



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

حل تمرین درس یادگیری ماشین

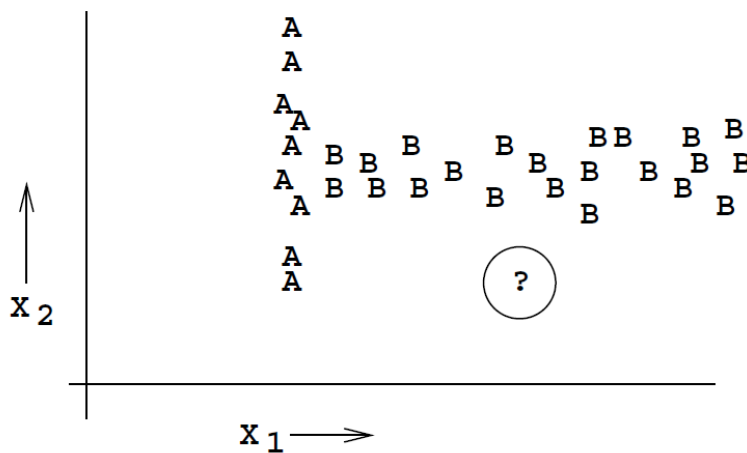
سری ۳

نام و نام خانوادگی دانشجو:

همایون حیدرزاده (۹۵۱۳۱۰۷۰)

نام استاد: دکتر ناظر فرد

آبان ۹۶



فرض میکنیم توزیع کلاس‌ها گوسی باشد بنابراین طبق شکل داریم:

$$p(A)=9/30=0.3, p(B)=0.7$$

? $\rightarrow z$

$$p(\text{class}|z)=\text{argmax } p(z|\text{class})p(\text{class})$$

$$p(A|z)=p(z|A)p(A)=p(z1|A)p(z2|A)p(A)$$

$$p(B|z)=p(z|B)p(B)=p(z1|B)p(z2|B)p(A)$$

چون در اینجا واریانس کلاس A برای ویژگی x_1 تقریباً صفر است بنابراین احتمال $p(A|z_1)$ بسیار کوچک می‌شود و این مقدار برای کلاس B نزدیک یک است چون نزدیک میانگین کلاس B است.

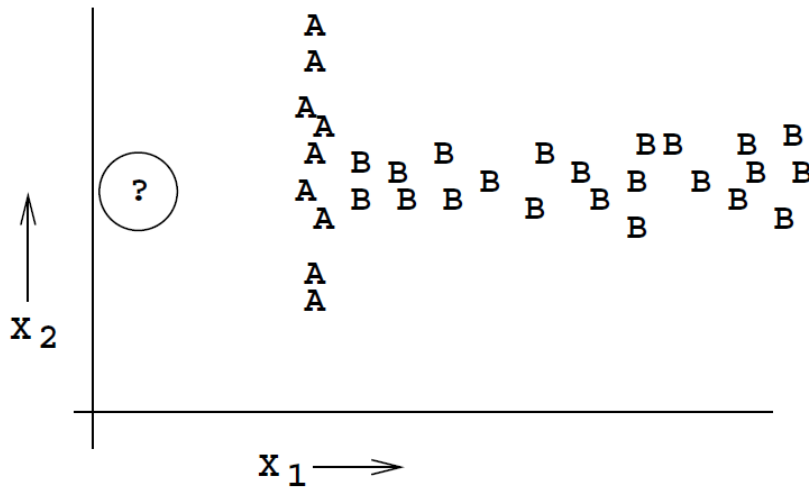
چون مقادیر احتمال برای کلاس B همگی بزرگ هستند بنابراین برچسب داده Z کلاس B می‌شود. یعنی داریم:

$$p(B)>p(A)$$

$$p(z1|B) \sim 1.0$$

$$p(z1|A) \sim 0.0$$

$$p(z2|B) \text{ کمی بیشتر از } p(z2|A)$$

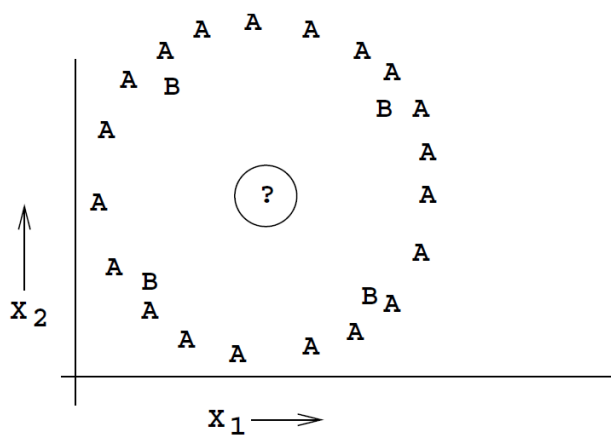


در این قسمت نیز مانند قسمت قبلی احتمال ها را مقایسه می کنیم:

