

به نام خدا



آزمایشگاه R درس تحلیل رگرسیون

دانشگاه صنعتی شریف

دانشکده ریاضی و علوم کامپیوتر

مدرس : دکتر میرصادقی - دکتر کدیور

تمرین سری اول

علی نوریان

۹۸۱۰۲۵۲۷

نیمسال ۰۲-۲

۱ بخش اول - پیش پردازش و آماده سازی دیتاست

در ابتدا کتابخانه‌های مورد نیاز و سپس دیتاست را لود می‌کنیم:

```
1 # ===== Import Libraries & Data =====
2 Sys.setlocale(locale = 'persian')
3
4 library(data.table)
5 library(ggplot2)
6
7 d = fread('iranprovs_mortality_monthly.csv', encoding = 'UTF-8')
```

سپس متغیر ym_num را توسط مقادیر سال و ماه به صورتی می‌سازیم که برای هر ماه هر سال یک مقدار یکتا داشته باشد. از این متغیر برای مرتب سازی و رسم نمودار استفاده می‌کنیم. همچنین دیتاست جدیدی را با حذف گروه سنی و جنسیت می‌سازیم. محتوای دیتاست جدید را آمار ۶ سال قبل تا به امروز قرار می‌دهیم زیرا داده‌های قبل ۶ سال با توجه به تغییر شرایط کشور و وضعیت بهداشت و جمعیتی ممکن است مناسب نباشند. دو متغیر *provs* و *months* را نیز برای داشتن مقادیر یکتای استان‌ها و ماه‌ها تعریف می‌کنیم:

```
1 # ===== pre-processing =====
2 d$ym_num = d$y + d$m / 12 - 1/24 # create year-month feature
3
4
5 ym_num_covid = 1398 + 10/12 - 1/24 # set start date of COVID-19
6 ym_num_start = ym_num_covid - 6 # set start date of task valid
7 data
8
9 ds = d[, .(n = sum(n)), .(y, m, ym_num, prov)] # remove age & sex features
10 ds = ds[ym_num > ym_num_start] # remove invalid data for this task
11
12 provs = unique(ds$prov) # provinces
13 months = unique(ds$m) # months
```

۲ بخش دوم - مدل خطی برای پیش‌بینی فوت عادی

برای بدست آوردن مقدار فوت اضافه ابتدا نیاز است مقدار فوت در شرایط عادی (بدون حضور بیماری) را پیش‌بینی کنیم. برای اینکار یک مدل خطی برای پیش‌بینی میزان فوت بدون حضور کرونا در بازه زمانی شیوع کرونا در نظر می‌گیریم. مدل خطی را با توجه به داده‌های قبل کرونا بدست آورده و برای بازه کرونا پیش‌بینی را انجام می‌دهیم.

با توجه به اینکه میزان فوتی در ماه‌های مختلف سال متغیر و وابسته به ماه است، مدل خطی را به تفکیک هر ماه برای هر استان بدست می‌آوریم. همچنین با توجه به مقدار p-value اگر مدل خطی مناسب نباشد از میانگین فوتی‌های قبل کرونا برای پیش‌بینی استفاده می‌کنیم.

برای هر مدل یک بازه اطمینان برای متغیر پاسخ (معادل با دو برابر انحراف معیار داده‌ها) در نظر می‌گیریم که اگر مقدار فوتی در زمان کرونا خارج از این بازه قرار گرفت به معنی تاثیر کرونا بر فوتی و ایجاد فوت اضافه است.

کد این قسمت در صفحه بعد آمده است.

```

1 # ===== fit linear model and obtaining the number of excess death =====
2
3 ds$n_predicted = 0                                # model predicted value for death
4 ds$upper_thresh = 0                               # upper threshold for normal death
5 ds$lower_thresh = 0                               # lower threshold for normal death
6
7 # fit model for each month of each province
8 for (this_prov in provs) {
9   for (this_month in months) {
10     condition = ds$prov == this_prov & ds$m == this_month
11
12     this_ds = ds[condition]
13     this_train_ds = this_ds[ym_num < ym_num_covid]
14
15     fit = lm(n ~ ym_num, this_train_ds)
16     pvalue = summary(fit)$coefficients[2,4] # calculating p-value
17
18     # choose model prediction for p-value < 0.1
19     # otherwise, choose average of samples instead of model prediction
20     if (pvalue < 0.1) {
21       n_predicted = pmax(predict(fit, this_ds), 0)
22       # if prediction gets negative then use min num of death instead
23       n_predicted[n_predicted == 0] = min(this_train_ds$n)
24     } else {
25       n_predicted = rep(mean(this_train_ds$n), length(this_ds$n))
26     }
27
28     sigma = sd(this_train_ds$n) # calculating standard-deviation
29
30     ds[condition]$n_predicted = n_predicted # save
31     predictions
32     # set upper threshold
33     ds[condition]$upper_thresh = n_predicted + 2 * sigma
34     # set lower threshold
35     ds[condition]$lower_thresh = n_predicted - 2 * sigma
36   }
37 }

