

به نام خدا

پروژه داده کاوی

استاد راهنما: دکتر فقیهی حبیب آبادی

عنوان پروژه: طرح هزینه و درآمد خانوار سال ۱۳۹۹

فهرست

پیشگفتار

طرح آمارگیری هزینه و درآمد خانوار از سال ۱۳۴۲ در مناطق روستایی و از سال ۱۳۴۷ در مناطق شهری به اجرا درآمده است. از سال ۱۳۵۳ علاوه بر هزینه‌های خانوار، اطلاعات درآمد نیز گردآوری شد. هدف کلی طرح آمارگیری هزینه و درآمد خانوار، برآورد متوسط هزینه‌ها و درآمد یک خانوار شهری و روستایی در سطح کشور و استان است. اهمیت این طرح در امکان بررسی روند مصرف کالاها و خدمات، مطالعه روابط متقابل ویژگی‌های اجتماعی اقتصادی خانوارها، ارزیابی آثار سیاست‌های اقتصادی در زمینه تأمین عدالت اجتماعی و بررسی توزیع درآمد و سرانجام نقش مهم آن در تأمین اطلاعات مورد نیاز حساب‌های ملی و منطقه‌ای، به خصوص در بررسی‌ها و برنامه‌ریزی‌های اقتصادی اجتماعی کشور است.

طرح آمارگیری هزینه‌های و درآمد خانوار سال ۱۳۹۹ بر اساس آمارگیری از ۱۹۳۰۶ خانوار نمونه در نقاط شهری و ۱۸۲۵۱ خانوار نمونه در نقاط روستایی کشور اجرا شده است.

نتایج تفصیلی این طرح شامل ۲۳۴ جدول در قالب ۲ نشریه به تفکیک شهری و روستایی توسط مرکز آمار ایران هر ساله منتشر می‌شود. همچنین نتایج تفصیلی طرح در قالب ۵۱۲ جدول، شامل جزئی‌ترین اطلاعات در سطح ۳۱ استان کشور و کل کشور برای استفاده از طریق واحد اطلاع‌رسانی مرکز آمار ایران در اختیار عموم علاقمندان، محققین، برنامه‌ریزان و پژوهشگران قرار می‌گیرد.

جامعه هدف این طرح شامل همه‌ی خانوارهای معمولی ساکن و گروهی در نقاط شهری یا روستایی است. این خانوارهای نمونه از ۳۸۷ شهرستان در مناطق شهری و از ۳۹۵ شهرستان در مناطق روستایی کل کشور انتخاب شده‌اند.^۱

فصل اول: معرفی مجموعه داده‌ها

مقدمه

به دلیل وسعت کشور ایران و وجود اقلیم‌های متنوع مطمئناً میزان هزینه و درآمد خانوار در همه استان‌ها و خانوارهای مختلف به یک میزان نمی‌باشد. پس برای اجرای عدالت اجتماعی اقتصادی در بین اقشار مختلف جامعه همواره باید الگویی خاص برقرار شود تا خانوارهایی که مستحق دریافت یارانه کمک معیشتی هستند از خانوارهایی که

پیشگفتار ۱

فصل اول: معرفی مجموعه داده‌ها ۱

مقدمه ۱

هدف ۲

خلاصه هدف: ۲

ضرورت هدف: ۲

معرفی مجموعه داده‌ها ۲

معرفی متغیرهای مجموعه داده‌ها ۳

نتیجه‌گیری فصل اول: ۴

فصل دوم: تصویری سازی ۵

توضیحات نمودارها به صورت کلی ۱۵

فصل سوم: درخت ۱۷

درخت رده‌بندی ۱۷

جنگل تصادفی ۱۹

فصل چهارم: رگرسیون لوژستیک ۲۰

ارزیابی عملکرد رده‌بندی ۲۱

فصل پنجم: k-نزدیک ترین همسایه (k-NN) ۲۱

فصل ششم: شبکه عصبی ۲۲

نتیجه گیری کلی ۲۳

سخن پایانی ۲۳

^۱ - بر اساس آخرین تقسیمات کشوری تا پایان سال ۱۳۹۵، کشور ایران از ۳۱ استان و ۴۲۹ شهرستان تشکیل شده است.