$p(B) > p(A)$
 $p(z1|B) \gg 0$
 $p(z1|A) \sim 0.0$
 $p(z2|B) \sim 1.0$
 $p(z2|A) \sim 1.0$ چون نزدیک میانگین هستند

چون در جهت x_2 داده نزدیک میانگین دو کلاس است پس احتمال آن تقریباً یک می شود (برای هر دو کلاس). ولی در جهت x_1 واریانس کلاس A بسیار پایین و نزدیک صفر است بنابراین احتمال $p(A|z1)$ باز هم نزدیک صفر است و این احتمال برای کلاس B خیلی بزرگتر است (واریانس بالا در جهت x_1). همچنین احتمال پیشین کلاس B نیز بیشتر است، بنابراین برچسب باز هم B می شود.

(ج)



داریم:

$$p(A)=4/24=0.16, p(B)=0.84$$

? -> z

$$p(\text{class}|z)=\text{argmax } p(z|\text{class})p(\text{class})$$

$$p(A|z)=p(z|A)p(A)=p(z1|A)p(z2|A)p(A)$$

$$p(B|z)=p(z|B)p(B)=p(z1|B)p(z2|B)p(B)$$

اگر توزیع‌ها را گوسی در نظر بگیریم چون داده نزدیک میانگین هر دو کلاس است پس احتمال‌های زیر تقریباً ۱ می‌شوند:

$$p(B)>p(A)$$

$$p(z1|B) \sim 1.0$$

$$p(z1|A) \sim 1.0$$

$$p(z2|B) \sim 1.0$$

$$p(z2|A) \sim 1.0$$

ولی چون احتمال پیشین کلاس B خیلی بیشتر از A است پس باز هم برچسب داده جدید کلاس B خواهد بود.

(۲)

(الف)

اگر استقلالی وجود نداشته باشد، مدل ما به صورت زیر خواهد بود:

$$\begin{aligned} p(C_k, x_1, \dots, x_n) &= p(x_1, \dots, x_n, C_k) \\ &= p(x_1 | x_2, \dots, x_n, C_k) p(x_2, \dots, x_n, C_k) \\ &= p(x_1 | x_2, \dots, x_n, C_k) p(x_2 | x_3, \dots, x_n, C_k) p(x_3, \dots, x_n, C_k) \\ &= \dots \\ &= p(x_1 | x_2, \dots, x_n, C_k) p(x_2 | x_3, \dots, x_n, C_k) \dots p(x_{n-1} | x_n, C_k) p(x_n | C_k) p(C_k) \end{aligned}$$

چون ویژگی‌های باینری هستند پس برای هر احتمال ۲ به توان تعداد متغیرها حالت در نظر بگیریم، بنابراین تعداد پارامترها برای حالت باینری:

$$2+4+\dots+2^{n+1} = 2^{n+2}-1-1=2(2^{n+1}-1)$$

(ب)

اگر ویژگی‌ها با دانستن دسته، از یکدیگر مستقل باشند به مدل بیز ساده می‌رسیم:

$$\begin{aligned} p(C_k | x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots \\ &\propto p(C_k) \prod_{i=1}^n p(x_i | C_k). \end{aligned}$$

که در آن تعداد پارامترهای به صورت زیر است:

$$2+4+\dots+4=2+4n=2(2n+1)$$

(ج)

در این حالت نیز مانند قسمت ب عمل می‌کنیم با این تفاوت که پارامترها به گونه دیگری هستند مثلاً اگر همه ویژگی‌ها نسبت به یک کلاس دارای m پارامتر باشند داریم:

$$2+2m+\dots+2m=2+2mn=2(mn+1)$$

$$P(B|D=T) = P(B=T|D=T) + P(B=F|D=T)$$

$$= \frac{P(D=T|B=T)P(B=T)}{P(D=T)} + \frac{P(D=T|B=F)P(B=F)}{P(D=T)}$$

$$= \frac{.24 \times .42 + .1772 \times .58}{.749944} = .939313983980$$

$$P(D=T) = P(D=T|B=T, C=T)P(B=T, C=T) + P(D=T|B=T, C=F)P(B=T, C=F)$$

$$+ P(D=T|B=F, C=F)P(B=F, C=F) + P(D=T|B=F, C=T)P(B=F, C=T)$$

$$P(B, C) = P(B, C|A=T)P(A=T) + P(B, C|A=F)P(A=F)$$

$$B \perp C | A \Rightarrow P(B, C) = P(B|A=T)P(C|A=T)P(A=T) + P(B|A=F)P(C|A=F)P(A=F)$$

