



به نام خدا
درس یادگیری
ماشین تمرین اول

حدیثه مصباح

810102253



پاسخ I. سوالاتی راجع M-Class Classifier

در یک مسئله طبقه‌بندی چند کلاسه:

الف) نشان دهید که تصمیم‌گیری به کمک روش Bayes احتمال خطا را کمینه می‌کند.

ب) ثابت کنید اگر M کلاس داشته باشیم، حد بالای خطا به صورت $p_e \leq \frac{M-1}{M}$ خواهد بود.

ج) راهی برای رسم نمودار ROC در حالت چند کلاسه پیشنهاد کنید.

د) توضیح دهید که در چه مجموعه داده‌هایی naïve Bayes عملکرد بهینه خواهد داشت. علت را به تفصیل

شرح دهید.

سوال ۱ (الف)

$$p_e = \sum_{i=1}^M p(\alpha \notin R_i, w_i)$$

$$p(\alpha \notin R_i, w_i) = P(\alpha \notin R_i | w_i) P(w_i) = \left(\int_{R_i} P(\alpha | w_i) d\alpha \right) P(w_i)$$

$$p_e \rightarrow \sum_{i=1}^m \left(\int_{R_i} P(\alpha | w_i) d\alpha \right) P(w_i) = \sum_{i=1}^m \left(\int_{R_i} P(\alpha | w_i) P(w_i) d\alpha \right)$$

برای مایموم کردن p_e داریم (چون m مقادیر داریم)

$$P(\alpha | w_i) P(w_i) > P(\alpha | w_j) P(w_j) \quad \forall j \neq i$$

با استفاده از قانون Bayes

$$P(w_i | \alpha) > P(w_j | \alpha) \quad \forall j \neq i$$

۱-۱

از آن جا که

$$\sum_{i=1}^m P(w_i | \alpha) = 1$$

حاصل می‌شود که

$$P(w_i | \alpha) \geq \frac{1}{m}$$

و استفاده از قانون Bayes داریم

$$P(w_i | \alpha) > \frac{1}{m} \rightarrow P(w_i | \alpha) = \max P(w_i | \alpha)$$

$$\rightarrow p_e = 1 - \max P(w_i | \alpha) \leq 1 - \frac{1}{m}$$



ج) برای استفاده از نمودار ROC در کلاسی که چند حالت باید اول فرضی را به صورت باینری تعریف کنیم بهترین روش مقایسه دو کلاسها با هم است. مثلاً ما یک کلاس را با هم کلاسها مقایسه می کنیم و بعد می بینیم سرانجام کلاس بعدی و به همین ترتیب ادامه می دهیم

د) حافظه می کنیم که داده ها کاملاً مستقل هستند (Complete independence) و در نتیجه naive bayes می تواند بهترین عملکرد را داشته باشد اما این به معنی وجود Conditional independence نیست

- ① New york is a crowded city
- ② New Cars! we offer the cheapest new Car!
- ③ The new pub 'pork slaughterhouse' opened today in york

با این به تمام این حالات مستقل از هم هستند ولی naive bayes هم با تمام عنوان New york این جملها را به دلیل وجود وزن می بیند در هر کلاس هم هست پس نمی شود فقط عملگر دهیم از من انتظار داشتیم در این موارد می توان از طبقه بندی logistic regression استفاده کرد که عملگر بهتری خواهد داشت

پاسخ 2. The corresponding area of the two classes in the Bayes class.

یک طبقه بند دو کلاسه با احتمال پیشین مساوی را در نظر بگیرید. فرض کنید داده های دو کلاس بر اساس

توزیع های زیر تولید می شوند:

$$p(x|y=1) = \frac{x}{\sigma^2} \exp\left(-\frac{x^2}{2\sigma^2}\right) \quad x \geq 0$$

$$p(x|y=2) = \theta x \exp(-\theta x) \quad x \leq 0$$

که $\sigma > 0$ و $\theta > 0$ پارامترهای مدل هستند. ناحیه مربوط به دو کلاس را در طبقه بند بیز به دست آورید.



