

بسم الله الرحمن الرحيم

پروژه نهایی درس مبانی علوم داده
دکتر بهرک و دکتر یعقوب زاده

محمد امین توانایی - ۸۱۰۱۰۱۳۹۶

سید علی تهامی - ۸۱۰۱۰۱۳۹۷

مهدی وجهی - ۸۱۰۱۰۱۵۵۸

فهرست

2.....	مقدمه
3.....	فاز صفرم : انتخاب داده ها
3.....	جمع آوری داده ها
3.....	پردازش داده ها
4.....	اضافه کردن فاصله از تعطیلی به داده ها
4.....	جمع آوری روز های تعطیل
6.....	محاسبه فاصله هر روز از تعطیلی
7.....	فاز اول : بررسی آماری
7.....	نگاه کلی
7.....	بررسی همبستگی داده ها
9.....	بررسی عوامل موثر بر عدم دریافت غذا
9.....	دانشکده های مختلف
10.....	خوابگاه و دانشکده
10.....	جنسیت
10.....	تعطیلی های هفته بعد، قبل و همان هفته
11.....	وعده غذایی
11.....	نوع غذا
12.....	بررسی داده های پرت
13.....	فاز دوم : پیش بینی
13.....	Feature Engineering and Selection
13.....	Dimensionality Reduction
13.....	Evaluation Metric
14.....	آموزش مدل ها
14.....	شبکه عصبی
14.....	سایر مدل ها
15.....	مقایسه
15.....	Feature Analysis
16.....	چالش ها

مقدمه

دانشگاه‌ها به عنوان مراکز علمی و آموزشی، نقش مهمی در تربیت نسل آینده جامعه ایفا می‌کنند. در کنار آموزش، دانشگاه‌ها خدمات رفاهی مختلفی از جمله ارائه غذا به دانشجویان را نیز ارائه می‌دهند. سلف‌های دانشگاهی به عنوان محل صرف غذا، نقش مهمی در تأمین نیازهای غذایی دانشجویان دارند. با این حال، در بسیاری از دانشگاه‌ها شاهد هدر رفتن مقادیر قابل توجهی غذا در سلف‌ها هستیم.

عدم دریافت غذای سلف دانشگاه، مشکلی است که از جنبه‌های مختلف مالی و زیست‌محیطی حائز اهمیت است. از نظر مالی، هدر رفتن غذا به معنای اتلاف منابع و سرمایه‌های دانشگاه و به تبع آن، دانشجویان است. هزینه‌های خرید مواد اولیه، پخت و پز و توزیع غذا، همگی صرف تهیه غذاهایی می‌شوند که در نهایت دور ریخته می‌شوند. این امر نه تنها به ضرر دانشگاه و دانشجویان است، بلکه به طور کلی اتلاف منابع ملی را به دنبال دارد. از نظر زیست‌محیطی نیز، هدر رفتن غذا پیامدهای منفی قابل توجهی دارد. تولید هر واحد غذا، مستلزم مصرف منابع طبیعی مانند آب، خاک و انرژی است. دور ریختن غذا به معنای هدر رفتن بیهوده این منابع ارزشمند است. علاوه بر این، تجزیه ضایعات غذایی در محل‌های دفن زباله، متان تولید می‌کند که یک گاز گلخانه‌ای قدرتمند است و به گرم شدن کره زمین دامن می‌زند.

بنابراین، پیش‌بینی میزان عدم دریافت غذا در سلف‌های دانشگاه و یافتن راهکارهایی برای کاهش آن، از اهمیت زیادی برخوردار است. با انجام این کار، می‌توان به حفظ منابع مالی و طبیعی، ارتقای سلامت جامعه و ایجاد محیطی پایدارتر کمک کرد.

در این گزارش، به بررسی روش‌های مختلف پیش‌بینی میزان عدم دریافت غذا در سلف‌های دانشگاه و ارائه راهکارهایی برای کاهش آن خواهیم پرداخت. امیدواریم این گزارش بتواند گامی در جهت حل این معضل مهم بردارد.



فاز صفرم : انتخاب داده ها

جمع آوری داده ها

با صحبت هایی که با دکتر یعقوب زاده انجام دادیم، ایشان نامه ای خطاب به معاونت دانشجویی و رونوشت به معاونت پژوهشی ارسال کردند و در آن خواستار دریافت اطلاعات سامانه تغذیه شدند. بعد از مدتی داده ای نمونه دریافت کردیم که ردیف های آن به ازای هر رزرو دانشجوی بود. سپس با پیگیری هایی که انجام شد داده ها را در بازه آذر ماه ۱۴۰۲ تا ابتدای خرداد ۱۴۰۳ در اختیار ما قرار دادند.

پردازش داده ها

ابتدا به قالب اولیه داده ها نگاهی بیندازیم.

PersonId	ReserveId	Gender	Price	DateReserve	Name	FoodName	Count	ReserveStatus	CollegeName	CollegeCode	FieldCode	FieldName	
0	44342	6381703	مرد_x000D_\n	100000	1402/11/1	ناهار	چلوکباب کوبیده+دوغ+کره (تک نفره)+نارنج (50 گرم)	1	دریافت شده	دامپزشکی	75.0	750112	دامپزشکی
1	141693	6405451	مرد_x000D_\n	100000	1402/11/1	شام	سبزی پلو یا تن ماهی+خرما+زیتون (40 گرم)	1	دریافت شده	دانشکده گان فنی	81.0	810164	مهندسی برق - سیستم های قدرت
2	112376	6352639	مرد_x000D_\n	100000	1402/11/1	ناهار	چلوکباب کوبیده+دوغ+کره (تک نفره)+نارنج (50 گرم)	1	دریافت شده	حقوق و علوم سیاسی	21.0	210412	حقوق