```

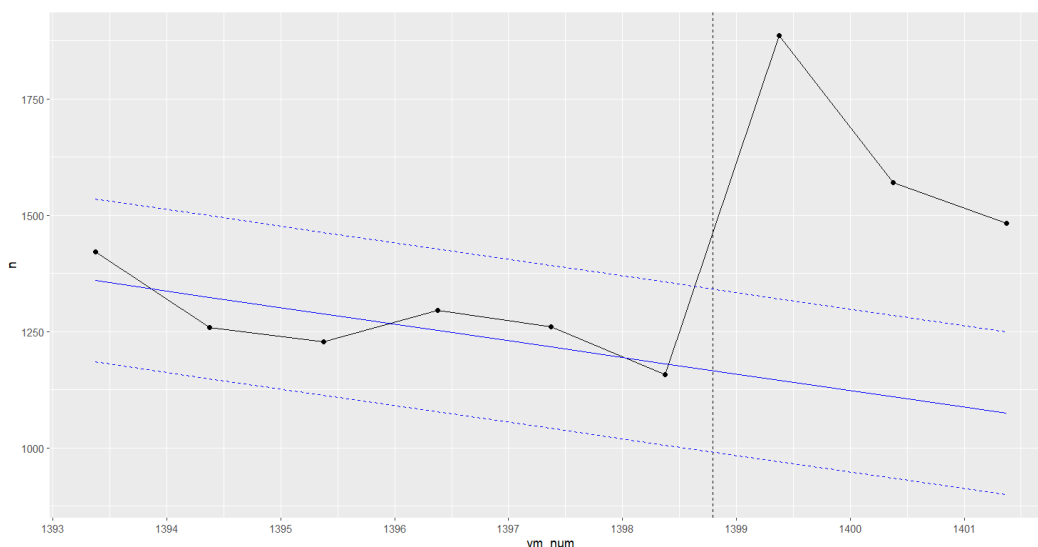
سپس با توجه به پیش‌بینی مقدار فوت در غیاب کرونا و همچنین آستانه بالا برای در نظر گرفتن تغییرات عادی در فوتی‌ها بدین صورت عمل می‌کنیم: اگر مقدار فوتی از آستانه بالا بیشتر شود، آن مقدار فوتی نامعقول بوده و اختلاف آن با مقدار فوت پیش‌بینی شده، مقدار فوت اضافه خواهد بود.

```

1 # calculating excess death according to predictions and upper thresholds
2   of death
3 ds$excess_death = round((ds$n - ds$n_predicted)*(ds$n > ds$upper_thresh))

```

نمونه‌ای از این اعمال مدل خطی در شکل زیر آماده است:



شکل ۱: فوتی‌های استان آذربایجان غربی برای ماه مرداد از سال ۱۳۹۳ تا ۱۴۰۱

داده‌های مشکل مقدار فوت ثبت شده در مرکز آمار، خط آبی، مدل خطی بدست آمده توسط داده‌های قبل کرونا، خط چین مشکی مشخص کننده زمان شروع کرونا و دو خط چین آبی نشان‌دهنده آستانه بالا و پایین برای میزان فوتی می‌باشد.

همانطور که در نمودار بالا نیز مشخص است بعد از شروع کرونا، میزان فوتی در ماه مرداد سال‌های ۱۳۹۹ تا ۱۴۰۱ تفاوت فاحشی با سال‌های قبل از کرونا دارد.

