

## گزارش CA3 آمار

نام : مصطفی کرمانی نیا

شماره دانشجویی : 810101575

### Q1

#### (Part1)

در این بخش هیستوگرام تعداد ورود metro یا BTR را کشیدیم (یعنی محور X همان تعداد ورود بوده و محور Y هم فراوانی آن عدد بعنوان تعداد ورود است) که در فایل کد ها و نمودار ها مشخص اند

#### (Part2)

توزیعی که X و Y دارند پواسون است، زیرا طبق تعریف پواسون میدانیم که پیشامدهای نادری که در یک بازه زمانی یا مکانی خاص رخ میدهند دارای توزیع پواسون است تعداد رخ دادنشان در نتیجه اگر بازه ی زمان را بسیار کوچک کنیم مشخصا تعداد گذرهای BRT و مترو خیلی نادر میشود در نتیجه این مسئله ماهیتش توزیع پواسون دارد. حالا میدانیم که در توزیع پواسون میانگین و واریانس توزیع پواسون با پارامتر  $\lambda$ ، همان  $\lambda$  است. در نتیجه جهت یافتن پارامتر این توزیع کفایت یا با مشاهده ی نمودار بطور چشمی میانگین را حساب کرده و یا با توابع موجود در پایتون همانطور که در کد های من می بینیم، پارامتر این توزیع را حساب کنیم. (نکته اینجاست که اتفاقا در کدی که نوشتم واریانس و میانگین برابر شدند حدودا که باز هم تاییدی بر پواسون بودن این توزیع ها است) و نهایتا به این نتایج رسیدیم:

metro passages is a poisson distribution with  $\lambda \approx 3.57$

BRT passages is a poisson distribution with  $\lambda \approx 2.07$

#### (Part3)

برای این بخش کار خاصی لازم نبود بکنیم چون تابع hist یک پارامتر density دارد و کافی بود با true کردن آن، دستور دهیم که نمودارهای قبلی، scale شده و به ما احتمال وقوع هر مقداری برای X یا Y را بدهند.

#### (Part4)

کد این بخش در کد ها قابل مشاهده است و با کشیدن نمودار های گفته شده به روی نمودار های قبلی و فیت شدن بسیار خوب هر دو جفت نمودار روی همدیگر ، فهمیدیم که تقریب درستی زده بودیم (در این بخش و بخش قبلی، علاوه بر مترو، BRT را هم نمودارهایش را کشیدم و چک کردم)

## (Part5)

بطور تئوری و با توجه به تایید های درس، میدانیم مجموع دو متغیر پواسون مستقل با  $\lambda_1$  و  $\lambda_2$  مساوی با یک متغیر دیگر با توزیع پواسون با پارامتر  $\lambda_1 + \lambda_2$  است، حالا وقتی متغیر های  $X$  و  $Y$  را جمع زده و نمودار آنها را میکشیم، متوجه می شویم که بخوبی با نمودار متغیر پواسون با پارامتر  $\lambda_X + \lambda_Y = 3.57 + 2.07$  فیت می شود که تاییدی بر این است که فرمول هایی که بطور تئوری یافته ایم، در دنیای واقعی هم با داده های واقعی درست کار می کنند.

## (Part6)

محاسبات دستی بصورت زیر است:

بخش 1 تا 4 فرمایم:

$$X = \text{تعداد زلزله های مترو} \rightarrow X \sim \text{Poi}(\lambda_1 = 3.57) \quad Y \perp X$$

$$Y = \text{تعداد زلزله های BRT} \rightarrow Y \sim \text{Poi}(\lambda_2 = 2.07)$$

بخش 5 فرمایم:

$$Z = X + Y = \text{تعداد زلزله ها} \rightarrow Z \sim \text{Poi}(\lambda_1 + \lambda_2 = 3.57 + 2.07 = 5.64)$$