DegreeCode	DegreeName	RestaurantName	GroupName	EducationSession	PersonType	Comment	Reception	ReceptionType
45	دکتری عمومی	دانشکده-سلف دامپزشکی تهران	دانشکده دامپزشکی	نوبت دوم	دانشجو	ممتنع	NaN	1
15	کارشناسی ارشد ناپیوسته	خوابگاه-کوی برادران	خوابگاه کوی دانشگاه	روزانه_x000D_\n	دانشجو	ممتنع	NaN	1
10	کارشناسی پیوسته	دانشکده-سلف مهر	دانشکده حقوق	روزانه_x000D_\n	دانشجو	بدون نظر	NaN	1

همانطور که می بینید این داده مشکلات زیر را دارد:

- ستون های بدون کاربرد
 - ستون ها با مقادیری به شکل نادرست
 - نام غذا ها همراه مخلفات است
 - وضعیت رزرو در حالت عدم دریافت ۳ مقدار دارد
- در نهایت در این بخش با رفع مشکلات بالا و موارد مشابه آن و همچنین اعمال تغییرات زیر به داده نهایی خود می رسیم.

- تبدیل وضعیت دریافت هر دانشجوی به دریافت با کارت یا کد
- ادغام خوابگاه های سطح شهر
- ایجاد ستون سلف خوابگاهی یا دانشگاهی
- یکسان کردن نام سلف های یکسان در اصلاح نام سلف ها
- مشخص کردن برنجی یا خوراک بودن غذا
- تبدیل ساختار تاریخ از متن به تاریخ

- اضافه کردن تاریخ میلادی

DateReserve	RestaurantName	RestaurantType	Meal	FoodName	FoodType	Gender	ReceiveWithCard	ReceiveWithCode	DontReceive	Reservation	DayOfWeek
1402/11/1	ابورحان	daneshgah	dinner	خوراک کوکو سیب زمینی	khoraak	man	14	0	0	14	1
1402/11/1	ابورحان	daneshgah	dinner	خوراک کوکو سیب زمینی	khoraak	woman	27	0	0	27	1
1402/11/1	ابورحان	daneshgah	dinner	سبزی بلو یا تن ماهی	berenji	man	126	0	0	126	1
1402/11/1	ابورحان	daneshgah	dinner	سبزی بلو یا تن ماهی	berenji	woman	110	0	0	110	1
1402/11/1	ابورحان	daneshgah	lunch	خوراک آبگوشت	khoraak	man	19	0	0	19	1

حال کار تمام شده اما یک عاملی که به نظر موثر می آید روز های تعطیل است پس باید به صورتی آنها را به پروژه اضافه کنیم.

اضافه کردن فاصله از تعطیلی به داده ها

جمع آوری روز های تعطیل

ابتدا یک جستجو می کنیم که کتابخانه ی مربوط به این موضوع را پیدا کنیم اما چیزی پیدا نمی کنیم. در مرحله بعد به دنبال داده های روز های تعطیل در اینترنت می گردیم ولی چیزی پیدا نمی کنیم در نهایت تصمیم گرفتیم از سایت time.ir تعطیلی ها را استخراج کنیم. برای این کار بخش تقویم سالانه را در سایت پیدا کردیم و تصمیم گرفتیم از آن صفحه برای استخراج استفاده کنیم. در ادامه تصویری از آن صفحه را مشاهده می کنید.



در این مرحله کافیت خانه های قرمز تقویم را انتخاب کنیم و بعد تاریخ آنها را بخوانیم. برای این کار به صورت زیر عمل می کنیم.

```
a = soup.find_all('div', class_='holiday')[1]
```

حال که آن ها را انتخاب کردیم باید مقادیر هر کدام را به دست بیاوریم که با توجه به ساختار سایت به صورت زیر عمل می کنیم.

```
d = a.find('div', class_='jalali').text
tmp = a.parent.parent.parent.parent.find('a', class_='jalali').text
m, y = tmp.split(' ')
```

که در اینجا ما روز ماه و سال را به دست آوردیم حال آن را به کلاس جلالی تبدیل می کنیم. برای این کار هم به صورت زیر عمل می کنیم.

```
month = {
    'فروردین': 1,
    'اردیبهشت': 2,
    'خرداد': 3,
    'تیر': 4,
    'مرداد': 5,
    'شهریور': 6,
    'مهر': 7,
    'آبان': 8,
    'آذر': 9,
    'دی': 10,
    'بهمن': 11,
    'اسفند': 12
}
JalaliDate(int(y), month[m], int(d))
```

محاسبه فاصله هر روز از تعطیلی

حال کلاسی تعریف می کنیم که همین کار با یکسری موارد تکمیلی انجام دهد. در نهایت خروجی آن لیستی از روز های تعطیل است. حال باید فاصله هر روز را با تعطیلی ها به دست بیاوریم برای این کار نیز کلاسی می نویسیم که خروجی ای مانند تصویر دارد.

