



دانشگاه تهران، دانشکده مهندسی برق و کامپیوتر آمار و احتمال مهندسی

تمرین کامپیوتری دوم - موضوع تمرین

طراح: متین بذرافشان - علی آریایی

سوپروایزر: علی محمدی

تاریخ تحویل: ۲۴ آذر ۱۴۰۲

(۴۰) نمره

۱. توزیع شرطی

دیتاست tarbiat.csv شامل تعداد گذرهای اتوبوس BRT و مترو از ایستگاه های تربیت مدرس در بازه های زمانی مشخص است.

۱. برای هر دو ستون، با استفاده از *matplotlib* یک هیستوگرام رسم کنید.

۲. اگر متغیر تصادفی X را تعداد گذرهای مترو و متغیر تصادفی Y را تعداد گذرهای BRT تعریف کنیم، با توجه به آموخته هایمان از درس و ماهیت مسئله و همچنین شکل نمودارهایی که رسم کرده اید، این متغیرها از چه توزیعی پیروی می کنند؟ پارامتر (پارامترهای) هر کدام را محاسبه کنید.

۳. برای مشاهده توزیع مقادیر گسسته، می توان از density histogram استفاده کرد. این نمودار را برای ستون metro رسم کنید.

۴. حال با استفاده از *Scipy*، نمودار توزیع X را با پارامترهایی که بدست آورده اید را به نمودار بالا اضافه کنید و از درستی توزیع و پارامترهایی که مسئله تان را با آن مدل کرده اید، مطمئن شوید.

۵. متغیر تصادفی $Z = X + Y$ تعریف می کنیم، این متغیر از چه توزیعی برخوردار است؟ پارامتر (پارامترهای) آن را بدست آورید. با استفاده از *Scipy* نمودار توزیع Z را رسم کنید. به این نمودار، نمودار density histogram مجموع دو ستون مترو و BRT را اضافه کنید و از درستی توزیع و پارامترهای بدست آمده اطمینان حاصل کنید.

۶. متغیر تصادفی $W \sim P(X|X + Y = n)$ را تعریف می کنیم. این توزیع و پارامترهای آن را به صورت تئوری محاسبه کنید.

۷. تابع جرمی احتمال W به ازای $n = 8$ ، با انتخاب توزیع مناسبی در *Scipy* و دادن پارامترهایی که بدست آورده اید به آن، رسم کنید.

۸. حال چگالی تعداد متروها را به شرط آنکه در مجموع ۸ مترو و BRT از ایستگاه های تربیت مدرس گذشته باشد را به نمودار بالا اضافه کنید. چه مشاهده می کنید؟

۲. تابع مولد گشتاور

(۴۰) نمره

مسئله جمع کننده کالابریک

مسئله‌ی جمع‌کننده‌ی کالابریک^۱ یکی از مسائل معروف در آمار و احتمال است. فرض کنید n نوع کالابریک داریم و از هرکالابریک نیز به اندازه کافی موجود است. حال متغیر تصادفی X را اینگونه تعریف می‌کنیم: تعداد کالابریک‌هایی که باید مشاهده کنیم به طوری که از هر نوع کالابریک حداقل یکی دیده باشیم. برخلاف ظاهر ساده مسئله، توزیع متغیر تصادفی X پیچیده می‌باشد.

روش مونته کارلو^۲: روش مونته کارلو، یک روش عددی در حل مسائل دارای محاسبات سخت و طولانی است. این روش با استفاده از نمونه‌گیری‌های تصادفی متعدد، پاسخ مسئله را پیدا می‌کند.

۱- می‌توان برای بدست آمدن میانگین X از روش مونته کارلو استفاده کرد. بدین منظور، تابعی بنویسید که n و k را دریافت کند، سپس k بار مسئله جمع‌کننده کالابریک را به ازای n حل کند، سپس میانگین جواب‌ها را برگرداند.

۲- به ازای $n = 10, 100, 1000$ و $k = 10$ ، تابع را اجرا کنید. مقادیر به چه عددی همگرا می‌شوند؟

Intro to SymPy

کتابخانه SymPy یکی از کتابخانه‌های قوی پایتون در زمینه ریاضیات و معادلات آن است. به‌وسیله این کتابخانه می‌توانید یک متغیر را به عنوان یک نماد^۳ تعریف کنید. یک نماد مقدار مشخصی ندارد در نتیجه می‌توان به کمک آن عملیات‌های ریاضی مانند مشتق و انتگرال و ... را صورت غیر عددی را انجام داد.

یک نماد را اینگونه تعریف می‌کنیم:

```
import sympy
x = sympy.symbols('x')
```

حال می‌توان به کمک آن نماد، یک معادله تشکیل داد:

```
equation = x ** 2 + 5 * x
```

می‌توان نماد را با هر عدد یا نماد دلخواهی جایگزین کرد:

```
equation.subs({x:10})
```

همچنین می‌توان عملیات‌های ریاضی را به صورت نمادین انجام داد، برای مثال می‌توان از یک معادله مشتق گرفت:

```
deriv = sympy.diff(equation, x)
```

و همچنین می‌توان مقدار آن را به ازای عدد یا نماد مورد نظر نمایش داد:

```
deriv.subs({x:10})
```

^۱ coupon collector's problem
^۲ Monte Carlo
^۳ symbol

یکی از روش‌های حل مسئله جمع‌کننده کالابریگ کاهش مسئله به n مسئله مستقل است. متغیر تصادفی X_i را تعداد کالابریگ‌هایی که باید مشاهده کنیم تا کالابریگ نوع i ام برای اولین بار مشاهده شود تعریف می‌کنیم، این متغیر دارای چه توزیعی است؟ با توجه به اینکه مسئله دارای ماهیت برنولی (موفقیت – شکست) بوده، آیا احتمال مشاهده کالابریگ نوع i برای اولین بار برای همه i های ۱ تا n یکسان است؟

$$X = X_1 + X_2 + X_3 + \dots + X_n$$

۳- با استفاده از *sympy*، متغیر s را به عنوان symbol تعریف کنید. تابع مولد گشتاور متغیر تصادفی X_i را به ازای $n = 10$ و i از ۱ تا n تعریف کنید.