۱۰۰۱۰۰۰۹۳۱۸	مرکزی	۱
۱۰۱۰۶۰۱۸۴۲۶	گیلان	۲
۱۰۲۱۰۰۳۱۴۱۳	مازندران	۳
۱۰۳۱۲۰۴۴۸۱۷	آذربایجان شرقی	۴
۱۰۴۰۱۰۴۸۸۲۹	آذربایجان غربی	۵
۱۰۵۰۵۰۶۴۱۱۸	کرمانشاه	۶
۱۰۶۰۳۰۷۳۹۲۳	خوزستان	۷
۱۰۷۲۴۰۹۴۹۲۶	فارس	۸
۱۰۸۰۶۰۹۹۴۲۶	کرمان	۹
۱۰۹۱۶۱۱۱۵۱۴	خراسان رضوی	۱۰
۱۱۰۱۲۱۳۰۸۲۹	اصفهان	۱۱
۱۱۱۰۶۱۴۳۴۱۵	سیستان و بلوچستان	۱۲
۱۱۲۰۴۱۵۳۱۲۹	کردستان	۱۳
۱۱۳۰۳۱۶۲۴۲۹	همدان	۱۴
۱۱۴۰۲۱۷۵۴۱۷	چهارمحال و بختیاری	۱۵
۱۱۵۰۹۱۹۱۵۲۳	لرستان	۱۶
۱۱۶۰۱۱۹۲۵۲۳	ایلام	۱۷
۱۱۷۰۶۲۰۸۸۱۷	کهگیلویه و بویراحمد	۱۸
۱۱۸۰۳۲۱۴۰۲۶	بوشهر	۱۹
۱۱۹۰۴۲۳۳۸۲۱	زنجان	۲۰
۱۲۰۰۲۲۳۱۸۲۲	سمنان	۲۱
۱۲۱۰۵۲۴۴۴۱۵	یزد	۲۲
۱۲۲۰۵۲۵۸۹۲۹	هرمزگان	۲۳
۱۲۳۰۱۲۷۷۹۱۹	تهران	۲۴
۱۲۴۰۱۲۹۶۷۲۶	اردبیل	۲۵
۱۲۵۰۱۳۰۴۴۱۹	قم	۲۶
۱۲۶۰۳۳۱۹۷۱۷	قزوین	۲۷
۱۲۷۱۴۳۴۴۸۲۰	گلستان	۲۸
۱۲۸۰۴۳۵۳۲۱۳	خراسان شمالی	۲۹
۱۲۹۰۱۳۵۸۱۲۱	خراسان جنوبی	۳۰
۱۳۰۰۱۳۷۶۱۲۳	البرز	۳۱

جدول ۱: نمونه شناسه‌های موجود در دیتا خام

این کدگذاری استان‌ها و شهرستان‌ها در درگاه ملی آمار ایران صورت گرفته است.^۳

نیازی به یارانه مذکور ندارند، تفکیک شوند، تا علاوه بر کاهش هزینه‌های اقتصادی کشور به بهبود برنامه‌ریزی اجتماعی اقتصادی منجر شود.

به گزارش ایسنا: "در ۱۹ اردیبهشت ماه سال ۱۴۰۱ وزیر امور اقتصادی و دارایی کشور از آغاز طرح واریز یارانه کمک معیشتی به میزان ۳۰۰ و ۴۰۰ هزار تومان برای تعداد اعضای خانوار خبر دادند."^۱

این مبالغ به حساب ۷۷ میلیون نفر از جمعیت کشور که تا پیش از این یارانه ماهیانه دریافت می‌کردند، به شرح زیر واریز شد:

سه دهک پایین جامعه که از درآمد کمی برخوردار بودند، هر یک ۴۰۰ هزار تومان، شش دهک میانی جامعه که از درآمد متوسط برخوردار بودند، هر یک ۳۰۰ هزار تومان دریافت نمودند و دهک دهم جامعه یارانه مذکور را دریافت ننمودند.