	HolidayInWeekCount	HolidayInPrevWeekCount	HolidayInNextWeekCount	NextHoliday	PreviousHoliday
1402-11-01	0	0	0	[21, 35]	[35, 110]
1402-11-10	0	0	0	[12, 26]	[44, 119]
1402-11-11	0	0	0	[11, 25]	[45, 120]
1402-11-12	0	0	0	[10, 24]	[46, 121]
1402-11-13	0	0	0	[9, 23]	[47, 122]
1402-11-14	0	0	1	[8, 22]	[48, 123]
1402-11-02	0	0	0	[20, 34]	[36, 111]
1402-11-21	1	0	0	[1, 15]	[55, 130]
1402-11-23	1	0	0	[13, 36]	[1, 57]
1402-11-24	1	0	0	[12, 35]	[2, 58]
1402-11-25	1	0	0	[11, 34]	[3, 59]
1402-11-26	1	0	0	[10, 33]	[4, 60]
1402-11-28	0	1	1	[8, 31]	[6, 62]
1402-11-29	0	1	1	[7, 30]	[7, 63]
1402-11-03	0	0	0	[19, 33]	[37, 112]
1402-11-30	0	1	1	[6, 29]	[8, 64]
1402-11-04	0	0	0	[18, 32]	[38, 113]
1402-11-05	0	0	0	[17, 31]	[39, 114]
1402-11-06	0	0	0	[16, 30]	[40, 115]
1402-11-07	0	0	0	[15, 29]	[41, 116]
1402-11-08	0	0	0	[14, 28]	[42, 117]
1402-11-09	0	0	0	[13, 27]	[43, 118]

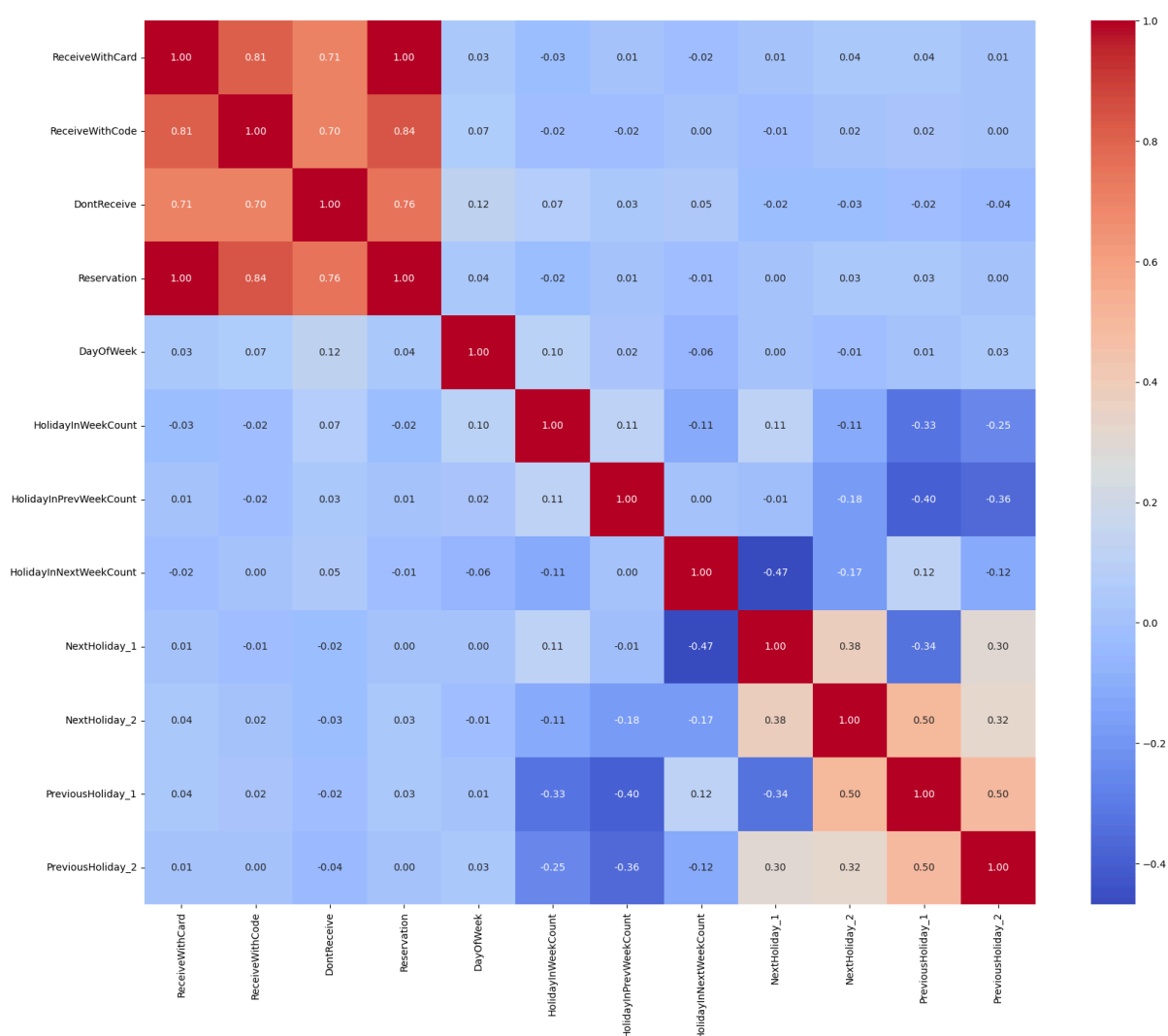
در آخر این موارد را با جدول خود ادغام می کنیم.