$$P(C) = P(C|A=T)P(A=T) + P(C|A=F)P(A=F)$$

$$P(B) = P(B|A=T)P(A=T) + P(B|A=F)P(A=F)$$

$$P(D|B) = P(D|B, C=T)P(C=T) + P(D|B, C=F)P(C=F)$$

احتمال هایی که زیر آن خط کشیده شده است می تواند کلی هستند و برای محاسبه نهایی استفاده می شوند \Rightarrow

$$P(B=F) = .43 \times .2 + .59 \times .8 \Rightarrow P(B=T) = .242$$

$$= .58$$

$$P(D=T|B=T) = .2 \times .26 + .1 \times .74$$

$$= .24$$

$$P(C=T) = .3 \times .2 + .25 \times .8$$

$$= .26$$

$$\Rightarrow P(C=F) = .74$$

$$P(D \neq T|B=F) = .47 \times .26 + .95 \times .74 = .8772$$

$$P(D=T) = .2 \times .0422$$

$$+ .1 \times .1772$$

$$+ .47 \times .0318$$

$$+ .95 \times .7494$$

$$= .749944$$

 \Leftarrow

$$P(B=T, C=T) = .3 \times .2 \times .26 + .25 \times .8 \times .1$$

$$= .0422$$

$$P(B=T, C=F) = .3 \times .74 \times .2 + .25 \times .58 \times .8$$

$$= .1772$$

$$P(B=F, C=T) = 1 - .26 = .74$$

$$= .0318$$

$$P(B=F, C=F) = .47 \times .26 + .95 \times .74 = .8772$$

منظور از مدل **Generative**، هر نوع فرایند یادگیری طبقه‌بندی با استفاده از تخمین احتمال مشترک $P(X,Y)$ یا هر نوع فرایند یادگیری طبقه‌بندی است که با استفاده از تخمین احتمال مقدم، که در آن Y یک طبقه و X توصیفی شئی است که باید طبقه‌بندی شود. با توجه به این مدل یا تخمین‌ها، می‌توان اشیای ترکیبی از توزیع مشترک ایجاد کرد. مدل **Generative** در مقابل مدل **Discriminative** است که در آن مدل یا تخمینی از $P(Y|X)$ بدون ارجاع به تخمینی صریح از $P(X)$ ، $P(Y,X)$ یا $P(Y|X)$ ایجاد می‌شود.

همچنین دسته‌بندی به عنوان رویکردهای افتراقی بر اساس تابع تصمیم نیز مرسوم است که به صورت مستقیم از ورودی X به خروجی Y نگاشت می‌کند (مثل ماشین‌های بردار پشتیبان، شبکه‌های عصبی و درخت‌های تصمیم)، که در آن ریسک تصمیم بدون تخمین $P(Y,X)$ ، $P(Y|X)$ یا $P(Y|X)$ کمینه می‌شود.

مثال استاندارد از مدل **Generative**، نایو بیز و از مدل **Discriminative** لاجستیک رگرسیون است.

مدل **Generative** هنگامی که نمونه‌ها کمیاب هستند خوب عمل می‌کند در حالی که مدل **Discriminative** عملکرد خطای مجانبی بهتری دارد.

مدل **Naive Bayes** یک مدل **Generative** است در صورتی که مدل **Logistic Regression** یک مدل **Discriminative** است.

به این معنی که مدل **Bayes** پارامترهای $P(Y)$ and $P(X|Y)$ را تخمین می‌زند، در صورتی که **Logit** پارامترهای $P(Y|X)$ را تخمین می‌زند.

بنابراین مدل **Logit** پیشفرضی بر روی توزیع داده ندارد و فقط سعی می‌کند **MCLE** احتمال $P(Y|X)$ را پیشینه کند. همچنین فرض استقلال شرطی در **Logit** به اندازه **Bayes** محکم و صریح نیست و **Logit** می‌تواند پارامترهایی متفاوت در این مواقع داشته باشد.

گر چه، اگر پارامترهای مدل **Logit** را با استفاده از مدل بیز ساده حساب شود، دو مدل یک دسته‌بند را ایجاد می‌کنند.