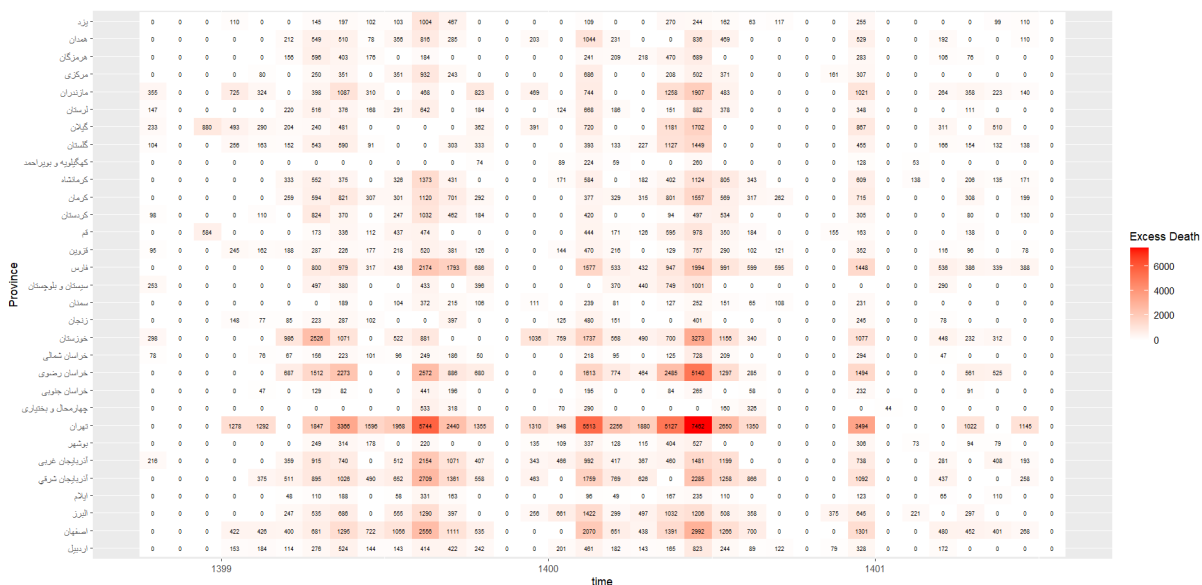
۳ بخش سوم - رسم نمودارهای فوت اضافه استان‌ها

در قسمت قبل مقدار فوت اضافه برای هر استان به تفکیک ماه‌های سال را بدست آوردیم. در این قسمت جدول مربوط به آمار فوتی اضافه در هر ماه از استان‌ها و همچنین نقشه حرارتی را برای نشان دادن شدت و ضعف کرونا در ماه‌های مختلف هر استان رسم می‌کنیم:

```
1 # ===== COVID-19 excess death =====
2
3 COVID19_ds = ds[ym_num >= ym_num_covid] # COVID-19 data
4 COVID19_ds$norm_excess_death = COVID19_ds$excess_death / COVID19_ds$n_
  predicted * 100
5
6 # plot heat map: percentage of excess death for each month of each
  province
7 ggplot(COVID19_ds, aes(x = ym_num, y = prov)) +
8   geom_raster(aes(fill = norm_excess_death)) +
9   scale_fill_gradient(low = "white", high = "red")+
10  labs(title = "COVID-19 Percentage of Excess Death", x = "time", y = "
    Province")
```