فاز اول : بررسی آماری

نگاه کلی

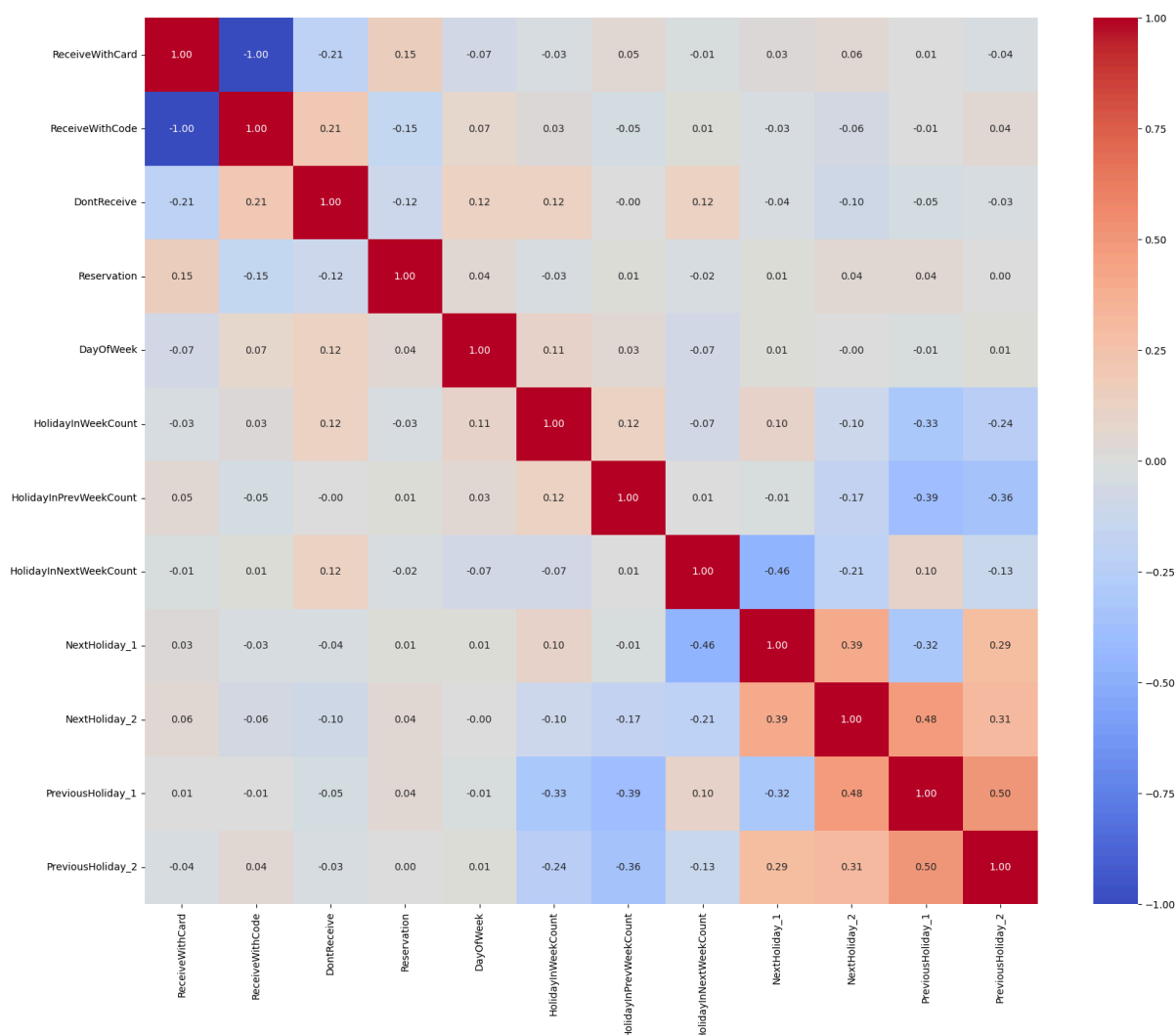
ابتدا یک بررسی کلی روی داده ها انجام می دهیم. در مجموع 162,429 پرس غذا گرفته نشده که در مدت ۵ ماه است. (از ۶ ماه داده یک ماه عید و رمضان بوده و در مجموع در ۱۴۲ روز غذا سرو شده) این عدد با توجه به این که قیمت غذا ها حدود ۸۵ هزار تومان برای دانشگاه می شود معادل ۱۳.۸ میلیارد تومان و با در نظر گرفتن حداقل وزن ۳۰۰ گرمی معادل ۴۹ تن غذا می شود. در نهایت این میزان تقریباً ۸ درصد کل غذاهای پخت شده است.

بررسی همبستگی داده ها



همانطور که می بینید بین تعداد رزرو ها و تعداد عدم دریافت ها و باقی موارد مربوط به دریافت همبستگی وجود دارد که این موضوع احتمالاً کمکی نکند زیرا واضح است که با افزایش تعداد رزرو ها هم تعداد دریافت با کارت و

هم تعداد دریافت با کد فراموشی افزایش می یابد. همین اتفاق برای تعطیلی ها هم می افتد. اما از این نمودار می توان رابطه ۱۲ درصدی روز هفته با عدم دریافت را متوجه شد. حال اگر سفارش ها را به صورت درصد نمایش دهیم به صورت زیر می شود.



همانطور که می بینید عدم دریافت ها با دریافت با کد رابطه ۲۱ درصدی دارد. دلیل این موضوع هم این است که وقتی بچه ها دانشگاه نمی آیند یا غذا را نمی گیرند و یا در کانال های کد فراموشی ارسال می کنند پس با افزایش افرادی که دانشگاه نیامدند این دو مقدار افزایش پیدا می کند. همچنین رابطه ۱۲ درصدی با تعطیلی هفته دارد که دلیل آن واضح است. جالب اینجاست که تعطیلی هفته بعد نیز به همان مقدار موثر است. یک مورد عجیب دیگر این است که فاصله دومین تعطیلی منفی با عدم دریافت منفی ۱۰ درصد است ولی اولین تعطیلی منفی ۴ درصد. در آخر هم می بینیم که درصد دریافت با کارت رابطه مثبت ۱۵ درصدی با تعداد رزرو ها دارد که این به این دلیل است که با افزایش تعداد رزرو یعنی روز کاری ای داریم و افراد در دانشگاه هستند.