پس از سیل عظیمی از اعتراضات منوط به عدم قرارگیری صحیح در دهک بندی صورت گرفته، تصمیم بر آن شد که این حمایت مالی از دهک بالای جامعه نیز انجام گیرد.

هدف

حال هدف این پروژه بر آن است تا مدلی ارائه نماید که این دهک بندی ها را به نحوی انجام دهد تا افراد قرار گرفته در دهک دهم تا جایی که ممکن است با خطای کمتری شناسایی گردند.

خلاصه هدف: شناسایی صحیح خانوارهای قرار گرفته در دهک دهم.

ضرورت هدف: اجرای صحیح دهک بندی و کاهش تعداد خانوارهای معترض در طرح یارانه کمک معیشتی و شناسایی خانوارهایی که نیازی به یارانه مذکور ندارند.

معرفی مجموعه داده‌ها

مجموعه داده‌های هزینه و درآمد خانوار سال ۱۳۹۹ مرکز آمار ایران^۲ در قالب یک فایل (Access) و توضیحات مرتبط با این فایل به همراه پرسشنامه‌ای است که این داده‌ها از آن استخراج گردیده است، می‌باشد.

در این مجموعه داده‌ها اسامی افراد خانوار و آدرس محل سکونت به دلیل رعایت حریم شخصی و حفظ حقوق شهروندان با شناسه‌هایی کدگذاری شده‌اند. تنها معیار برای افراز این خانوارها سه رقم اول این کدگذاری می‌باشد، که نشان دهنده استان محل سکونت خانوار است.^۳ (جدول شماره ۱)

ردیف	نام استان	نمونه شناسه کدگذاری شده
------	-----------	-------------------------

^۱ - برای جزئیات خبر می‌توانید از طریق پیوند زیر اطلاعات بیشتری کسب نمایید.

^۲ - پیوند مجموعه داده‌ها: (<https://drive.google.com/file/d/1QEcXCzVIWOjp1IRGqNnrkpryrsNBxppa/view?usp=sharing>)

^۳ - فایل تقسیمات کشوری سال ۹۹ را می‌توانید از طریق پیوند زیر مشاهده نمایید. (<https://www.amar.org.ir/Portals/0/Geo/GEO99.xlsx>)

همانطور که در قسمت مقدمه ذکر شد و از آنجایی که میزان هزینه و درآمد خانوار در استان‌ها متفاوت است پس استان‌های کشور را در ۷ گروه‌بندی کلی به شرح زیر در نظر می‌گیریم. (جدول شماره ۲)

گروه	استان‌های قرار گرفته در گروه
۱	کرمان، خراسان رضوی، خراسان شمالی، خراسان جنوبی، سیستان و بلوچستان
۲	گیلان، مازندران، گلستان
۳	تهران، قزوین، البرز
۴	مرکزی، اصفهان، همدان، سمنان، قم، یزد
۵	آذربایجان شرقی، آذربایجان غربی، زنجان، اردبیل
۶	کرمانشاه، کردستان، چهارمحال و بختیاری، لرستان، ایلام، کهگیلویه و بویر- احمد
۷	خوزستان، فارس، بوشهر، هرمزگان

جدول ۲: گروه بندی استان‌های همسایه

مطابق گروه‌بندی صورت گرفته در کل کشور، استان‌های همسایه با یک دیگر در یک گروه قرار گرفته‌اند.