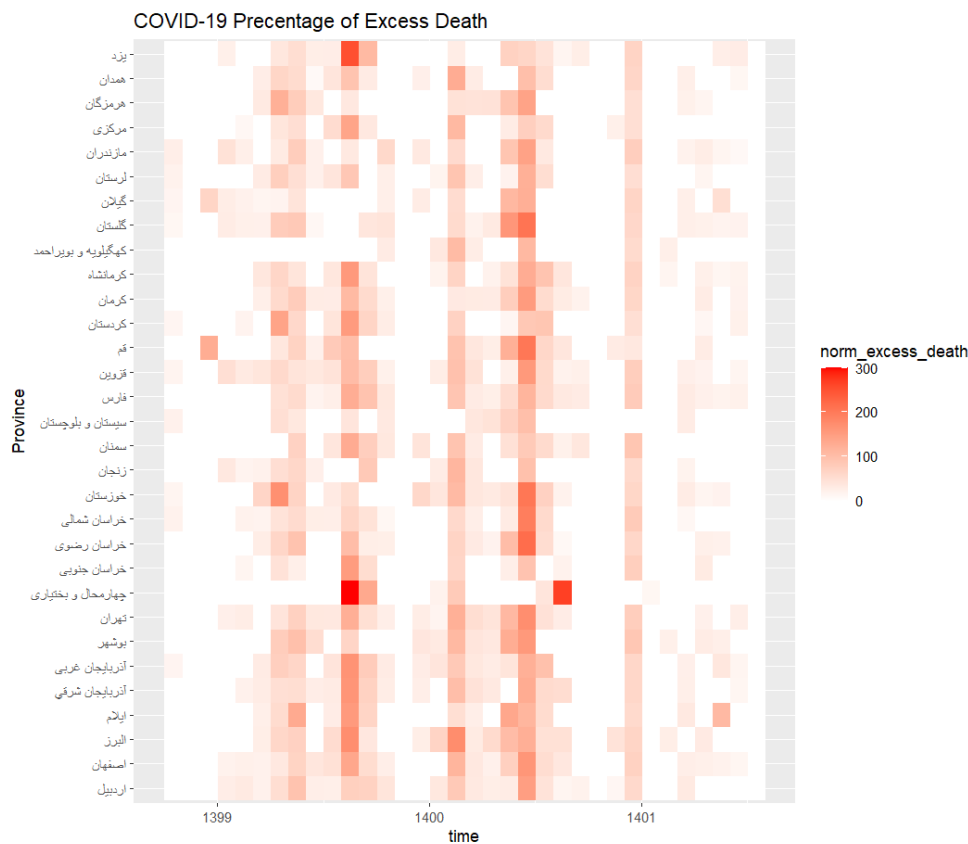
پس از بدست آوردن مقدار فوت اضافه در هر ماه از هر استان، با توجه به اینکه جمعیت استان‌ها با یکدیگر متفاوت است، در نقشه حرارتی برای نشان دادن شدت و ضعف کرونا از نسبت میزان فوت اضافه به فوت معقول (مقدار پیش‌بینی شده) استفاده می‌کنیم. بنابراین مقادیر موجود در نقشه حرارتی درصد تغییر فوتی‌ها را در اثر کرونا نشان می‌دهد. البته در جدول صفحه بعد مقدار فوت اضافه در هر خانه جدول نمایش داده شده است. خروجی کد فوق در ادامه آمده است.

جدول فوتی‌های اضافه:



شکل ۲: مقدار فوت اضافه در ماه‌های مختلف در استان‌های کشور

نقشه حرارتی شدت و ضعف کرونا:



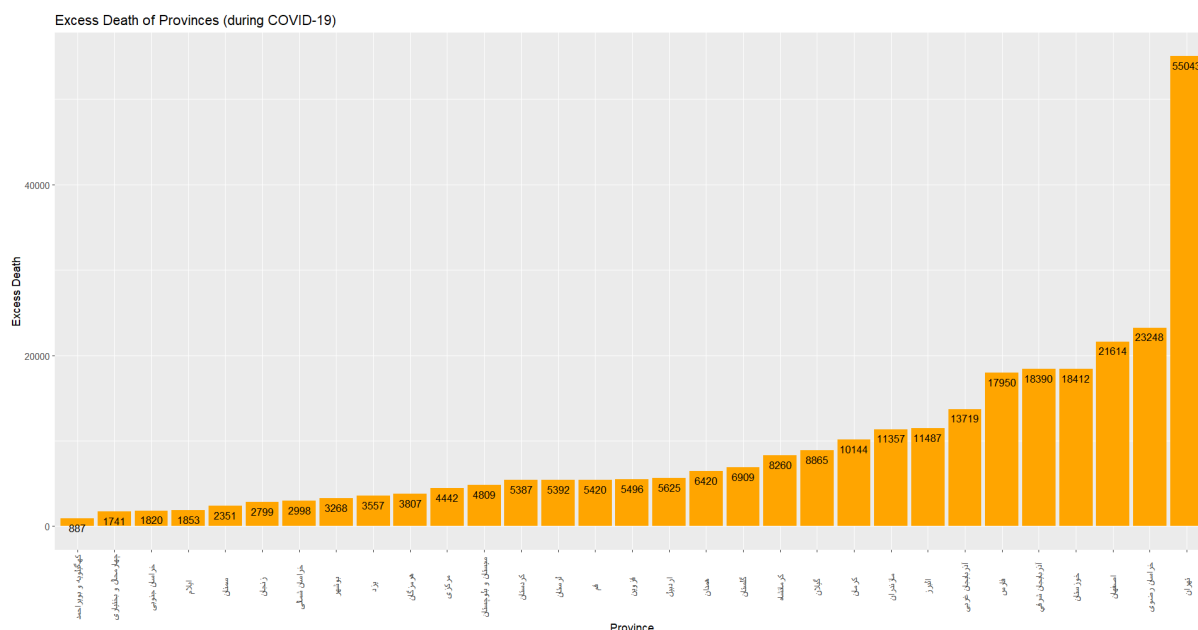
شکل ۳: نقشه حرارتی شدت و ضعف کرونا در استان‌های کشور