Prior Probability $\rightarrow P(y=1) \cdot P(y=2)$ ①

② $\frac{P(x|y=1)}{P(x|y=2)} \geq \frac{P(y=2)}{P(y=1)}$ $\xrightarrow{1,2} \frac{\frac{x}{\sigma^2} \exp(-\frac{x^2}{2\sigma^2})}{\theta x e^{\theta x}} \geq 1$

$\frac{1}{\theta \sigma^2} e^{\frac{-x^2}{2\sigma^2} + \theta x} \geq 1$ $\xrightarrow{\text{از طرفین } \ln \text{ می گیریم}} \ln\left(\frac{1}{\theta \sigma^2}\right) + \ln\left(e^{\frac{-x^2}{2\sigma^2} + \theta x}\right) = \ln 1$

$-\ln(\theta \sigma^2) + \theta x - \frac{x^2}{2\sigma^2} = 0$ $\xrightarrow{\text{در معادله}} \Delta = \theta^2 - \frac{4}{2\sigma^2} \ln(\theta \sigma^2)$

$x_{1,2} = \frac{+\theta \pm \sqrt{\theta^2 - \frac{2}{\sigma^2} \ln(\theta \sigma^2)}}{\frac{1}{\sigma^2}}$

چون اطلاعاتی از θ و σ نداریم بهترین نقطه قطع مقدار + قابل قبول هستیم یا مقدار منفی.

پاسخ 3. RISK MATRIX

ماتریس ریسک زیر را در نظر بگیرید.

$$\begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

الف) نشان دهید با این شرایط مرز تصمیم شرط زیر را ارضا می کند.

$$\int_{R_2} p(x|\omega_1) dx = \int_{R_1} p(x|\omega_2) dx$$

ب) آیا این پاسخ همواره یکتاست؟ در غیر این صورت یک مثال نقض بزنید.



سوال چهارم ()

$$p(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$
 معده تقاطع / تقاطع
 $N(\mu, \sigma^2)$

★ ما نسیم برین MAP / Maximum a Posteriori / Maximum Likelihood

$$\mu_{MAP} = \arg \max_{\mu} \frac{P(x|\mu)P(\mu)}{P(x)}$$

→ $\arg \max_{\mu} P(x|\mu)P(\mu)$

likelihood از داده های قبلی x_i ها، μ و σ^2 است.

$$L(\mu) = \left(\prod_{i=1}^N P(x_i|\mu, \sigma^2) \right) P(\mu)$$

log likelihood

$$L(\mu) = \ln(L(\mu)) = \ln(P(\mu)) + \sum_{k=1}^N \ln(p(x_k|\mu, \sigma^2))$$

$$\frac{\mu \exp\left(-\frac{\mu^2}{2\sigma_{\mu}^2}\right)}{\sigma_{\mu}^2}$$

محاسبه $\ln(P(\mu))$: $\ln \mu - \ln(\sigma_{\mu}^2) - \frac{\mu^2}{2\sigma_{\mu}^2}$

محاسبه $\ln(p(x_k|\mu, \sigma^2)) = \sum_{k=1}^N \ln\left(\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x_k-\mu}{\sigma}\right)^2}\right) = N \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{1}{2\sigma^2} \sum_{k=1}^N (x_k - \mu)^2$

ما $\frac{dL(\mu)}{d\mu}$ را میگیریم، Max place در μ قرار می دهیم تا $L(\mu)$ بیشترین مقدار داشته باشد.

$$\frac{1}{\mu} - \frac{\mu}{\sigma_{\mu}^2} + \frac{1}{\sigma^2} \sum_{k=1}^N (x_k - \mu) = 0$$

$$-\frac{N\mu^2}{\sigma^2} - \frac{\mu^2}{2\sigma^2} + \frac{\mu}{\sigma^2} \sum_{k=1}^N x_k + 1 = 0$$

$\sum_{k=1}^N x_k = \sum_{k=1}^N \mu \rightarrow N\mu$

$$\mu^2 \left(\frac{1}{\sigma_{\mu}^2} + \frac{N}{\sigma^2} \right) - \left(\frac{1}{\sigma^2} \sum_{k=1}^N x_k \right) \mu + 1 = 0$$

$A = b^2 - 4ac \rightarrow \Delta = 2^2 + 4R$

$$x_1, x_2 = \frac{-2 \pm \sqrt{2^2 + 4R}}{2R}$$

ما Δ را میگیریم

$$\frac{2(1 + \sqrt{1 + 4R})}{2R} \quad \checkmark$$



پاسخ 5. Log likelihood Vs. Maximum A Posteriori

فرض کنید تابع چگالی احتمال متغیر تصادفی Y به شکل زیر است.

$$f_Y(y|\theta) = \begin{cases} \frac{1}{\theta} r y^{r-1} e^{-\frac{y^r}{\theta}}, & \theta > 0, y > 0 \\ 0, & \text{elsewhere} \end{cases}$$

که r یک ثابت مثبت است.