در این پروژه قصد داریم هزینه و درآمد خانوار استان‌های گروه شماره ۶ را مورد بررسی قرار دهیم. پس نسبت به جداسازی این گروه اقدامات لازم را انجام می‌دهیم. همچنین با توجه به اینکه در هر ثبت داده خام، متغیرهایی اعم از ریز هزینه‌های خوراک، پوشاک، مسکن و ... وجود دارد پس تمام این ریز هزینه‌ها را جمع کرده و در متغیر مربوطه با یک عنوان وارد می‌نماییم.

در نتیجه مجموعه داده‌های گروه ۶ با ۳۱۰۷ ثبت و ۶۹ متغیر در دسترس قرار گرفته است.

معرفی متغیرهای مجموعه داده‌ها ۱

حال به معرفی متغیرهای مجموعه داده می‌پردازیم. (جدول شماره ۳)

ر	نام متغیر	معرفی و توضیحات
۱	X	اندیس شماره سطر
۲	Ostan	استان محل سکونت
۳	gorooh	گروه بندی صورت گرفته
۴	address	شناسه خانوار
۵	Nbkhanevar	نوع خانوار (۱: معمولی ساکن، ۲: گروهی)
۶	TedadAza	تعداد اعضا خانوار (عدد صحیح)
۷	SarparstJensiat	جنسیت سرپرست (۱: مرد، ۲: زن)
۸	SarparstSen	سن سرپرست (عدد صحیح)
۹	SarparstSavad	سواد سرپرست (۱: با سواد، ۲: بیسواد)
۱۰	SarparstTahsil	در حال تحصیل سرپرست (۱: بله، ۲: خیر)
۱۱	SarparstMadrak	مدرک تحصیلی سرپرست (۱: ابتدایی، ۲: راهنمایی، ۳: متوسطه، ۴: دیپلم، ۵: کاردانی، ۶: کارشناسی، ۷: کارشناسی ارشد، ۸: دکتری تخصصی، ۹: سایر و غیر رسمی)
۱۲	SarparstFaaliat	فعالیت سرپرست (۱: شاغل، ۲: بیکار (جویای کار)، ۳: دارای درآمد بدون کار، ۴: محصل، ۵: خانه دار، ۶: سایر)