حالا پرسش این است که متغیرهای زیر، از چه توزیعی پیروی می کنند:

$$W \equiv X | X+Y$$

کانتی Pmf آنرا بیاییم با توزیع آن آنگاه است:

$$P_{X|X+Y}(X=k | X+Y=n) = \frac{P_{(X=k, X+Y=n)}}{P_{(X+Y=n)}} = \frac{P_{(X=k, Y=n-k)}}{P_{(Z=n)}} \stackrel{X \perp Y}{=} \frac{P_{X(k)} P_{Y(n-k)}}{P_{Z(n)}}$$

$$\stackrel{X, Y, Z \sim \text{Poi}}{=} \frac{\frac{e^{-\lambda_1} \lambda_1^k}{k!} \times \frac{e^{-\lambda_2} \lambda_2^{n-k}}{(n-k)!}}{\frac{e^{-(\lambda_1+\lambda_2)} (\lambda_1+\lambda_2)^n}{n!}} = \frac{n!}{k!(n-k)!} \times \frac{\lambda_1^k \lambda_2^{n-k}}{(\lambda_1+\lambda_2)^n} = \binom{n}{k} \left( \frac{\lambda_1}{\lambda_1+\lambda_2} \right)^k \left( \frac{\lambda_2}{\lambda_1+\lambda_2} \right)^{n-k}$$

این عبارت دقیقاً به فرمت Pmf یک متغیر دوجمله ای است از این داریم:

$$\text{if } S \sim \text{Bin}(n, p) \rightarrow P_{S(k)} = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\Rightarrow P_{(X=k | X+Y=n)} = P_{W(X=k)} = \binom{n}{k} \left( \frac{\lambda_1}{\lambda_1+\lambda_2} \right)^k \left( 1 - \frac{\lambda_1}{\lambda_1+\lambda_2} \right)^{n-k} \Rightarrow$$

$$W = X | X+Y \sim \text{Bin}\left(n, \frac{\lambda_1}{\lambda_1+\lambda_2}\right)$$

## (Part7)

با توجه به نتیجه ای که در محاسبات دستی گرفتیم، اکنون باید  $w \sim \text{Bin}(n, \frac{\lambda_1}{\lambda_1 + \lambda_2})$  را رسم کنیم که رسم کردم. که همانطور که در دوجمله ای انتظار داریم دارای یک پیک در نزدیکی  $np=5.06$  بود.

## (Part8)

همانطور که انتظار میرفت، نموداری که با محاسبات تئوری و بصورت دوجمله ای کشیده بودیم، با نموداری که حاصل از محاسبه ی دستی و عملی روی داده ها بود، به خوبی فیت شدند (محاسبه ی دستی روی داده ها هم به این صورت بود که ابتدا سطرهایی از دیتافریم را که مجموع مترو و BRT هایشات 8 میشد را جدا کرده و سپس تعداد مترو های آنها را خوانده و بین آنها هیستوگرام را رسم کردم که در کد نیز مشخص است)

نکته: در تمام بخش های سوال، در جاهایی که باید نمودار محاسبات عملی و تئوری روی هم فیت میشدند، در نواحی ای که مقدار نمودار ها به پیک نزدیک تر میشدند مقدار خطا بیشتر بود که طبیعی هم هست زیرا در کل درصد خطا ثابت است و وقتی مقدار احتمال و طول میله ی هیستوگرام زیادتر میشود طبیعی است که اندازه ی خطا هم بیشتر به چشم بیاید مگر نه درصد کمی از خطا داشتیم در کل.

## Q2

## (Part1)

کد جوابی که داده ام موجود است اما روند کلی حل اینگونه بود: ابتدا یک دیکشنری با  $n$  عضو میسازم که کلید هایش اعداد 0 تا  $n-1$  هستند که نشانگر نوع هر کوپن هستند، و  $value$  ها نیز بولین هستند که اگر False باشند یعنی آن کوپن دیده نشده و اگر True باشند دیده شده است، حالا هر بار یک نوع کوپن رندوم را برداشته و مقدار آنرا  $true$  میکنیم و این کار را تا جایی تکرار میکنیم که تمامی کوپن ها در دیکشنری مقدارشان  $true$  شود و سپس تعداد تکرار هایمان برای اینکار را گزارش میکنیم. و در تابعی دیگر این روند را  $k$  بار تکرار کرده و میانگین میگیریم.

## (Part2)

با چند بار خروجی گرفتن حدودا متوجه شدم به عدد 29.2 جواب ها میل میکنند (البته با بزرگتر کردن  $k$  همانطور که در کد های کامنت شده موجود است متوجه شدم 29.3 تقریب بهتری است اما با  $K$  هایی که خود سوال گفته بود، حدود 29 تا 29.2 تقریبی بود که بنظر آمد)

## (Hand calculate)

در اینجا بطور تئوری و بدون استفاده از مولد گشتاور یکبار مسئله را حل کردم و توضیحات دقیقتری راجع به فرضیاتمان از متغیر های تصادفی دادم:

$X = \text{شمار مشاهده لازم}$ ، برای بین تمام کوبین ها

$X_i = \text{شمار مشاهده لازم برای بین کوبین } i\text{ام}$  این از بین  $(i-1)$  کوبین مختلف اول.

$$X = \sum_{i=1}^n X_i \quad \Leftarrow$$

«کل  $n$  کوبین داریم، پس احتمال بین کوبین خاص،  $\frac{1}{n}$  است اما احتمال بین کوبین جدید، بعد از بین  $i-1$  کوبین مختلف، برابر است با  $\frac{n-(i-1)}{n}$  که تعداد انواع دیده نشده است»

برای  $X_i$  رسماً باید انتظار از بین کوبین با احتمال  $P_i = \frac{n-(i-1)}{n}$

توزیع  $P_i$  انجام دهیم تا به کوبین جدید برسیم. در نتیجه  $X_i$  توزیع هندسی با پارامتر  $P_i$  دارد

$$X_i \sim \text{Geo}(P_i) = \text{Geo}\left(\frac{n-(i-1)}{n}\right) \Rightarrow E[X_i] = \frac{1}{P_i} = \frac{n}{n-i+1}$$

$$E[X] = E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i] = \sum_{i=1}^n \frac{1}{P_i} = \sum_{i=1}^n \frac{n}{n-i+1} = \frac{n}{n} + \frac{n}{n-1} + \dots + \frac{n}{1}$$

حالا باید به خواص امید ریاضی:

$$\rightarrow E[X] = n\left(\frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n}\right) \approx n \ln(n)$$

$$\text{if } n=10 \rightarrow E_{(X)} = 10\left(\frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{10}\right) = 29.28968254 \dots$$

پس دیدیم که  $X_i$  ها توزیع هندسی دارند اما پارامتر این توزیع های هندسی با هم متفاوت است سوالی که در متن پرسیده شده بود هم کمی ابهام داشت، اگر منظور این است که موقع جستجو برای کوپن جدید، احتمال دیدن هر کدام از کوپن های دیده نشده با هم برابر است، بله همینطور است، اما اگر منظور این است که مثلاً وقتی هیچ کوپنی ندیده ایم احتمال دیدن کوپن جدید با وقتی که مثلاً 5 تا کوپن دیده ایم برابر است، باید گفت که خیر، اینطور نیست و همانطور که در محاسبات دستی من مشخص است، احتمال دیدن کوپن جدید از عدد  $1/1$  تا عدد  $1/n$  متغیر است.

### (Part3)

$$\begin{aligned}\phi_{X(s)} &= E(e^{sX}) = \sum_i e^{su_i} P_{X(u_i)} \\ \rightarrow \phi_{X_i(s)} &= E(e^{sX_i}) = \sum_K e^{sK} P_{X_i(K)} \\ X_i &\sim \text{Geo}\left(\frac{n-(i-1)}{n}\right) \Rightarrow P_{X_i(K)} = P_{X_i=K} = \left(\frac{n-(i-1)}{n}\right) \left(1 - \frac{n-(i-1)}{n}\right)^{K-1} \\ \Rightarrow P_{X_i(K)} &= \left(\frac{n-i+1}{n}\right) \left(\frac{n-n+1-1}{n}\right)^{K-1} = \left(\frac{n-i+1}{n}\right) \left(\frac{i-1}{n}\right)^{K-1} \\ \Rightarrow \phi_{X_i(s)} &= \sum_K e^{sK} \left(\frac{n-i+1}{n}\right) \left(\frac{i-1}{n}\right)^{K-1} = \frac{n-i+1}{n} \sum_{K=1}^{\infty} e^{sK} \left(\frac{i-1}{n}\right)^{K-1} \\ &= \frac{n-i+1}{i-1} \sum_{K=1}^{\infty} \left(\frac{(i-1)e^s}{n}\right)^K \xrightarrow{n \rightarrow \infty} \frac{n-i+1}{i-1} \sum_{K=1}^{\infty} \left(\frac{(i-1)e^s}{10}\right)^K\end{aligned}$$

میتوانم سیگما ها را دستی حساب کرده و فقط برای گرفتن نتیجه ی نهایی از sympy استفاده کنم اما ترجیح دادم تمام کارها را به خودش بسپرم پس در نهایت این عبارت را بعنوان مولد گشتاور هر کدام از متغیر ها در لیست گشتاور ها قرار میگیرد بطور پارامتری وارد کردم و لیستی از گشتاور ها ساختم:

$$\phi_{X_i(s)} = \sum_K e^{sK} \left(\frac{n-i+1}{n}\right) \left(\frac{i-1}{n}\right)^{K-1}$$