الف) تابع log likelihood را به دست بیاورید.

ب) توضیح دهید در چه حالتی و چرا تخمین گر MAP به ML میل می کند.

سوال پنج) $L(y_1, y_2, \dots, y_n | \theta) = \prod_{i=1}^n f(y_i | \theta) = \left(\frac{1}{\theta^n} \times r^n \times e^{-\sum_{i=1}^n y_i^r / \theta} \times \prod_{i=1}^n y_i^{r-1} \right)$

الف) $\mathcal{L}(\theta) \rightarrow -n \ln \theta - \frac{1}{\theta} \sum_{i=1}^n y_i^r + \ln \left(r^n \times \left(\prod_{i=1}^n y_i \right)^{r-1} \right)$
 $n \ln r + (r-1) \sum_{i=1}^n \ln y_i$

ب) ما نسیم $\hat{\theta}$ به این معنی است که $\hat{\theta}$ همان θ است که $\mathcal{L}(\theta)$ را بیشینه می کند.
 $\mathcal{L}'(\theta) = -n\theta^{-1} + \theta^{-2} \sum_{i=1}^n y_i^r = 0 \Rightarrow \hat{\theta} = \frac{1}{n} \sum_{i=1}^n y_i^r$

در حالتی که θ را $\frac{1}{n}$ فرض کنیم، توزیع یونیفر دارد.
 و در حالتی که θ را $\frac{1}{n}$ فرض کنیم، MAP و ML یکی می شوند. در این حالت هیچ دیتای Prior وجود ندارد.
 یا حتی کم وضعیف یا یونیفرم باشد چون این است که شود تابع خاصی ندارد. این Prior روی تخمین و حالتی که در حالیکه کم کردن ندارد و حتی ضربه می کشد.



پاسخ 6. Naïve bayes classification Vs. optimal classification

هدف از این سوال آشنایی و پیاده سازی طبقه‌بند naïve bayes است.

آ) در ابتدا در مورد طبقه‌بند naïve bayes توضیح دهید و تفاوت ساختاری آن را با یک طبقه‌بند بیزی بیان کنید. توضیح دهید که چرا به جای طبقه‌بند بیز از این طبقه‌بند استفاده می‌کنیم، هزینه‌ای که می‌دهیم چیست و در چه زمان‌هایی استفاده از این طبقه‌بند کاری منطقی است. مجموعه داده Breast Cancer Wisconsin را از لینک زیر دانلود کنید.

<https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data>

توضیحات مربوط به این مجموعه داده به طور کامل در صفحه فوق وجود دارد؛ لطفاً قبل از شروع به انجام تمرین توضیحات را مطالعه نمایید.

در ابتدا در صورت نیاز روی داده‌ها پیش‌پردازش انجام دهید. (هر پیش‌پردازشی که روی داده‌ها انجام می‌دهید را باید با ذکر دلیل توضیح دهید).

ب) این مجموعه داده شامل دو کلاس است. یک طبقه‌بند naïve bayes را از پایه و بدون استفاده از کتابخانه پیاده‌سازی کنید. و طبقه‌بندی که طراحی کردید استفاده کنید. دقت، precision، Recall و ماتریس آشفتگی^۱ را بررسی و تحلیل نمایید.

پ) مورد ب را به کمک کتابخانه SKLEARN انجام دهید. نتایج دو بخش را مقایسه کنید.

قسمت الف)

Naïve bayes: روشی برای دسته‌بندی است که از قضیه احتمالی بیز استفاده می‌کند که فرض می‌کند feature ها از هم مستقل هستند و هیچ ترتیب یا تعاملی میان آن‌ها وجود ندارد. این تبسیت، محاسبات را ساده‌تر می‌کند. و به طور کلی برچسبی را به داده اختصاص می‌دهد که احتمالش بیشتر باشد.

$$y = \arg \max_y P(y) \prod_{i=1}^n P(x_i|y)$$



Optimal bayes: مثل Naïve bayes است فقط فرض مستقل بودن feature ها از هم را ندارد، و کواریانس بین آن‌ها را هم حساب می‌کند. پس نمی‌توان آن را مثل فرمول بالا نوشت. در این طبقه‌بند معمولاً از توزیع‌های احتمالی پیچیده‌تری برای مدل‌سازی استفاده می‌شود. این بدین معناست که طبقه‌بند بیزی می‌تواند تعامل‌ها و ترتیب‌های پیچیده‌تری میان ویژگی‌ها را مدل کند.

