرار بگیرد و فرض کنیم این تعداد مورد نظر f باشد. امید ریاضی طول ضلع این مکعب از رابطه زیر به دست می‌آید:

$$e_D(f) = f^{\frac{1}{D}}$$

این رابطه نشان می‌دهد که سائز ضلع مکعب مورد نیاز برای اینکه تعداد مشخصی داده آموزشی در آن قرار بگیرد، با افزایش ابعاد داده‌ها به شدت افزایش پیدا می‌کند. شکل زیر این مورد را نشان می‌دهد:

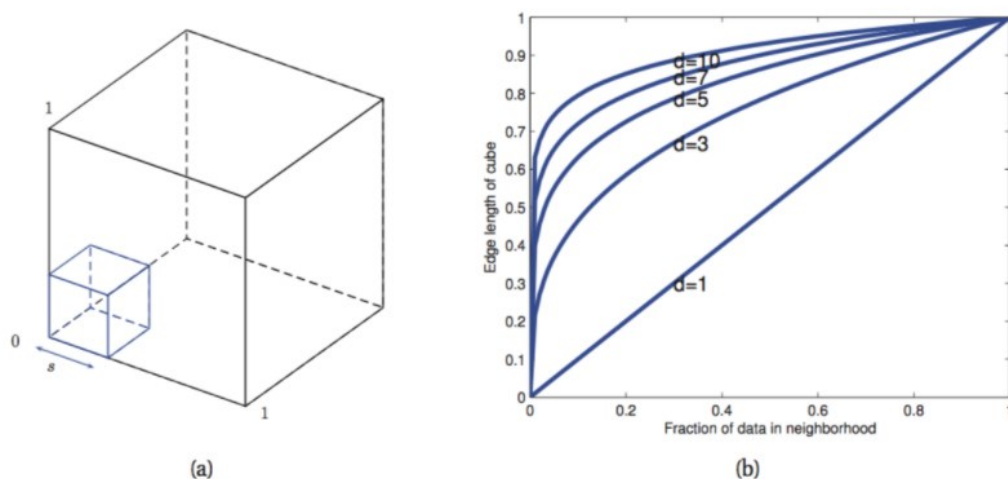


Figure 1.16 Illustration of the curse of dimensionality. (a) We embed a small cube of side s inside a larger unit cube. (b) We plot the edge length of a cube needed to cover a given volume of the unit cube as a function of the number of dimensions. Based on Figure 2.6 from (Hastie et al. 2009). Figure generated by `curseDimensionality`.

سوال (۴)

مدل های مولد توزیع واقعی داده ها را مدل می کنند و مدل ها تمایزگر یک مرز تصمیم را مدل میکنند. در الگوریتم های مولد دسته بند مقادیر $p(Y|X)$ را محاسبه می کند. برای محاسبه ی این احتمال $P(X|Y)P(Y)$ را بدست می آورد. با این محاسبات الگوریتم توزیع داده ها را فرا می گیرد و می تواند خودش داده ها را دوباره تولید کند. مدل تمایزگر (Discriminative) برای دسته بندی تنها یک مرز تصمیم را می آموزد و بر اساس آن دسته بندی را انجام می دهد. مدل مولد احتمال توزیع توام $p(x,y)$ را می آموزد. این مدل احتمال شرطی را با کمک قضیه بیز پیش بینی می کند اما مدل تمایزگر احتمال توزیع شرطی $p(y|x)$ را فرا میگیرد.

نمونه هایی از دسته بندهایی که مولد هستند:

- Naïve Bayes
- Bayesian networks
- Markov random fields
- Hidden Markov Models (HMM)

نمونه هایی از دسته بند های تمایزگر:

- Logistic regression
- Scalar Vector Machine
- Traditional neural networks
- Nearest neighbour
- Conditional Random Fields (CRF)s

سوال (۵)

این دو دسته بندی کننده از دو نوع مختلف دسته بندی کننده های generative و discriminative هستند. بیز ساده از نوع generative؛ یعنی مولد است. در الگوریتم های مولد دسته بند مقادیر $p(Y|X)$ را محاسبه می کند. برای محاسبه ی این احتمال $P(X|Y)P(Y)$ را بدست می آورد. با این محاسبات الگوریتم توزیع داده ها را فرا می گیرد و می تواند خودش داده ها را دوباره تولید کند. همانطور که میدانیم، الگوریتم بیز ساده این احتمال ها را بدست می آورد.

رگرسیون لاجستیک از نوع discriminative؛ یعنی تمایزدهنده است. به این معنا که برای دسته بندی تنها یک مرز تصمیم را می آموزد و بر اساس آن دسته بندی را انجام می دهد. در بیز ساده ابتدا احتمال $p(Y|X)$ به شکل زیر بازنویسی می شود. سپس مقادیر $P(X|Y)$ از داده های آموزشی محاسبه می شود.

$$\mathbf{X} = \langle X_1, X_2, \dots, X_n \rangle, \quad X_i \perp\!\!\!\perp X_j | Y \quad (s.t: i, j = 1, 2, \dots, n \text{ \& } i \neq j)$$

$$y_{MLE} = \operatorname{argmax}_{y \in \{+, -\}} p(\mathbf{X}, Y = y) = \operatorname{argmax}_{y \in \{+, -\}} p(\mathbf{X} | Y = y) p(Y = y)$$

$$= \operatorname{argmax}_{y \in \{+, -\}} p(X_1 | Y = y) p(X_2 | Y = y) \dots p(X_n | Y = y) p(Y = y)$$

$$= \operatorname{argmax}_{y \in \{+, -\}} p(Y = y) \prod_{i=1}^n p(X_i = x_i | Y = y)$$