۴ بخش چهارم - مجموع تعداد فوت اضافه در استان‌ها و کشور

در قسمت قبل مقدار فوت اضافه را برای هر ماه از هر استان حساب کردیم. در این قسمت نمودار مجموع فوتی‌های اضافه استان‌ها را نشان می‌دهیم:

```
1 # calculate the total excess death for each province
2 prov_excess_death = COVID19_ds[, .(excess_death=sum(excess_death)), .(y,
   prov)]
3
4 # plot the total excess death of provinces
5 ggplot(prov_excess_death, aes(x = reorder(prov, excess_death), y =
   excess_death, fill=factor(y))) +
6 geom_bar(stat="identity")+
7 labs(title = "Excess Death of Provinces (during COVID-19)", x = "
   Province", y = "Excess Death")+
8 scale_fill_manual(values=rainbow(20))+
9 theme(axis.text.x = element_text(angle = 90, vjust = 0.5))
10
11 all_excess_death = sum(prov_excess_death$excess_death)
12 sistan_excess_death =sum(prov_excess_death[prov==provs[16]]$excess_death)
13 yazd_excess_death = sum(prov_excess_death[prov==provs[31]]$excess_death)
```

نتایج کد فوق به صورت زیر است:



شکل ۴: فوتی‌های اضافه هر استان

تعداد فوتی‌های اضافه کل کشور و تعداد فوتی‌های استان سیستان و بلوچستان و استان یزد (خواسته سوال) به صورت زیر است:

تعداد فوتی‌های اضافه کل کشور: ۲۹۳۴۷۰

تعداد فوتی‌های اضافه استان سیستان و بلوچستان: ۴۸۰۹

تعداد فوتی‌های اضافه استان یزد: ۳۵۵۷

۵ بخش پنجم - مقایسه عملکرد استان‌ها برای کنترل کرونا

اکنون می‌خواهیم عملکرد استان‌های مختلف در برابر کنترل و مهار کرونا را بررسی کنیم. ایده اصلی بدست آوردن سیر تغییر فوتی‌ها در هر استان است. بدین صورت که هر استانی کاهش فوت بیشتری را تجربه کرده باشد یا به عبارت دیگر شیب فوتی‌های آن استان کمتر از سایر استان‌ها باشد، در نتیجه عملکرد بهتری در مهار داشته است.

با توجه به آنکه میزان فوتی در سنین مختلف متفاوت است یا به بیان دیگر تاثیر کرونا بر سنین مختلف متفاوت بوده است، ابتدا جمعیت را به سه بازه سنی ۰ تا ۲۵ سال، ۲۵ تا ۶۵ سال، بالای ۶۵ سال تقسیم می‌کنیم و ابتدا عملکرد استان‌ها را در سه بازه سنی معرفی شده به تفکیک بررسی می‌کنیم.

البته توجه به نکته خوب است که بازه سنی زیر ۱۰ سال در برابر این ویروس مقاومت بیشتری داشته‌اند و درگیری این سنین با کرونا بسیار کم گزارش شده است. همچنین در سنین بسیار پایین دلایل فوت دلایلی غیر از کرونا دارد. بنابراین بازه سنی زیر ۱۰ سال را نیز از دیتا حذف کرده‌ایم.