دلایل استفاده از Naïve bayes

سرعت و کارایی: طبقه‌بند Naïve bayes به دلیل ساختار ساده‌تر و محاسبات سریع‌تر معمولاً در مسائلی که دارای تعداد زیادی ویژگی هستند، مانند پردازش متن و معنایی، بسیار مؤثر و سریع است.

مقاومت در برابر ابعاد بالا: در مسائل با فضای ویژگی‌های بزرگ، نیاز به تعداد زیادی داده آموزشی ندارد و با دیتای کم هم نتیجه قابل قبولی می‌دهد.

مثال‌هایی برای استفاده از Naïve bayes

طبقه‌بندی متن: معمولاً برای شناسایی ایمیل‌های هرزنامه، تجزیه و تحلیل احساسات و طبقه‌بندی اسناد استفاده می‌شود، جایی که ویژگی‌ها (کلمات) می‌توانند به عنوان مستقل مشروط در نظر گرفته شوند.

فضاهای ویژگی بزرگ: زمانی که یک فضای ویژگی با ابعاد بالا دارید و الگوریتمی می‌خواهید که بتواند آن را به طور مؤثر مدیریت کند.

نمونه سازی سریع: برای نمونه سازی سریع و به عنوان یک مدل پایه برای مقایسه با الگوریتم‌های پیچیده تر مفید است.

طبقه‌بند Naïve bayes فرض‌های ساده‌تری نسبت به واقعیت دارد و در مواردی که تعاملات پیچیده‌تری بین ویژگی‌ها وجود دارد، عملکرد ضعیف‌تری داشته دارد با این حال، محدودیت‌هایی وجود دارد و کم پیش میاد که داده‌های ما به طور کلی و تماماً از هم مستقل باشند. در چنین مواردی، مدل‌های پیچیده‌تر مانند طبقه‌بند بیزی، درخت‌های تصمیم‌گیری یا شبکه‌های عصبی ممکن است مناسب‌تر باشند و با استفاده از Naïve bayes احتمال خطا بسیار بالا می‌رود و احتمال دسته بندی اشتباه وجود دارد.

(ب)

Recall score: 0.9166666666666666

یعنی در 93 درصد بیمار بودن به درستی تشخیص می‌دهد.