بررسی عوامل موثر بر عدم دریافت غذا

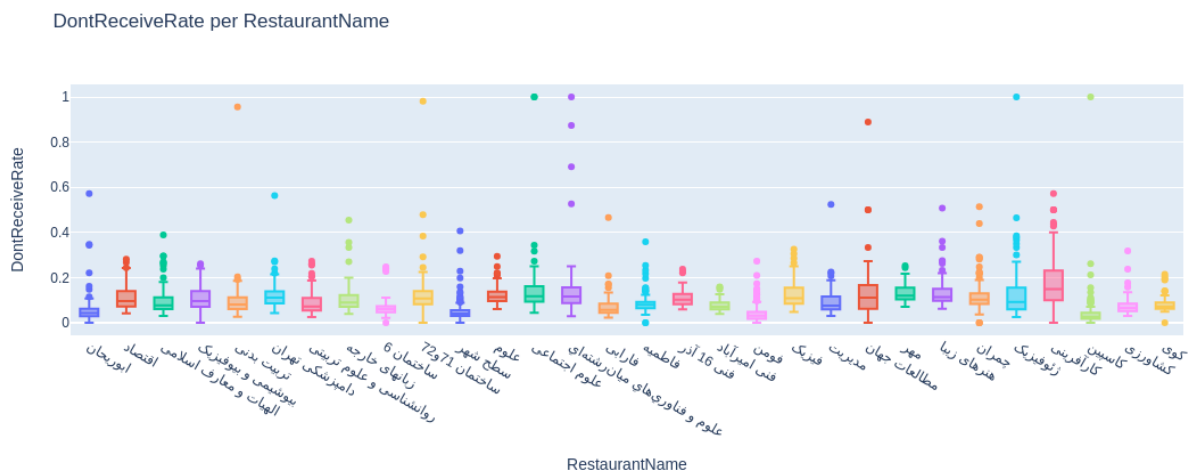
حال به بررسی عوامل مختلف و تاثیر آنها بر عدم دریافت غذا می پردازیم. فقط قبل آن لازم است توضیح دهیم که تست فرض چگونه تعریف می کنیم. در هر قسمت روی دو بخش مجموعه داده تست فرض را به شکل زیر تعریف می کنیم:

$$H_0 : Rate_x - Rate_y = 0$$

$$H_A : Rate_x - Rate_y < 0$$

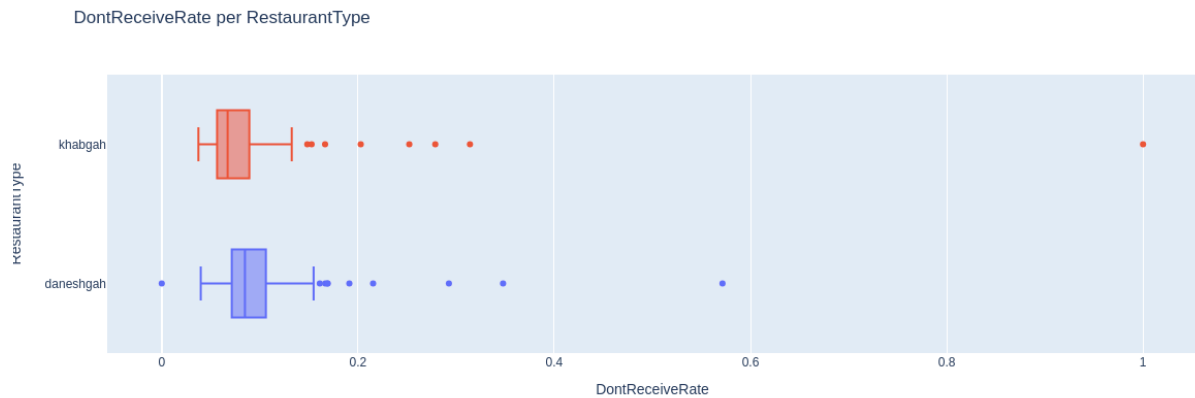
یعنی در فرض اولیه در نظر می گیریم که دسته ها تفاوتی از نظر درصد عدم دریافت تفاوتی ندارند ولی در فرض ثانویه می گوییم که تفاوت دارند.

دانشکده های مختلف



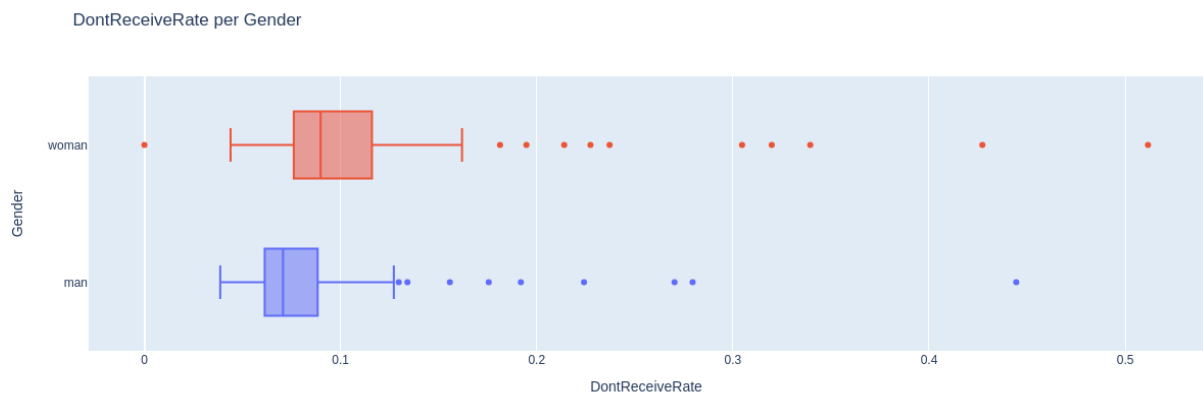
همانطور که می بینید رفتار دانشکده های مختلف متفاوت است مثلاً دانشکده کارآفرینی میانه ۱۴ درصد و چارک سوم ۲۳ درصد عدم دریافت غذا دارد که زیاد است. اما از آن طرف دانشکده ای مثل فارابی میانه ۶ درصد دارد که مقدار کمی به شمار می رود.