در اینجا یک مشکل وجود دارد، به ازای  $i = 1$  مولد گشتاور درستی دریافت نمیکنیم چون اولین بار که یک نوع جدید کارت میخواهیم به احتمال 1 قرار است کارتی که برمیداریم جدید باشد اما این متغیر تصادفی را به چشم یک متغیر تصادفی هندسی بهش نگاه کرده ایم در نتیجه در محاسبه ی احتمال  $K = 1$  به عبارت صفر به توان صفر میخوریم، در نتیجه ترجیح دادم محاسبات مربوط به مولد گشتاور متغیر تصادفی  $X_1$  را بطور دستی محاسبه کرده و نهایتاً آنرا وارد لیست مولد گشتاورها بکنم:

$$P_{X_i(k)} = \begin{cases} 1 & k=1 \\ 0 & k \neq 1 \end{cases}$$

$$\phi_{X_i(s)} = \sum_{k=1}^{\infty} P_{X_i(k)} e^{sk} = 1 \times e^s = e^s$$

$$\Rightarrow \underline{\alpha} = [\phi_{X_1(s)}, \phi_{X_2(s)}, \dots, \phi_{X_n(s)}]$$

$$= \left[ e^s, \sum_{k=1}^{\infty} \frac{n-1}{n} \left( \frac{1}{n} \right)^{k-1} e^{sk}, \sum_{k=1}^{\infty} \frac{n-2}{n} \left( \frac{2}{n} \right)^{k-1} e^{sk}, \dots, \sum_{k=1}^{\infty} \frac{n-1}{n} \left( \frac{1}{n} \right)^{k-1} e^{sk} \right]$$

$$\stackrel{n=10}{=} \left[ e^s, \sum_{k=1}^{\infty} \frac{9}{10} \left( \frac{1}{10} \right)^{k-1} e^{sk}, \sum_{k=1}^{\infty} \frac{8}{10} \left( \frac{2}{10} \right)^{k-1} e^{sk}, \dots, \sum_{k=1}^{\infty} \frac{1}{10} \left( \frac{9}{10} \right)^{k-1} e^{sk} \right]$$

#### (Part4

همانطور که در درس دیده ایم، مولد گشتاور مجموع دو یا چند متغیر مستقل، برابر حاصلضرب مولد های آنها است، اینجا نیز  $X_i$ ها از هم مستقل اند زیرا مثلاً اینکه برای یافتن کوپن سوم، 2 تلاش نا موفق کرده باشیم یا 15 تلاش ناموفق کرده باشیم، فرقی به حال کوپن چهارم نمیکند و باز هم برای یافتن کوپن چهارم همان فرمول های قبلی برقرار است در نتیجه دانستن هر کدام از این  $X_i$ ها هیچ اثری روی بقیه ندارد در نتیجه همگی آنها از همدیگر مستقل اند.

$X_1, X_2, \dots, X_n$  are independent because  
for every sublist of  $X_i$  we have this:

$$P(X_{i_1} = u_1, X_{i_2} = u_2, \dots, X_{i_r} = u_r) = \prod_{k=1}^r P(X_{i_k} = u_k)$$

از دانشن هیچ دسای از  $X_i$  ها، هیچ دیکنی راجع ب  $X_i$  ندی نی «مد».

پس فقط کافی بود مولد ها را در هم ضرب کنیم که بدون زدن حلقه با استفاده از تابعی در نامپای هندل شد

**(Part5)**

نهایتا با توجه به فرمول زیر

$$m_n = E(X^n) = \Phi^{(n)}(0)$$

کافیست یکبار مشتق گرفته و مقدار  $s$  را صفر بگذاریم تا به امید ریاضی یا میانگین برسیم که پاسخ نهایی ای که گرفتیم این بود:

29.2896825396825

که بسیار نزدیک به پاسخ عملی (حدود 29.3) بوده و نشان می دهد نتایج تئوری و عملی با هم تطابق خوبی دارند.

### Q3

**(Part1)**

کد این بخش و نتیجه در کدها موجود است

**(Part2)**

برای اینکار ابتدا یک کپی از ستون label گرفتن سپس تمامی اعداد را تبدیل کردم سپس ستون label را به مقدار اولیه اش برگرداندم.