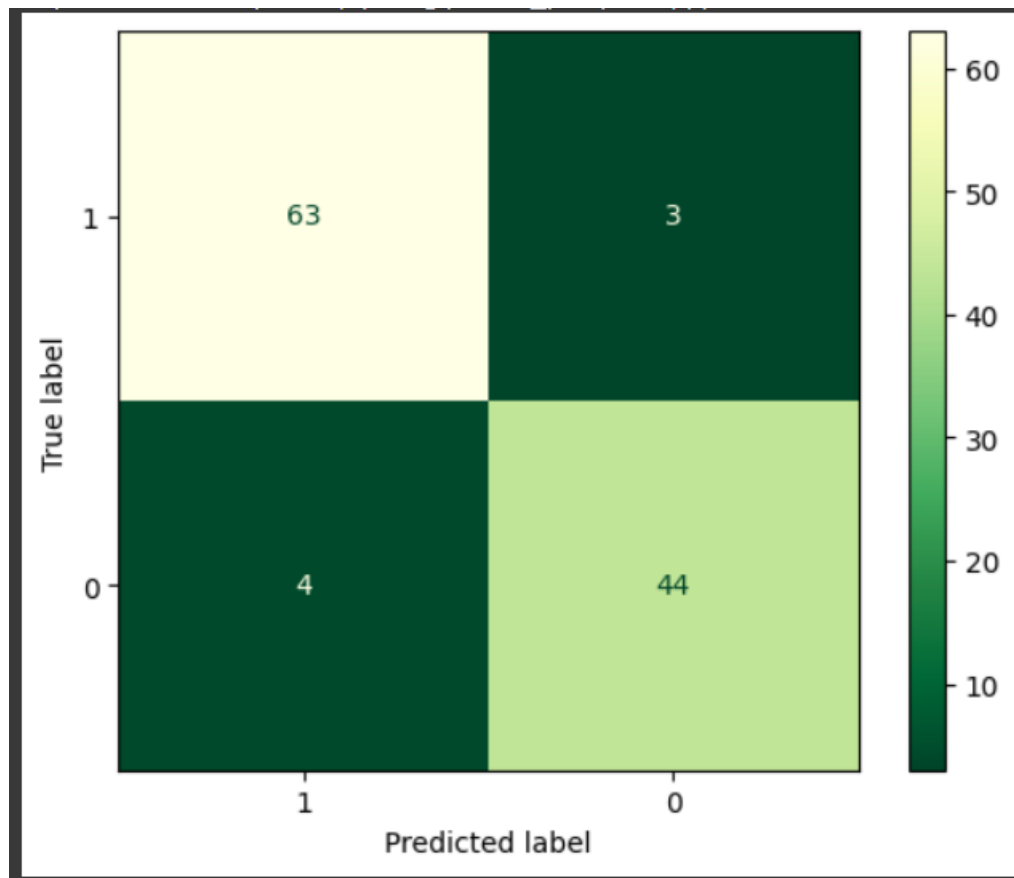
Precision score: 0.9361702127659575

یعنی در 93 درصد مواقعی که طبقه‌بند ما تشخیص می‌دهد کیس بیمار است، درست گفته است.



Naive Bayes classification accuracy 0.9385964912280702

یعنی در 92 درصد مواقع بیمار بودن یا نبودن به درستی تشخیص داده شده اند.



(پ)

(I برابر malignant و 0 برابر benign).

Recall: 0.82758

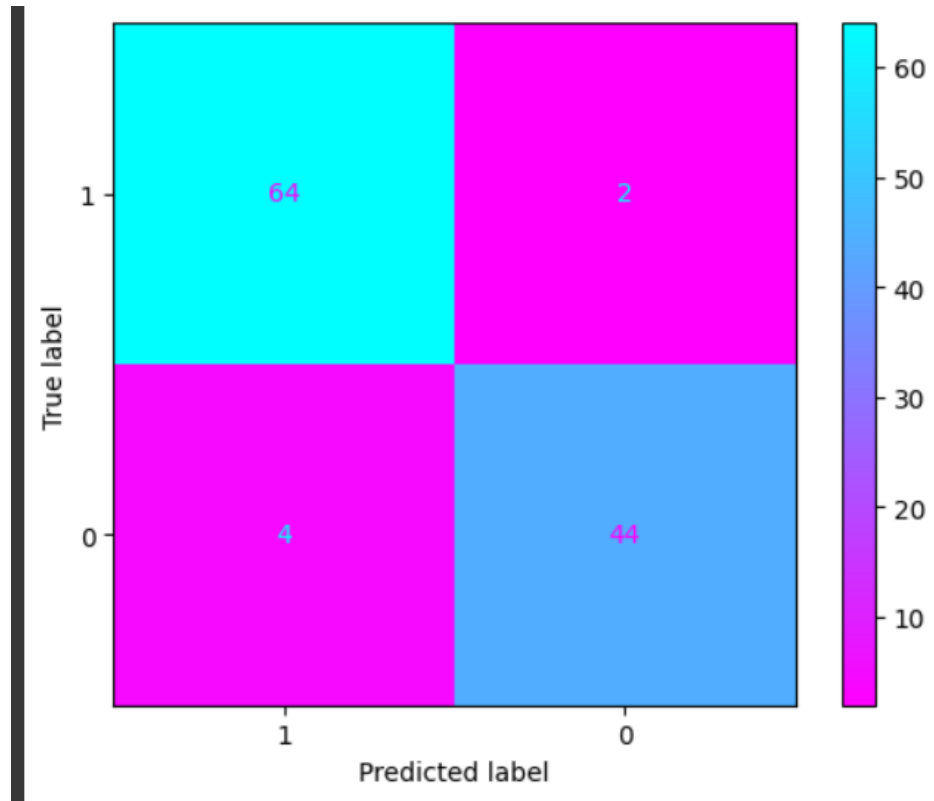
یعنی در 82 درصد بیمار بودن به درستی تشخیص میدهد.

Precision: 0.9795

یعنی در 97 درصد مواقعی که طبقه‌بند ما تشخیص میدهد کیس بیمار است، درست گفته است.