با توجه به اینکه توزیع جمعیتی در استان‌های مختلف متفاوت است، داده‌های هر استان را نرمالایز می‌کنیم تا شیب‌های بدست آمده برای استان‌های مختلف قابل مقایسه با یکدیگر باشند. برای اینکار داده‌های قبل کرونا را جدا کرده و مشابه همان کاری که در بخش ۲ انجام دادیم، میزان فوت عادی در ماه‌های مختلف هر استان را بدست می‌آوریم. سپس داده‌های فوتی بعد از کرونا را با مقدار فوتی عادی منتظر با آن نرمالایز می‌کنیم.

موارد فوق را برای هر بازه سنی به صورت جداگانه اعمال شده است. سپس بر روی داده‌های بدست آمده که نرمال شده‌ی آمار فوتی بعد از کرونا هستند، مدل خطی اعمال می‌کنیم و شیب آن را برای آن استان ذخیره می‌کنیم. بنابراین در نهایت ۳ شیب برای هر استان که هر کدام متعلق به یک گروه سنی هستند بدست می‌آید.

برای نمونه کد مربوط به عملیات برای یک گروه سنی در زیر آماده است:

```
1 # ===== Best Provinces in Controlling COVID-19 =====
2 new_ds = d[, .(n = sum(n)), .(ym_num, y, m, prov, age_group)]
3
4 # "bc" stands for "before COVID19"
5 before_COVID19_ds = new_ds[ym_num >= ym_num_start & ym_num < ym_num_covid
6   ]
7 bc_young_ds = before_COVID19_ds[age_group >= "10-14" & age_group <= "
8   20-24"]
9 bc_young_ds = bc_young_ds[, .(n=mean(n)), .(m, prov)]
10 after_COVID19_ds = new_ds[ym_num >= ym_num_covid]
11
12 young_ds = after_COVID19_ds[age_group>="10-14" & age_group<="20-24"]
13 young_ds = young_ds[, .(n=sum(n)), .(ym_num, y, m, prov)]
14
15 young_ds$n2 = 0
16
17 # normalizing the number of deaths for each months of each province
18 for (this_prov in provs) {
19   for (this_month in months) {
20     young_ds[m == this_month & prov == this_prov]$n2 =
21     young_ds[m == this_month & prov == this_prov]$n / bc_young_ds[m ==
22       this_month & prov == this_prov]$n
23   }
24 }
```