برای آموزش **Logit** به $O(n)$ داده نیاز دارد (n تعداد ویژگی) در حالی که **Bayes** با $O(\lg n)$ داده قابل آموزش است و در عمل نشان داده شده است زمانی که داده زیاد است **Logit** بهتر عمل می‌کند و زمانی که داده کم باشد **Bayes** دسته بند بهتری ایجاد می‌کند (**Generalization** بهتری دارد).

مدل Bayes ، بایاس بیشتر و واریانس کمتری نسبت به مدل Logit دارد، ولی اگر تعداد داده‌ها به بی‌نهایت میل کند هر دو دسته‌بند برابر هم خواهند شد

پیاده‌سازی

Logistic Regression (الف)

کد در فایل main_lr.ipnb موجود است. در تمامی قسمت‌ها از ضریب یادگیری ۰.۰۱ استفاده شد.

(ب)

لیست مقادیر خطای fold-۱۰ برای لامبدهای مختلف به صورت زیر به دست آمد:

Lambda	Error
[2.	0.25597484]
[1.	0.25974843]
[0.	0.25597484]
[-1.	0.25786164]
[-2.	0.25660377]
[-3.	0.26100629]
[-4.	0.26100629]
[-5.	0.26100629]

Picking best lambda --> 2

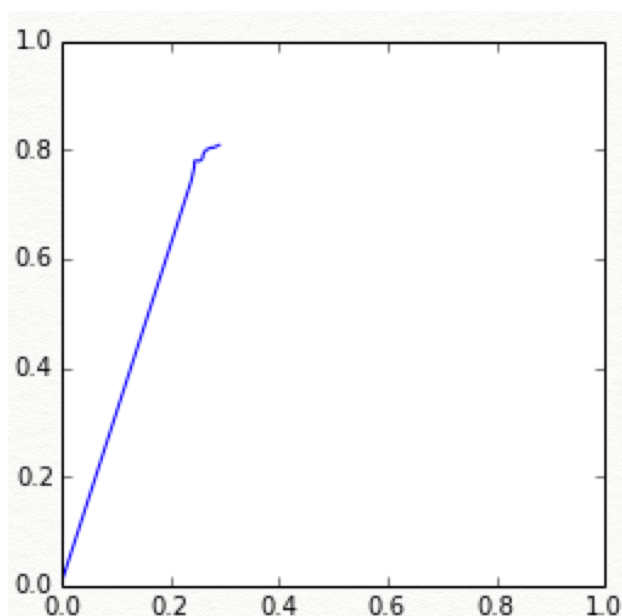
بنابراین لامبدا بهینه ۱۰ به توان ۲ (است. تعداد قدم‌ها برای هر مرحله ۲۰۰۰ قدم بود).

(ج)

در این قسمت confusion matrix و خطای داده آموزشی و آزمایشی برای لامبدا ۱۰ به توان ۲ و با ۱۰۰۰۰ قدم به صورت زیر محاسبه گردید:

```
Test confusion matrix:
[[ 146.   45.]
 [  50.  159.]]
Test error:
0.2375
Train confusion matrix:
[[ 637.  149.]
 [ 167.  647.]]
Train error:
0.1975
```


نمودار ROC نیز با جابجا کردن T به صورت زیر به دست آمد:



Naïve Bayes

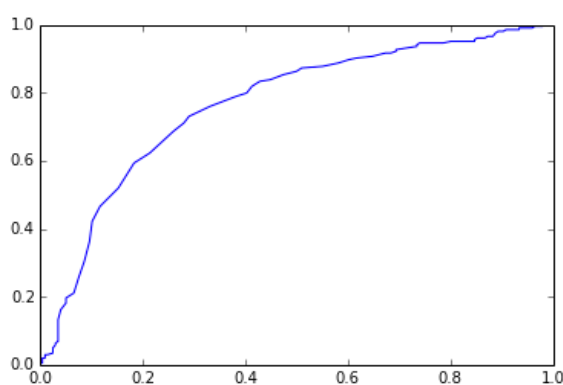
(الف)

کد در فایل main_naive.ipnb موجود است.

Test error	0.240625
Train error	0.295

(ب)

برای به دست آوردن نمودار ROC ابتدا تفاضل احتمال‌های پسین محاسبه شد و با قرار داده T از مینیمم تا ماکسیمم مقدار تفاضل، specificity و sensitivity-۱ محاسبه گردید:



با توجه به نمودار ROC و AUC متناظر با آن این مدل نسبتاً خوب است و خیلی بهتر از یک مدل تصادفی عمل کرده است.