Accuracy: 0.92307

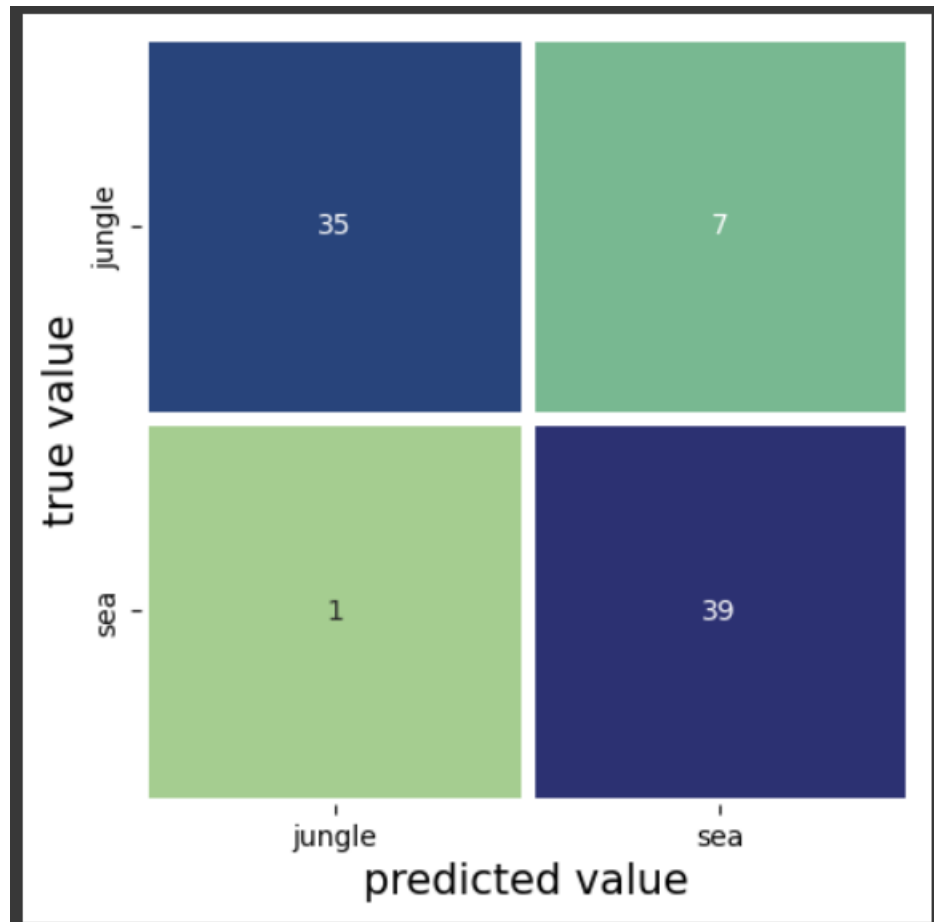
یعنی در 92 درصد مواقع بیمار بودن یا نبودن به درستی تشخیص داده شده اند.



همان‌طور که می‌انتظار داشتیم، نتایج به شدت به یکدیگر نزدیک هستند. (چه با استفاده از کتابخانه چه نوشتن از اول کد این طبقه بند (این نکته نشان می‌دهد که پیاده‌سازی الگوریتم Naïve Bayes بدون استفاده از کتابخانه به درستی کار می‌کند. به لحاظ مهندسی نتیجه‌های ما در این طبقه بند بسیار خوب است؛ اما به لحاظ انسانی، این الگوریتم باید در تشخیص بیماری دقیق‌تر عمل کند. به عبارت دیگر، باید بیشتر به سمت تشخیص بیماری متمایل باشد تا این که بگوید سالم، و بعداً این تشخیص می‌تواند در آزمون‌های بعدی مثل نمونه برداری یا اسکن‌های پزشکی تأیید یا رد شود. (2 نمونه با این که باید سرطانی تشخیص داده می‌شدن و لیبل یک می‌گرفتن لیبل غیر سرطانی گرفته اند) بنابراین، دقت (Precision) در این مسئله باید ۱۰۰ درصد باشد. (I برابر malignant و 0 برابر benign)



پاسخ 7. Implementation of binary classifier



Recall: 0.8333333333333334

یعنی در 82 درصد بیمار بودن به درستی تشخیص میدهد.

Precision: 0.9722222222222222

یعنی در 97 درصد مواقعی که طبقه‌بند ما تشخیص میدهد کیس بیمار است، درست گفته است.

Accuracy: 0.9024390243902439

یعنی در 92 درصد مواقع بیمار بودن یا نبودن به درستی تشخیص داده شده اند.

پی نوشت : کد هر دو سوال اخر ضمیمه شده است.