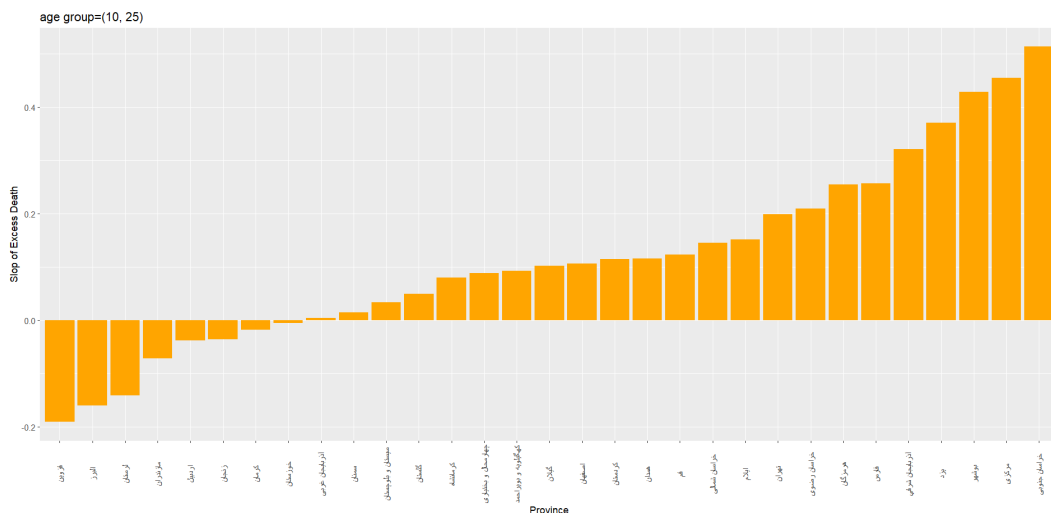
ادامه کد صفحه قبل:

```

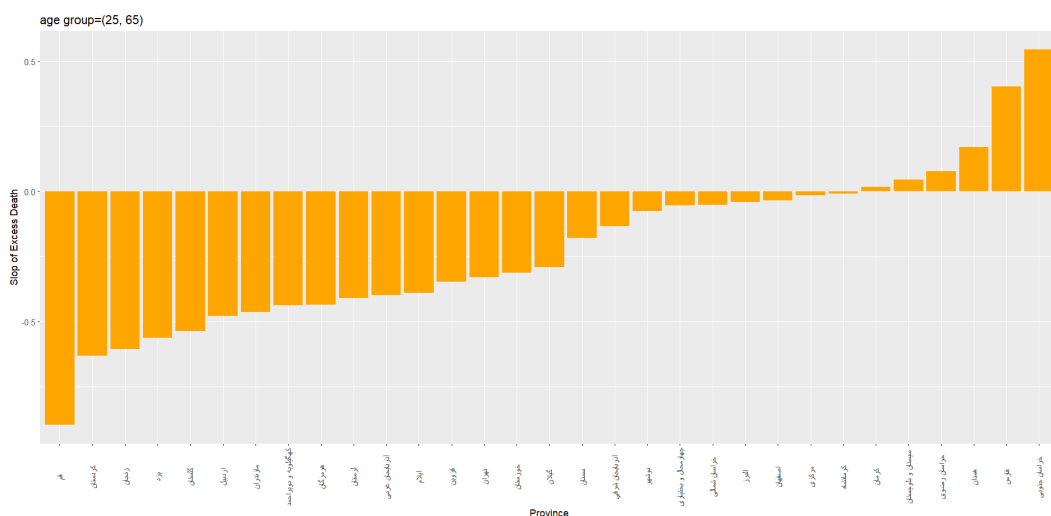
1 young_ds$n_preidicted = 0
2 slop_result = new_ds[, .(n=sum(n)), .(prov)]
3 slop_result$young_slop = 0
4 for (this_prov in provs) {
5   condition = young_ds$prov == this_prov
6   this_young_ds = young_ds[condition]
7   fit = lm(n2 ~ ym num, this_young_ds)
8   slop_result[prov==this_prov]$young_slop=summary(fit)$coefficients[2,1]
9   young_ds[condition]$n_preidicted = predict(fit, this_young_ds)
10 }

```

مشابه آنچه برای گروه سنی جوان آمده است و در کنار همین عملیات، عملیات مربوط دو رده سنی بزرگسال و پیر نیز انجام می‌شود.
خروجی مقایسه کاهش فوتی‌ها برای سه گروه سنی جوان، بزرگسال و پیر به تفکیک در ادامه آمده است:

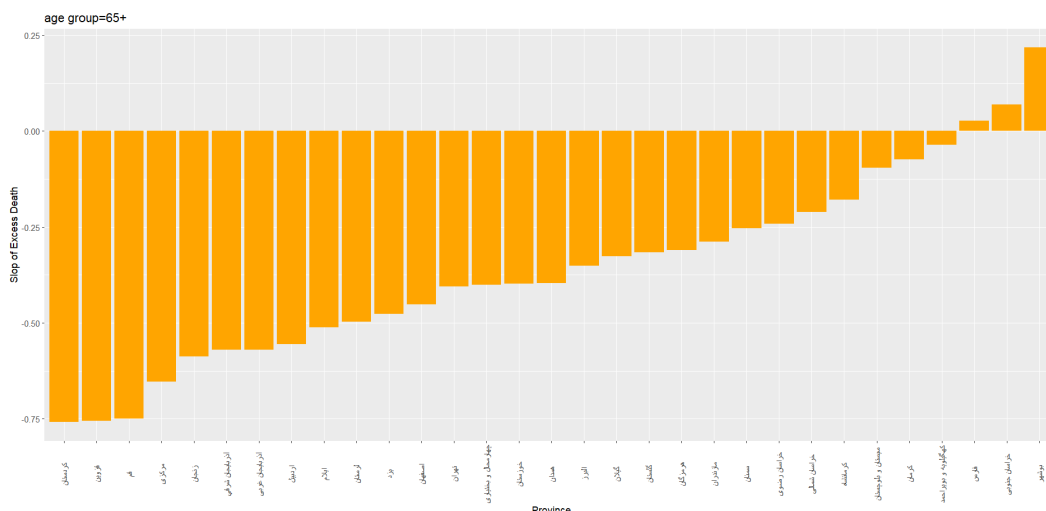


شکل ۵: شیب تغییر فوتی‌ها در استان‌های مختلف در گروه سنی جوان



شکل ۶: شیب تغییر فوتی‌ها در استان‌های مختلف در گروه سنی بزرگسال

ذکر این نکته خوب است که استان قم عملکرد خوب قابل توجهی نسبت به سایر استان‌ها در بازه سنی بزرگسال دارد.