۱۳	SarparstZanashoyi	زن‌اشویی سرپرست (۱: دارای همسر، ۲: بی‌همسر بر اثر فوت همسر، ۳: طلاق، ۴: هرگز ازدواج نکرده)
۱۴	NahveTasarof	نحوه تصرف محل سکونت (۱: ملکی عرصه و اعیان، ۲: ملکی اعیان، ۳: اجاری، ۴: رهن، ۵: در برابر خدمت، ۶: رایگان)
۱۵	TedadOtagh	تعداد اتاق (عدد صحیح)
۱۶	SatheZirbana	مساحت زیر بنا (برحسب متر مربع، عدد صحیح)
۱۷	NbeEskellet	نوع اسکلت (۱: فلزی، ۲: بتون ارمه، ۳: سایر)
۱۸	mashin	خودرو (۱: دارد، ۲: ندارد)
۱۹	motor	موتورسیکلت (۱: دارد، ۲: ندارد)
۲۰	docharakhe	دوچرخه (۱: دارد، ۲: ندارد)
۲۱	radio	رادیو (۱: دارد، ۲: ندارد)
۲۲	radiozabt	ضبط (۱: دارد، ۲: ندارد)
۲۳	tv	تلویزیون سیاه سفید (۱: دارد، ۲: ندارد)
۲۴	tvrangj	تلویزیون رنگی (۱: دارد، ۲: ندارد)
۲۵	video	ویدیو (۱: دارد، ۲: ندارد)
۲۶	computer	رایانه و تبلت (۱: دارد، ۲: ندارد)
۲۷	mobile	گوشی همراه (۱: دارد، ۲: ندارد)
۲۸	freezer	فریزر (۱: دارد، ۲: ندارد)
۲۹	yakhchal	یخچال (۱: دارد، ۲: ندارد)
۳۰	yakhchalfreezer	یخچال فریزر (۱: دارد، ۲: ندارد)
۳۱	qjaghgaz	اجاق گاز (۱: دارد، ۲: ندارد)
۳۲	jandbarghi	جاروبرقی (۱: دارد، ۲: ندارد)
۳۳	lebasshoyi	لباسشویی (۱: دارد، ۲: ندارد)
۳۴	khayati	چرخ خیاطی (۱: دارد، ۲: ندارد)
۳۵	panke	پنکه (۱: دارد، ۲: ندارد)
۳۶	coolerabimoteharek	کولر ابی متحرک (۱: دارد، ۲: ندارد)
۳۷	coolergazimoteharek	کولر گازی متحرک (۱: دارد، ۲: ندارد)
۳۸	zarfshoyi	ظرفشویی (۱: دارد، ۲: ندارد)
۳۹	microfer	مایکروفر (۱: دارد، ۲: ندارد)
۴۰	hichkodam	هیچکدام (۱: دارد، ۲: ندارد)
۴۱	ldekeshiab	لوله کشی اب (۱: دارد، ۲: ندارد)
۴۲	bargh	برقی (۱: دارد، ۲: ندارد)
۴۳	gazldekeshi	لوله کشی گاز (۱: دارد، ۲: ندارد)
۴۴	telephone	تلفن (۱: دارد، ۲: ندارد)
۴۵	internet	اینترنت (۱: دارد، ۲: ندارد)
۴۶	hamam	حمام (۱: دارد، ۲: ندارد)
۴۷	ashpazkhane	اشپزخانه (۱: دارد، ۲: ندارد)
۴۸	coolerabisabet	کولر ابی ثابت (۱: دارد، ۲: ندارد)
۴۹	borodatmarkazi	سرمایش مرکزی (۱: دارد، ۲: ندارد)
۵۰	hararatmarkazi	گرمایش مرکزی (۱: دارد، ۲: ندارد)
۵۱	package	پکیج (۱: دارد، ۲: ندارد)
۵۲	coolergazisabet	کولر گازی ثابت (۱: دارد، ۲: ندارد)
۵۳	Fazelabshahri	فاضلاب شهری (۱: دارد، ۲: ندارد)
۵۴	nsokhtpkhtpaz	سوخت پخت و پز (۳: گاز مایع، ۴: گاز طبیعی شبکه)
۵۵	nsokhtgamma	سوخت گرمایش (۱: نفت سفید، ۱۳: گاز مایع، ۱۴: گاز - طبیعی شبکه)
۵۶	nsokhtabgarm	سوخت آب گرم (۱: نفت سفید، ۲۳: گاز مایع، ۲۴: گاز طبیعی شبکه)
۵۷	Hazinekhoraki	هزینه خوراک (برحسب ریال، عددی)

مقادیر ناموجود به متغیر **SarparastSavad** که با سواد بودن یا بی سواد بودن سرپرست خانوار را تعیین میکند، ارتباط دارد. این مقادیر به دلیل بی سواد بودن سرپرست خانوار می باشد. با توجه به تعداد زیاد این مقادیر و اهمیت این متغیر این مقادیر ناموجود را با مقدار "۰" جایگزین میکنیم.

- متغیر **Hchkodam** به دلیل یکسان بودن در تمامی ثبت‌ها و همچنین با توجه به سوال پرسشنامه اطلاعاتی در اختیار ما قرار نمیدهد. پس این متغیر نیز باید حذف گردد.
- متغیر **bargh** نیز همانند متغیر فوق در تمامی ثبت‌ها دارای مقدار ۱ می باشد. این نشان دهنده آن است که در این نمونه گیری هیچ خانوار بدون اشتراک برق وجود ندارد و از آنجایی که این نمونه گیری از تمامی بخش‌های استان‌های مذکور جمع آوری شده است، پس نتیجه میگیریم در تمام بخش-های استان‌های مذکور برق شبکه خانگی وجود دارد. پس این متغیر نیز در ارائه مدل هیچ تأثیری ندارد.