Bayes Network

الف) کد این قسمت در فایل `main_bayes_net.ipnb` موجود است.

مدل‌های بیز ساده زیر استفاده شد:

Model\Accuracy	
Naive1	0.854221
Naive2	0.689259
Naive3	0.814276

هر مدل بیز ساده با حذف یکی از ویژگی‌ها ساخته شد.

ب)

ویژگی‌های موجود برای این مجموعه داده به صورت زیر مشخص شده است:

class values:

unacc, acc, good, vgood

attributes:

buying: vhigh, high, med, low.

maint: vhigh, high, med, low.

doors: 2, 3, 4, 5more.

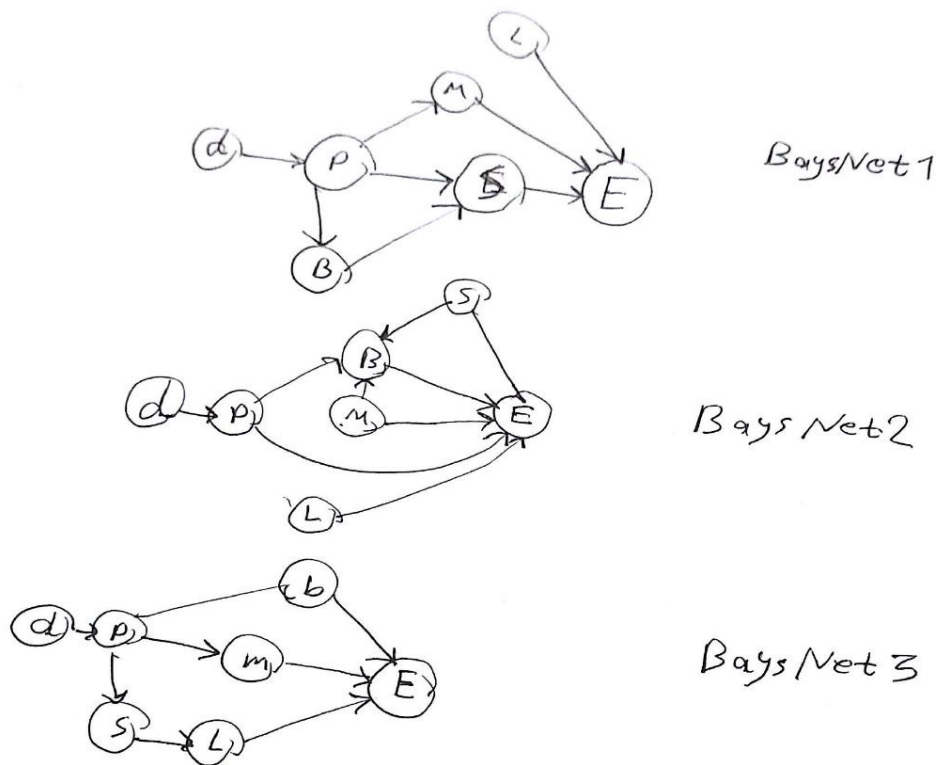
persons: 2, 4, more.

lug_boot: small, med, big.

safety: low, med, high.

با توجه به ویژگی‌های بالا گراف مدل‌ها در تصویر زیر آورده شده است:

(هر گره شبکه بیزین اول حرف ویژگی متناظر است)



توجیح مدل‌ها در زیر آورده شده است:

Model\Reason	بعضی از فرض‌ها
BayesNet1	فرض شده است که تعداد سرنشین ماشین به تعداد درهای ماشین وابستگی دارد. امنیت به تعداد سرنشین وابسته است. امنیت ماشین به قیمت ماشین وابستگی دارد.
BayesNet2	قیمت ماشین به امنیت ماشین وابستگی دارد. قیمت ماشین به تعداد سرنشین وابسته است. قیمت ماشین به هزینه نگهداری ماشین وابسته است.
BayesNet2	یک فرض این قسمت وابستگی اندازه جای پا به امنیت ماشین است، یعنی هر چه جای پا بیشتر باشد احتمالاً ماشین ایمن‌تر بوده است!

دقت مدل‌ها در جدول زیر آمده است:

Model\Accuracy	
BayesNet1	0.671881
BayesNet2	0.939824
BayesNet3	0.697933

مشاهده می‌شود که ساختار دوم بالاترین دقت را داشته است. بنابراین وابستگی بین ویژگی‌ها بسیار خوب در نظر گرفته شده است. . بنابراین برای داده‌های ناشناخته استخراج ویژگی بهتری صورت گرفته است.