رگرسیون لاجستیک برای دسته‌بندی مقدار $p(Y|\mathbf{X})$ را مستقیماً محاسبه می‌کند و بر اساس مقدار آن دسته‌بندی را انجام می‌دهد. در ادامه توابعی که احتمال از آنها محاسبه میشود برای حالت دو کلاسه آورده شده است. رگرسیون لاجستیک از تابع سیگموئید برای محاسبه کلاس مثبت استفاده می‌کند و ورودی این تابع یک ترکیب خطی از ویژگی‌های مختلف داده‌هاست. پس در واقع لاجستیک رگرسیون یک دسته‌بندی کننده خطی است. وزنهای این ترکیب خطی به کمک داده‌های آموزشی و روش‌های بهینه‌سازی بدست می‌آیند. (البته لاجستیک غیرخطی نیز وجود دارد که ترکیب غیر خطی ویژگی‌ها به عنوان ورودی تابع سیگموئید در نظر گرفته می‌شود.)

$$P(Y = 1 | \mathbf{X}) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

$$P(Y = 0 | \mathbf{X}) = \frac{\exp(w_0 + \sum_{i=1}^n w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

مقایسه عملکرد مدل‌های لاجستیک و بیز ساده‌ی گوسی:

بیز ساده بایاس زیادی نسبت به داده‌های آموزش دارد. رگرسیون لاجستیک واریانس زیادی دارد که می‌توان میزان آن را با تغییر در تابع بهینه‌سازی بهبود بخشید. اگر تعداد داده‌های آموزش به سمت بی‌نهایت برود و بیز ساده مفروضاتی مانند برابر بودن پراکندگی داده‌ها در کلاس‌های مختلف داشته باشد، نتایج دو مدل همگرا است یعنی هر دو مدل یک نتیجه را می‌دهند. اما در غیر این صورت رگرسیون لاجستیک به دلیل اینکه پارامترهای خود را براساس تمام داده‌ها محاسبه می‌کند دقت بیشتری نسبت به بیز ساده دارد.

$$P(A) = 0.12 \quad P(A') = 0.18$$

سوال ۹

$$P(C) = P(C|A)P(A) + P(C|\bar{A})P(\bar{A}) = (0.13 \times 0.12) + (0.15 \times 0.18) = 0.124$$

$$P(\bar{C}) = 1 - P(C) = 1 - 0.124 = 0.876$$

$$P(D|B=T) = P(D|B,C)P(C) + P(D|B,\bar{C})P(\bar{C})$$

$$= (0.15 \times 0.124) + (0.15 \times 0.876) = 0.141$$

Scanned with CamScanner

سوال ۷)

شکل ۳:

برچسب داده ی تست A است. میانگین داده های A و B در راستای X_1 و X_2 تقریباً با هم برابر است. اما واریانس داده های A در هر دو راستا از B بیشتر است. هم چنین احتمال اولیه ی تعلق به کلاس A بیشتر است چون نسبت تعداد داده های A به کل داده ها از نسبت تعداد داده های B به کل داده ها بیشتر است. بنابراین ۲ مورد از ۳ موردی که در هم ضرب می شوند برای A بزرگتر و مورد سوم هم تقریباً با هم برابر است و احتمال تعلق داده تست به کلاس A بیشتر خواهد شد.

شکل ۴:

برچسب داده ی تست B است. چون در راستای X_1 انحراف معیار B از A بیشتر است و احتمال متعلق بودن داده ی تست به کلاس B بیشتر مساوی از احتمال متعلق بودن آن به کلاس A است. از طرف دیگر در راستای X_2 نیز احتمال متعلق بودن داده ی تست به کلاس B بیشتر است. چون اگر فرض کنیم میانگین هر دو کلاس در یک نقطه است و داده ی تست روی میانگین قرار دارد، چون انحراف معیار A بیشتر از B است پس ارتفاع قله ی B بیشتر از A است و در نتیجه احتمال متعلق بودن داده ی تست به کلاس B بیشتر است. و دلیل آخر اینکه احتمال اولیه ی کلاس B نیز

بیشتر از کلاس A است. به همین سه دلیل زمانی که احتمال‌ها در هم ضرب می‌شوند حاصل نهایی برای کلاس B بیشتر می‌شود و داده‌ی تست برچسب کلاس B را می‌گیرد.

سوال ۸

عمل **regularization** موجب می‌شود که بایاس مدل افزایش یابد و **overfitting** اتفاق نیفتد. وقتی لاندا را افزایش می‌دهیم در واقع این بایاس را بیشتر می‌کنیم اما وقتی پارامتر لاندا را بیش از اندازه بزرگ کنیم میزان فیت شدن آن به داده به حدی کاهش می‌یابد که موجب می‌شود مدل به درستی داده‌های آموزش را فرا نگیرد، در نتیجه عملکرد تست و **train** هر دو بدتر می‌شود و بهم نزدیک می‌شود.