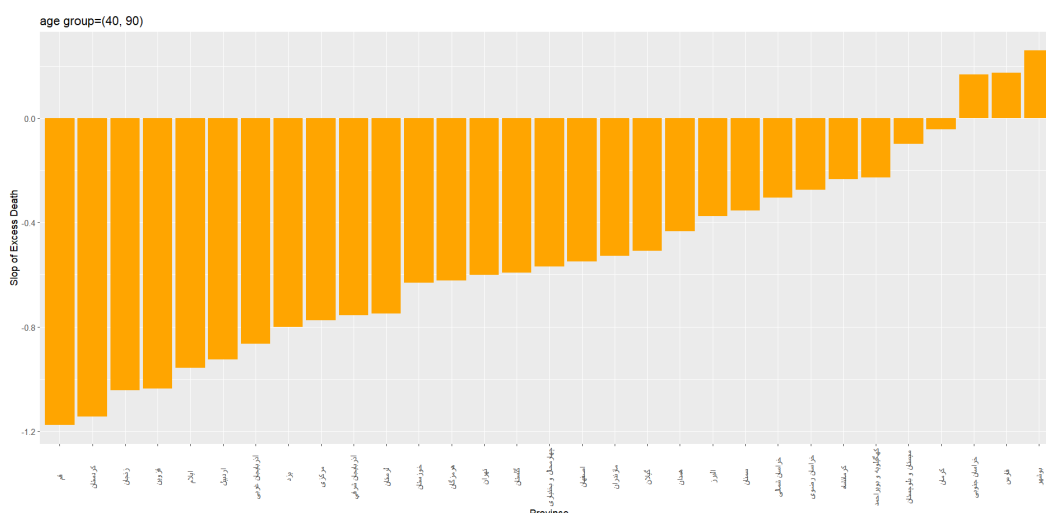


شکل ۷: شیب تغییر فوتی‌ها در استان‌های مختلف در گروه سنی پیر

همانطور که نمودارها مشخص است بهترین استان‌ها در کنترل کرونا در بازه‌های سنی مشخص شده به صورت زیر است:
 بازه سنی جوان (۱۰ تا ۲۵ سال): قزوین - البرز - لرستان - مازندران - اردبیل
 بازه سنی بزرگسال (۲۵ تا ۶۵ سال): قم - کردستان - زنجان - یزد - گلستان
 بازه سنی پیر (بالای ۶۵ سال): کردستان - قزوین - قم - مرکزی - زنجان
 عملکرد سایر استان‌ها نیز در نمودارها مشخص است.

اکنون با توجه به اینکه سنین بالاتر بیشتر در معرض این ویروس بودند بنابراین عملکرد هر استان در سنین بالاتر بیشتر از سنین پایین شایان توجه است. همچنین افراد با سن بسیار بالا (بالای ۹۰ سال) بسیار آسیب‌پذیر هستند و کادر درمان نمی‌تواند تاثیر جدی‌ای برای درمان این افراد داشته باشد. از این رو استان‌هایی که میانگین سنی بالاتری دارند فوتی بیشتری را تجربه می‌کنند. بنابراین بازه سنی ۴۰ تا ۹۰ سال را برای بررسی عملکرد کلی یک استان در نظر می‌گیریم.

کد این قسمت کاملاً مانند کد ذکر شده در ابتدای این قسمت است. نتیجه عملکرد استان‌ها به صورت کلی در شکل زیر قابل مشاهده است:



شکل ۸: شیب تغییر فوتی‌ها در استان‌های مختلف

میتوان ترکیبی از سه شیب بدست آمده از گروه سنی‌های مختلف را ارائه داد. نتایج این کار نیز با تقریب خوبی مشابه نمودار فوق است. در نتیجه ۵ استان دارای بهترین عملکرد در کنترل کرونا قم، کردستان، زنجان، قزوین و ایلام می‌باشند.