خوابگاه و دانشکده



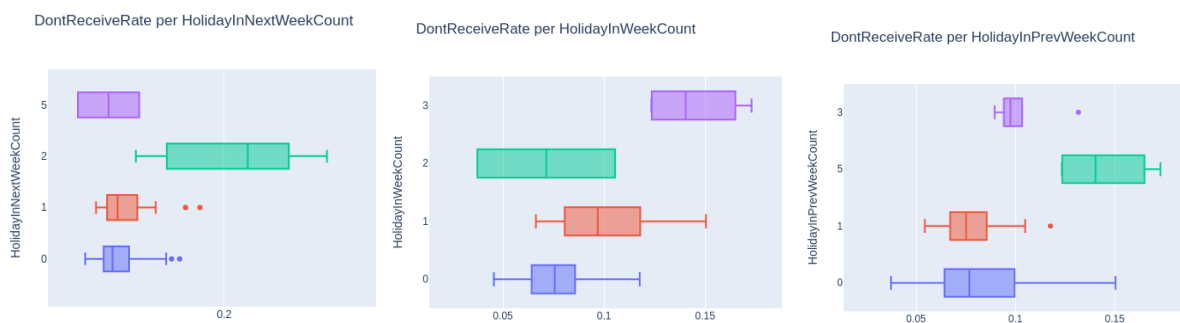
به طرز عجیبی طبق تست فرضی که زدیم مشاهده کردیم که تفاوتی میان میزان عدم دریافت غذا در خوابگاه و دانشگاه وجود ندارد البته که با توجه به $p\text{-value } 6.1\%$ این موضوع قویا رد نمی شود و خیلی هم موثر نیست.

جنسیت



با توجه به مقدار $p\text{-value } 0.000192$ می توان گفت که جنسیت بر عدم دریافت غذا قویا موثر است.

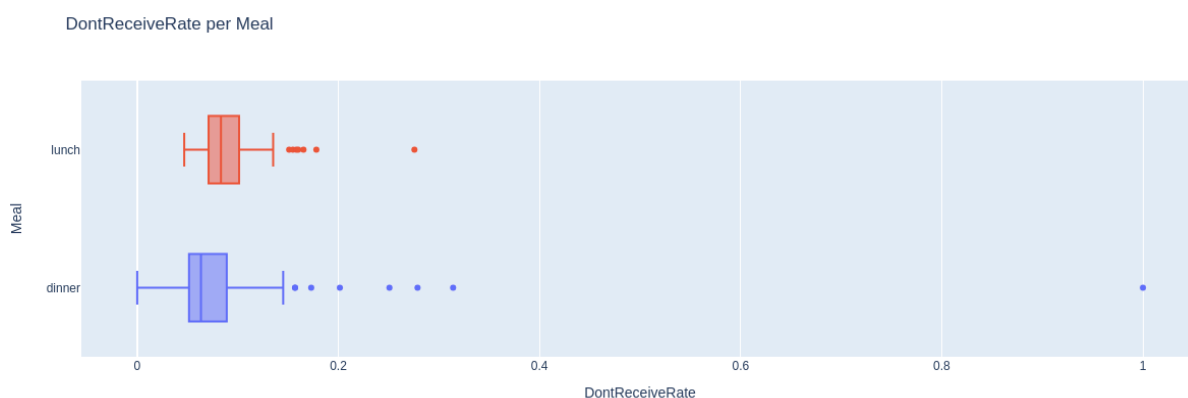
تعطیلی های هفته بعد، قبل و همان هفته



هفته	p-value
قبل	0.6598
همان هفته	0.00155
بعد	0.0009

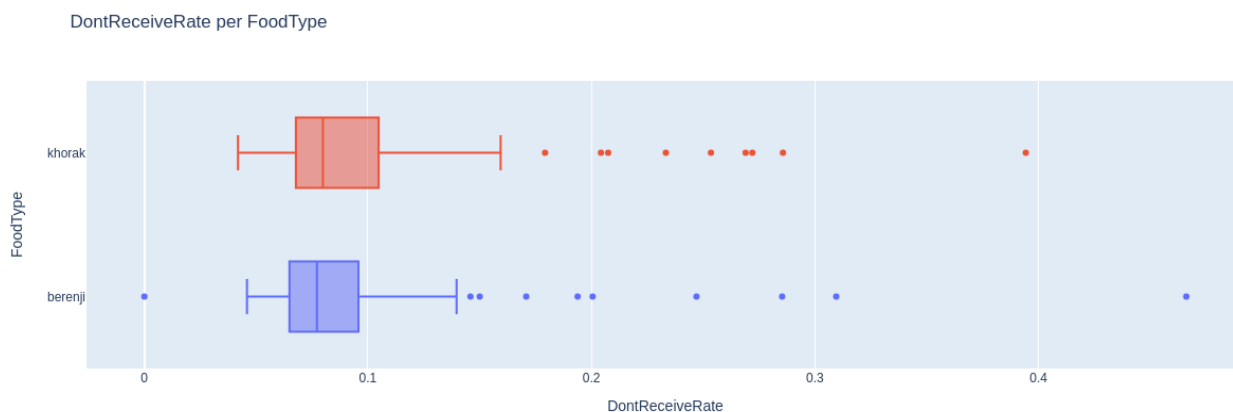
همانطور که می بینید برخلاف تعطیلی های هفته گذشته که بر عدم دریافت موثر نیست هفته بعد و همان هفته بر عدم دریافت موثر است.