۴- با استفاده از خواص تابع مولد گشتاور، تابع مولد گشتاور X را بدست آورید.

۵- با استفاده از خواص تابع مولد گشتاور، میانگین متغیر تصادفی X را بدست آورید. پاسخ شما باید با پاسخ قسمت ۲ یکسان شود.

۳. تخمین بیزی و استنباط بیزی

(۲۰+۲۰) نمره

دیتاست digits.csv دارای ۲۰۲ عدد دست نویس (۱۰۱ عدد ۸ و ۱۰۱ عدد ۹) به اندازه 28×28 می باشد، به طوریکه مقدار هر پیکسل (که عددی بین ۰ و ۲۵۵ می باشد) آمده است. همچنین عدد متناظر نیز در ستون label آمده است.

۱. اعداد ۲۰۱ ام و ۲۰۲ ام را در متغیرهایی ذخیره کنید و سپس از دیتافریم حذف کنید.

۲. برای ساده سازی، با قرار دادن threshold به اندازه ۱۲۸ اعداد را binary کنید به طوری ۱ به معنی روشن بودن و ۰ به معنی خاموش بودن هر پیکسل باشد. (توجه داشته باشید، label ها باید بدون تغییر باقی بمانند.)

۳. یک عدد به طور دلخواه انتخاب کنید و با استفاده از reshape آن را به آرایه دو بعدی با اندازه 28×28 تبدیل کرده و سپس با استفاده از تابع *matplotlib.pyplot.imshow* آن را نمایش دهید.

۴. متغیر تصادفی Y را احتمال روشن بودن پیکسل ۴۰۴ ام به شرط $\text{label} = 8$ و متغیر تصادفی برنولی N را روشن بودن یا نبودن پیکسل ۴۰۴ ام تعریف میکنیم. حال می‌خواهیم به وسیله تخمین بیزی، تغییرات توزیع Y را پس از مشاهده هر نمونه مشاهده کنیم. ابتدا به دلیل نداشتن اطلاعات اولیه، Y دارای توزیع یکنواخت (توزیع بتا با $a = 1$ و $b = 1$) می‌باشد. حال با مشاهده هر داده با توجه به توزیع پیشین (prior) متغیر تصادفی Y و همچنین مشاهدات خود از نمونه جدید (likelihood)، توزیع پسین (posterior) متغیر تصادفی Y را به دست آوریم. سپس در مشاهده داده بعدی، توزیع پسین به دست آمده، خود توزیع پیشین می‌شود و این روند تا دیدن آخرین داده ادامه دارد. بدین منظور شما کد کافیست کد الحاق شده را کامل کنید و سپس آن را اجرا کنید تا تغییرات توزیع متغیر تصادفی Y را مشاهده کنید. روابط مورد نیاز، در پایین آمده است. پس از اتمام، اگر بخواهیم یک عدد را به عنوان احتمال روشن بودن این پیکسل اعلام کنید، چه عددی را بیان می‌کنید؟ این عدد همان مد توزیع بتا پسین یا همان میانگین این پیکسل است.

$$f(Y|N) = \frac{P(N=n|Y=p)f(Y)}{\int_0^1 P(N=n|Y=p)f(Y)}$$

$$\int_a^b f(x)dx \approx \frac{b-a}{N} \sum_{i=1}^N f(x_i); x_i = a + i \frac{b-a}{N}$$

امتیازی:

۵. متغیر تصادفی X را به صورت برداری از متغیرهای مستقل تعریف میکنیم به طوری که هر x_i احتمال روشن بودن پیکسل i ام را مشخص می کند. حال می‌خواهیم $P(X|\text{label})$ را به دست آوریم. برای این کار می‌توانید، میانگین پیکسل ها را برای هر پیکسل

منحصر بفرد، یکبار برای داده ها با $label = ۸$ و یکبار برای داده ها با $label = ۹$ به دست آورید. در انتها شما باید دو آرایه با اندازه ۷۸۴ داشته باشید.

۶. حال می‌خواهیم با استفاده قانون بیز، $P(label|X)$ را برای اعداد ۲۰۱م و ۲۰۲م به دست آوریم. چهار مقدار $p(label = ۸|X^{۲۰۲})$ و $p(label = ۸|X^{۲۰۱})$ و $p(label = ۹|X^{۲۰۱})$ و $p(label = ۹|X^{۲۰۲})$ را محاسبه کنید. با توجه به این که پیکسل‌ها را مستقل از هم در نظر گرفتیم و همچنین برنولی بودن توزیع هر x_i می‌توان از فرمول زیر استفاده کرد:

$$p(label|X) = \frac{p(X|label)P(label)}{\sum_{label_i} p(X|label_i)P(label_i)}$$

$$p(X|label) = \prod_{i=1}^{784} p(x_i|label)^{x_i} (1 - p(x_i|label))^{(1-x_i)}$$