نتیجه گیری فصل اول:

بنا به توضیحات فوق متغیرهای **Ostan** و **gorooch** و **address** و **Hchkodam** و **bargh** حذف گردید و مقادیر **NA** با صفر جایگزین گردید، در شکل شماره ۱ خلاصه‌ای از هر متغیر را مشاهده می کنید. همچنین برای درک بهتر مدل، از طریق دهک‌بندی، خانوارهای قرار گرفته در رده دهم را (با توجه به هدف پروژه) با یک و خانوارهای دیگر را با صفر نشان می دهیم.

همانطور که در خروجی زیر مشاهده میکنید دهک دهم، درآمدی بین ۶۵۴۸۳۳۳۳۳ ریال و ۷۰۸۸۲۶۶۷ ریال دارند.

```
> quantile(final.df$Target, probs = seq(.1, .9, by = .1))
 10%    20%    30%    40%    50%    60%    70%    80%    90%
8690000 15113333 19694000 23881333 28956667 34389667 41790333 51916667 70882667
> #calculate deciles of dataset
> quantile(final.df$Target, probs = seq(.9, by = .1))
 90%    100%
70882667 654833333
```

و همچنین تعداد خانوار قرار گرفته در هر رده به شرح زیر است.

```
> as.data.frame(table(final.df$decile))
  Var1 Freq
1     1   311
2     2   311
3     3   311
4     4   311
5     5   311
6     6   311
7     7   311
8     8   310
9     9   310
10    10   310
```

پس مطابق گفته‌های بالا داریم.

```
> as.data.frame(table(final.df$decile))
  Var1 Freq
1     0 2797
2     1  310
```

هزینه نوشیدنی (برحسب ریال ، عددی)	HazineNshidani	۵۸
هزینه پوشاک (برحسب ریال ، عددی)	HazinePushak	۵۹
هزینه مسکن (برحسب ریال ، عددی)	HazineMaskan	۶۰
هزینه لوازم خانگی (برحسب ریال ، عددی)	HazineLavazemKhanegi	۶۱
هزینه درمانی (برحسب ریال ، عددی)	HazineDarmani	۶۲
هزینه حمل و نقل (برحسب ریال ، عددی)	HazineHamtonaghl	۶۳
هزینه ارتباطات (برحسب ریال ، عددی)	HazineErtebatat	۶۴
هزینه تفریحات فرهنگی (برحسب ریال ، عددی)	HazineTafrihatFarhangi	۶۵
هزینه غذا آماده (برحسب ریال ، عددی)	HazineGhazaAmade	۶۶
هزینه کالا متفرقه (برحسب ریال ، عددی)	HazineKalaMotefaregheh	۶۷
هزینه کالا بادوام (برحسب ریال ، عددی)	HazineKalaBadavam	۶۸
درآمد (برحسب ریال ، عددی)	Target	۶۹

جدول ۳: معرفی متغیرهای مجموعه داده

همان طور که مشاهده میکنید اکثر متغیرهای فوق نیازی به توضیحات خاصی ندارند. متغیرهای **Hazine** همگی مرتبط به ماه گذشته می باشد به غیر از **HazineKalaBadavam** که مرتبط با هزینه ها در یکسال گذشته می باشد.

تنها متغیر **Target** که از مشورت و همفکری دانشجویان دانشگاه شهید بهشتی طبق معادله زیر ایجاد گردیده است.

$$Target = (SarmayeGozari + DaramadNKH12 + AzadDaramadKh + MotefaregheHoghogh + Motefareghe Ejareh + MotefareghePasandaz + MotefaregheKomakhazine + MotefaregheForoshMasnoat + MotefaregheDaryafiti)/12$$

حال به بررسی دقیق تر برخی از متغیرها می پردازیم.