وعده غذایی



طبق تست فرض ارتباطی بین نوع غذا و عدم دریافت وجود ندارد. (p-value 0.144)

نوع غذا



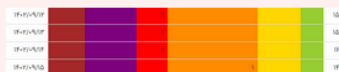
نوع غذا نیز بر عدم دریافت موثر نیست. (p-value 0.143)

بررسی داده های پرت

در ادامه ۳ روز که بیشترین عدم دریافت است را یک بار بر اساس درصد مشاهده می کنیم و یکبار بر اساس تعداد.

TOP-3 NOT RECEIVED BASED ON NUMBER

وضعیت آب و هوا، امروز ۸ اسفند ۱۴۰۲ / تداوم بارش برف و باران / شب گذشته سردترین شب ۶ ماه اخیر
یک گزارش هواشناسی گفت: شب گذشته در بسیاری از نقاط کشور سردترین شب سال در ۶ ماه اخیر بود.
پایگاه خبری تحلیلی انتخاب (Entekhab.ir)



2,661

1402-12-08
**COLD WEATHER,
SNOWFALL, HOLIDAYS**

2,443

1402-09-(13-14-15)
AIR POLLUTION

2,319

1403-03-01
**MARTYRDOM OF
PRESIDENT**

TOP-3 NOT RECEIVED BASED ON PERCENTAGE

ماه	هفته	شنبه	یکشنبه
شماره	ثبت نام ۱۴۰۲-۲	(۰) شروع ثبت نام و انتخاب واحد	
۱	هفته اول	۲۶ شروع کلاس ها	۲۲ پروازی انقلاب اسلامی



47%

1402-11-21
**THE FIRST DAY OF THE
TERM AND THE FIRST DAY
OF THE WEEK AND
BETWEEN HOLIDAYS**

25%

1402-12-(21,25)
LAST DAYS OF THE YEAR

16%

1403-01-(14-15)
**TUESDAY MORNING,
FARVARDIN 14**

همانطور که می بیند موارد درصدی به علت تعطیلی هستند و کاملاً قابل پیشبینی و در موارد تعدادی همه رویداد های غیر منتظره هستند.

فاز دوم : پیش بینی

Feature Engineering and Selection

در این قسمت به دلیل اینکه در آن روز دیتای دریافتی با کد و دریافتی با کارت را نداریم که بخواهیم پیشبینی کنیم به همین دلیل نمی توان از آنها استفاده کرد.

سپس با طبق تست فرض هایی که در فاز اول پروژه زدیم متوجه شدیم که ستون Holiday In Pre Week Count با ستون متغیر رزرو نشده ها ارتباطی نداشته و میتوان آن را حذف کرد.

سپس برای اینکه ستون هایی که به تعطیلی ها ربط داشتند را به دلیل اینکه در دیتا فریم ما دارای معنا و مفهوم باشند به صورت بولین تعریف کردیم زیرا برای مثال فرقی ندارد فاصله تعطیلی بعدی ۳۹ باشد یا ۴۰ هر دو یک معنا و مفهوم را میسرانند.

سپس ستون های categorical را one hot کردم.

یک دیتایی که بسیار مفید است تعداد رزرو دانشکده در روز های قبل و بعد آن روز است برای همین این دو فیچر را نیز به هر ردیف اضافه کردم.

Dimensionality Reduction

Pca با تبدیل به دو فیچر ۴.۹۱ درصد واریانس را پوشش میدهد.

Pca با نگه داشتن ۹۵ درصد واریانس ۱۳۷ فیچر را به ۱۱۸ فیچر تبدیل کرد.

Evaluation Metric

دلیل استفاده از r^2_score اینست که چون معیاری است که همیشه بین صفر و یک است برای ارزیابی مدل های مختلف و مقایسه عملکرد آنها بهتر است.

دلیل استفاده از $r^2_score_adjust$ که برای اضافه کردن فیچر های مختلف اگر با متغیر پیشبینی شونده ارتباطی نداشته باشند جریمه در نظر میگیرد و با بالا رفتن فیچر ها عدد آن لزوما بهتر نمیشود.

از $rmse$ استفاده کردیم زیرا واحد آن با واحد متغیر پیشبینی شونده یکی است و میتواند به ما مقدار خطا به ازای هر داده را توضیح دهد.

آموزش مدل ها

شبکه عصبی

معماری شبکه: مدل یک مدل Sequential با چندین لایه Dense است که هر کدام از آن‌ها توسط یک لایه Dropout دنبال می‌شود تا از overfitting جلوگیری کند. معماری با یک لایه 256 نرونی شروع می‌شود و به تدریج تعداد نرون‌ها در لایه‌های بعدی کاهش می‌یابد.

تابع فعال‌سازی: PReLU (واحد خطی اصلاح‌شده پارامتری) به عنوان تابع فعال‌سازی برای لایه‌های پنهان استفاده می‌شود. این انتخاب به مدل اجازه می‌دهد تا ضرایب rectifiers را یاد بگیرد، که انعطاف‌پذیری کمی را نسبت به تابع ReLU استاندارد فراهم می‌کند.

لایه خروجی: لایه نهایی دارای یک نرون با تابع فعال‌سازی خطی است، که برای وظایف رگرسیونی که خروجی آن‌ها یک مقدار پیوسته است، معمول است.

تابع زیان: MSE (Mean Squared Error) به عنوان تابع زیان استفاده می‌شود، که برای مسائل رگرسیونی استاندارد است زیرا اشتباهات بزرگ‌تر را بیشتر جریمه می‌کند و هدف آن کمینه کردن میانگین تفاوت مربعی بین مقادیر پیش‌بینی‌شده و واقعی است.