### (Part3)

ستون label را جدا کردم و reshape کرده و با تابع گفته شده نمایش دادم که در کد ها موجود است.

### (Part4)

برای نوشتن pny اولین راهی که به ذهنم رسید این است

$$pny = (p^{**n}) * (1-p)^{**(1-n)}$$

اما در مواقعی که  $p = 0, 1$  میشود از لحاظ ریاضی ممکن است به مشکل بخوریم، پس با if و else عادی نوشتیم آنرا که باز هم درست است میشد یک خطی و به این صورت بنویسم:

$$pny = p \text{ if } n \text{ else } 1-p$$

اما ترجیح دادم بطور گستره بنویسم تا انواع مقادیر ممکن برای n هم بیان شوند.

برای integral هم بطور کامل مسئله و کدی که برایش نوشته شده بود را بررسی کردم و نتیجه این شد:

$$f_{Y|N(n)} = \frac{P_{(N=n|Y=y)} f_{Y|Y}}{P_{(N=n)}} = \frac{P_{(N=n|Y=y)} f_{Y|Y}}{\int P_{(N=n|Y=y)} f_{Y|Y} dy}$$

$$N|Y \sim \text{Bernoulli}(Y) \Rightarrow P_{N|Y(n|y)} = P^n (1-p)^{1-n} = \begin{cases} P & n=1 \\ 1-P & n=0 \end{cases}$$

$$\int_0^1 P_{N|Y} f_Y dy = \frac{1-p}{N} \sum_{i=0}^N (P_{N|Y} f_Y)_{(y=\frac{i}{N})} = \frac{1}{N} \sum_{i=0}^N (P_{N|Y} f_Y)_{(y=\frac{i}{N})}$$

$$\frac{1}{t} \sum_{i=0}^t (P_{N|Y} f_Y)_{(y=\frac{i}{t})}$$

$$f_Y \sim \text{Beta}(1,1) \rightarrow P = (1, 1, 1, \dots, 1)$$

$$P_{N|Y} = P_{(N=n|Y=y)} = 404$$

$$\Rightarrow P_{N|Y} = \begin{cases} P & n=1 \\ 1-P & n=0 \end{cases} \rightarrow \begin{cases} (1, 1, 1, \dots, 1) & n=1 \\ (1, 0, 0, \dots, 0) & n=0 \end{cases}$$

$$\Rightarrow \frac{1}{N} \sum_{i=0}^N (P_{N|Y} f_Y)_{(y=\frac{i}{N})} = \frac{1}{N} P_{N|Y} \cdot f_Y$$

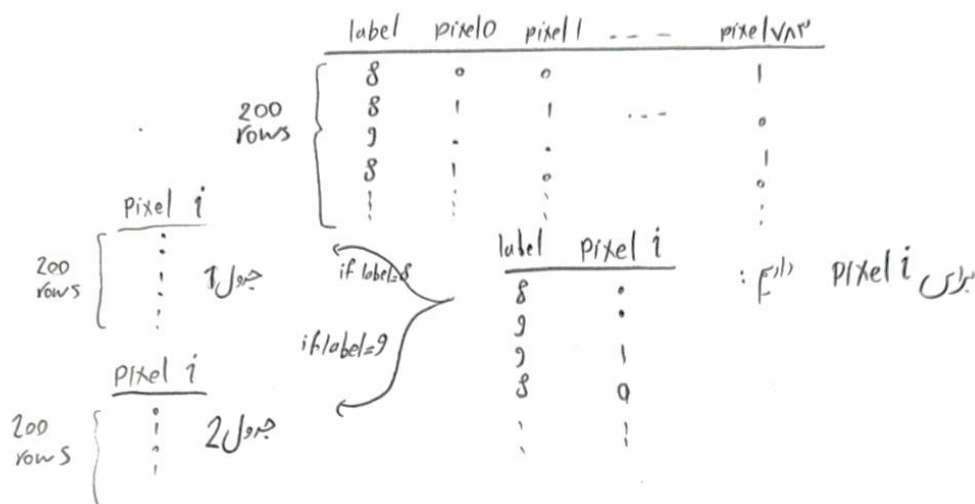


پس همانطور که کدش را نوشتیم نیاز است ضرب داخلی بردار های  $fy$  و  $pny$  را بگیریم که در نامپای به راحتی انجام می شود. یعنی ما بطور تقریبی توزیع پیوسته را گسسته کردیم و برای 1000 تا  $y$  مختلف مقدار  $fy$  و  $pny$  را یافتیم و حالا هم برای تمام آن  $y$  مقدار  $post$  را می یابیم که یعنی بطور حدودی تابع  $Y|N$  را یافته ایم

برای  $post$  هم دیگر کاری نمی ماند و لازم است تقسیم کنیم.