- متغیر **Ostan** مرتبط با محل سکونت خانوار می باشد و تا قبل از گروه بندی کلی مورد استفاده گرفته است و با توجه اینکه در ادامه پروژه اطلاعات مفیدی در اختیار ما نمیگذارد و تمام اطلاعات لازم در متغیر آدرس موجود است، پس می بایست حذف گردد.
- متغیر **gorooch** نیز مرتبط با موضوع فوق است و با توجه به یکسان بودن آن برای هر سطر نیز در ادامه هیچ اطلاعاتی برای ما ندارد.
- متغیر **address** نیز با توجه به اینکه شناسه هر خانوار است و برای هر خانوار یکتاست پس همانند دو متغیر فوق در ادامه پروژه کمک شایانی نمیکند.
- متغیر **SarparastTahsil** و **SarparastMadrak** هر کدام دارای ۶۰۸ مقدار **NA** یا همان مقادیر ناموجود می باشند. این

متغیر decile متغیر برآمد است.

microfer	tolekeshiab	gaztolekeshi	X	NoeKhanevar	TedadAza
Min.: 0.00000	Min.: 0.000	Min.: 0.0000	Min.: 685	Min.: 1.000	Min.: 1.000
1st Qu.: 0.00000	1st Qu.: 1.000	1st Qu.: 1.0000	1st Qu.: 5548	1st Qu.: 1.000	1st Qu.: 3.000
Median: 0.00000	Median: 1.000	Median: 1.0000	Median: 10366	Median: 1.000	Median: 4.000
Mean: 0.04731	Mean: 0.999	Mean: 0.9865	Mean: 9359	Mean: 1.001	Mean: 3.551
3rd Qu.: 0.00000	3rd Qu.: 1.000	3rd Qu.: 1.0000	3rd Qu.: 12280	3rd Qu.: 1.000	3rd Qu.: 4.000
Max.: 1.00000	Max.: 1.000	Max.: 1.0000	Max.: 17117	Max.: 2.000	Max.: 11.000
telephone	internet	hamam	SarparstJensiat	SarparstSen	SarparstSavad
Min.: 0.0000	Min.: 0.0000	Min.: 0.0000	Min.: 1.000	Min.: 18.00	Min.: 1.000
1st Qu.: 0.0000	1st Qu.: 1.0000	1st Qu.: 1.0000	1st Qu.: 1.000	1st Qu.: 39.00	1st Qu.: 1.000
Median: 0.0000	Median: 1.0000	Median: 1.0000	Median: 1.000	Median: 49.00	Median: 1.000
Mean: 0.3949	Mean: 0.7512	Mean: 0.9958	Mean: 1.128	Mean: 51.04	Mean: 1.196
3rd Qu.: 1.0000	3rd Qu.: 1.0000	3rd Qu.: 1.0000	3rd Qu.: 1.000	3rd Qu.: 62.00	3rd Qu.: 1.000
Max.: 1.0000	Max.: 1.0000	Max.: 1.0000	Max.: 2.000	Max.: 99.00	Max.: 2.000
ashpazkhane	coolerabisabet	borodatmarkazi	SarparstTahsil	SarparstMadrak	SarparstFaaliat
Min.: 0.0000	Min.: 0.0000	Min.: 0.00000000	Min.: 0.000	Min.: 0.0	Min.: 1.00
1st Qu.: 1.0000	1st Qu.: 0.0000	1st Qu.: 0.00000000	1st Qu.: 2.000	1st Qu.: 1.0	1st Qu.: 1.00
Median: 1.0000	Median: 1.0000	Median: 0.00000000	Median: 2.000	Median: 2.0	Median: 1.00