بهینه‌ساز: بهینه‌ساز Adam به دلیل کارایی آن در مدیریت گرادیان‌های پراکنده و توانایی‌های نرخ یادگیری تطبیق‌پذیر انتخاب شده است، که می‌تواند منجر به همگرایی سریع‌تر شود.

Callbacks: EarlyStopping, ModelCheckpoint, و ReduceLROnPlateau به عنوان callbacks استفاده می‌شوند تا از overfitting جلوگیری کنند، بهترین مدل را ذخیره کنند، و نرخ یادگیری را کاهش دهند اگر زیان اعتبارسنجی بر روی یک سطح ثابت باقی بماند، به ترتیب.

دلایل این انتخاب‌ها بر اساس ماهیت داده‌ها، پیچیدگی مسئله و نیاز به مدلی است که بتواند به خوبی تعمیم داده شود، بدون اینکه به داده‌های آموزشی بیش‌برازش کند. استفاده از PReLU نشان‌دهنده تمایل به ضبط روابط غیرخطی بدون مشکل ReLU مرده است، و استفاده از Dropout نشان‌دهنده تلاش برای افزایش توانایی تعمیم مدل است. بهینه‌ساز Adam یک انتخاب منطقی پیش‌فرض برای بسیاری از کاربردهای شبکه عصبی به دلیل نرخ یادگیری تطبیق‌پذیر و عملکرد کلی خوب در طیف وسیعی از وظایف است.

سایر مدل ها

در اینجا ما دو مدل SVM و درخت تصمیم را انتخاب کردیم. به این دلیل svm را انتخاب کردیم که در فضاهای غیر خطی و فضاهای دارای نویز پر قدرت عمل میکند. انتخاب درخت تصمیم هم به این علت بود که در واقع چون داده‌ها به دانشکده‌های مختلف، غذا‌های مختلف و... تقسیم می‌شد احتمالاً درخت تصمیم به خوبی می‌توانست آن‌ها را تقسیم کند.

در این دو متد کاری که انجام دادیم این بود که هایپر پارامتر ها را روی داده validation فیت کردیم و سپس دقت مدل را روی داده تست بدست آوردیم.

مقایسه

دقت مدل ها روی داده تست:

	R2_SCORE	R2_SCORE ADJ	RMSE
NN	0.810547	-	-
SVR	0.777646	0.765432	8.064536
DECISION TREE	0.731059	0.716287	8.869207

روش شبکه های عصبی روی این داده بهتر جواب میدهد. به چند دلیل:

1. درک روابط پیچیده بین فیچر ها و تولید فیچر های پیچیده تر . در صورتی که مدل هایی مثل درخت تصمیم چون مدل های white box ای هستند توان تشخیص فیچر های پیچیده را ندارند.
2. همچنین در روش های قسمت دوم ممکن است یک مساله روی یک متد به خوبی عمل نکند یا واریانس زیادی داشته باشد برای همین از روش های ensemble استفاده میشود که چند مدل مانند knn , decision tree , svm روی یک داده فیت میشوند و با استفاده از روش رای گیری پیشبینی را انجام میدهند که این باعث کاهش واریانس و افزایش قدرت میشود.
3. پارامتر های عصبی توسط خود مدل با استفاده از روش های optimization پیدا میشود اما پارامتر های درخت تصمیم و SVM را باید با آزمون و خطا پیدا کرد.

Feature Analysis

در داده های ما، بعد از استفاده از PCA، مقدار r^2 squared ما کاهش می یابد و از مقدار 0.81 به مقدار 0.77 میرسد.

کاهش عملکرد: کاهش امتیاز R-squared نشان دهندهی آن است که توانایی مدل در توضیح تغییرات متغیر هدف پس از کاهش بعد با PCA کاهش یافته است.

تأثیر PCA: این ممکن است نشان دهندهی آن باشد که برخی اطلاعات پیش بینی کنندهی مهم در فرآیند PCA از دست رفته اند، که می تواند اتفاق بیفتد اگر اجزای حذف شده شامل ویژگی های مرتبط باشند.

تغییرات داده: در نظر گرفتن میزان تغییرات دادهی اولیه که پس از PCA حفظ شده است، حیاتی است. اگر هدف حفظ درصد بالایی از تغییرات (مثلاً 95٪) بوده و ابعاد نتیجه گیری شده آن را در بر نگرفته اند، ممکن است منجر به عملکرد ضعیف تر مدل شود.

اهمیت features: اینجا به این نکته پی بردیم که بعد از استفاده از PCA ، دقت ما کاهش یافت و این به این دلیل است که بعد از کاهش ابعاد، برخی ویژگی ها از داده های ما حذف شده است که این واژه ها در مدل ما تاثیر گذار بوده اند و فقدان آنها باعث ضعیف تر شدن مدل میشود.

چالش ها

داده ها وابستگی زمانی داشتند و باید سعی می کردیم این موضوع را روی داده بازنمایی کنیم. در دریافت داده ها از مراجع ذی ربط به مشکل خوردیم. مسئولین مربوطه در اول کار با بی اهمیتی و بی توجهی داده های کوچکی را در اختیار ما قرار دادند که حاوی اطلاعات کمی بود اما با پافشاری توانستیم داده های لازم را دریافت کنیم. یکی از فیچر های موثر فاصله تا تعطیلی ها بود که برای این کار مجبور شدیم تقویم شمسی را crawl کنیم و کلاس بزرگی برای محاسبه تعطیلی ها درست کنیم. داده های دریافتی به صورت log رزرو غذا بود که باید بر حسب مواردی که می خواستیم فیلتر و جمع می کردیم.