حالا در تابعی که آماده شده بود برای ما، هر بار  $fy$  وارد تابع اپدیت شده و رسماً  $pre$  جدید همنا  $post$  قبلی است و در نهایت در نمودارمان می بینیم که پیک نمودار مرحله به مرحله بطور واضح تری به یک مقدار نهایی اشاره می کند که همان مد نهایی توزیع بتای پسین بوده و احتمالی است که توسط تخمین بیزی به دست ما میرسد (توزیع پسین هم بتا است و بتا یک توزیع مزدوج برای توزیع بتا است) و پس از 100 بار انجام عملیات اپدیت، به مقدار حدودی 0.67 بعنوان احتمال روشن بودن پیکسل 404 میرسیم چون پیک نمودار است و یعنی به احتمال بسیار زیاد، احتمال روشن بودن 404 به شرط  $label = 8$  مساوی 0.67 است.

## (Part5



$$\Rightarrow P(x_i=1 | label=8) = \frac{\text{تعداد 1 های جدول 1}}{\text{تعداد کل سطرهای جدول 1}}$$

$$P(x_i=1 | label=9) = \frac{\text{تعداد 1 های جدول 2}}{\text{تعداد کل سطرهای جدول 2}}$$

$$\rightarrow P(x | label=8) = (P_{x_{18}}, P_{x_{19}}, \dots, P_{x_{10008}})$$

$$P(x | label=9) = (P_{x_{19}}, P_{x_{19}}, \dots, P_{x_{10009}})$$

$$\rightarrow P(x | \text{label}=8) = (P_{u_1|8}, P_{u_2|8}, \dots, P_{u_{n_{\text{vec}}}|8})$$

$$P(x | \text{label}=9) = (P_{u_1|9}, P_{u_2|9}, \dots, P_{u_{n_{\text{vec}}}|9})$$

$$P(\text{label} | x) = \frac{P(x | \text{label}) P(\text{label})}{P(x | 8) P(8) + P(x | 9) P(9)}$$

$$P(\text{label}=8) = P(\text{label}=9) = 1/2 \rightarrow \text{از 200 کلمه داده می‌انیم، دقیقاً 100 کلمه عدد 8 بود و دقیقاً 100 کلمه عدد 9 بود.}$$

$$\rightarrow P(\text{label} | x) = \frac{P(x | \text{label})}{P(x | 8) + P(x | 9)}$$

ما یک ترکیب رندی بنام  $x$  داریم:  $(u_1, u_2, \dots, u_{n_{\text{vec}}})$   
حالا احتمال 8 بودن آنرا می‌خوایم:

$$P(\text{label}=8 | x) \propto P(x = (u_1, \dots, u_{n_{\text{vec}}}) | \text{label}=8)$$

پس احتمال آن ترکیب رندی به شش 8 بستن دلیل راس خواتیم پس در این احتمال برای  $\text{label}=8$  هست و به ازای هر  $u_i$  دیتا هست  $x$ ، اگر  $u_i=8$  بود که خود کلمه موجود در دیتا که احتمال 1 شدن هست را فریب می‌دهیم، و اگر  $u_i=9$  بود  $1-P$  که  $P$  موجود در دیتا احتمال است را فریب می‌دهیم.

اما چون اعداد اولیه این 255 بوده اند، یکبار  $\text{thershold}$  را احتمال می‌دهیم و پس  $x$  خوانده شده

در نهایت به این نتیجه رسیدیم:

for 201th pic, probability of being 8 is about 0.9997406457898158

and probability of being 9 is about 0.00025935421018423725

so This picture is 8

for 202th pic, probability of being 8 is about 5.915116612170155e-19

and probability of being 9 is about 1.0

so This picture is 9

دلیل اینکه در 202امین عکس، احتمال 9 بودن را 1.0 گزارش کرده است این است که آن عدد انقدر نزدیک به 1 بوده است که پایتون توانایی بیان آنرا نداشته و تقریباً 1 در نظر گرفته است، مگر نه چون عکس هایی دقیقاً مساوی با آنهایکه داشتیم، جزو داده های train نبودند، پاسخگویی به احتمال 1 تقریباً غیرممکن است.