Mean: 0.9907	Mean: 0.5401	Mean: 0.0009656	Mean: 1.603	Mean: 2.7	Mean: 1.72
3rd Qu.: 1.0000	3rd Qu.: 1.0000	3rd Qu.: 0.00000000	3rd Qu.: 2.000	3rd Qu.: 4.0	3rd Qu.: 3.00
Max.: 1.0000	Max.: 1.0000	Max.: 1.00000000	Max.: 2.000	Max.: 9.0	Max.: 6.00
hararatmarkazi	package	coolergazisabet	SarparstZanashoyi	NahveTasarof	TedadOtagh
Min.: 0.0000000	Min.: 0.000000	Min.: 0.0000	Min.: 1.000	Min.: 1.000	Min.: 1.00
1st Qu.: 0.000000	1st Qu.: 0.000000	1st Qu.: 0.0000	1st Qu.: 1.000	1st Qu.: 1.000	1st Qu.: 3.00
Median: 0.000000	Median: 0.000000	Median: 0.0000	Median: 1.000	Median: 1.000	Median: 4.00
Mean: 0.001931	Mean: 0.03991	Mean: 0.1381	Mean: 1.195	Mean: 1.943	Mean: 3.76
3rd Qu.: 0.000000	3rd Qu.: 0.000000	3rd Qu.: 0.0000	3rd Qu.: 1.000	3rd Qu.: 3.000	3rd Qu.: 4.00
Max.: 1.000000	Max.: 1.000000	Max.: 1.0000	Max.: 4.000	Max.: 6.000	Max.: 8.00
fazelabshahri	nsokhtpokhtpaz	nsokhtgarma	SatheZirbana	NoeEskelet	maskhin
Min.: 0.0000	Min.: 0.0000	Min.: 11.00	Min.: 12.0	Min.: 1.000	Min.: 0.0000
1st Qu.: 0.0000	1st Qu.: 4.000	1st Qu.: 14.00	1st Qu.: 80.0	1st Qu.: 2.000	1st Qu.: 0.0000
Median: 1.0000	Median: 4.000	Median: 14.00	Median: 100.0	Median: 3.000	Median: 0.0000
Mean: 0.5452	Mean: 3.985	Mean: 13.96	Mean: 107.1	Mean: 2.556	Mean: 0.4493
3rd Qu.: 1.0000	3rd Qu.: 4.000	3rd Qu.: 14.00	3rd Qu.: 125.0	3rd Qu.: 3.000	3rd Qu.: 1.0000
Max.: 1.0000	Max.: 4.000	Max.: 14.00	Max.: 350.0	Max.: 3.000	Max.: 1.0000
nsokhtabgarm	HazineKhoraki	HazineNoshidani	motor	docharkhe	radio
Min.: 21.00	Min.: 600000	Min.: 0	Min.: 0.0000	Min.: 0.0000	Min.: 0.000000
1st Qu.: 24.00	1st Qu.: 8426250	1st Qu.: 0	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.: 0.000000
Median: 24.00	Median: 12199000	Median: 0	Median: 0.0000	Median: 0.0000	Median: 0.000000
Mean: 23.98	Mean: 13623076	Mean: 319900	Mean: 0.1007	Mean: 0.1435	Mean: 0.09495
3rd Qu.: 24.00	3rd Qu.: 16815000	3rd Qu.: 0	3rd Qu.: 0.0000	3rd Qu.: 0.0000	3rd Qu.: 0.000000
Max.: 25.00	Max.: 112235000	Max.: 9000000	Max.: 1.0000	Max.: 1.0000	Max.: 1.000000
HazinePoshak	HazineMaskan		radiozabt	tv	tvrangei
Min.: 0	Min.: 400000		Min.: 0.00000	Min.: 0.000000	Min.: 0.0000
1st Qu.: 0	1st Qu.: 5900000		1st Qu.: 0.0000	1st Qu.: 0.000000	1st Qu.: 1.0000
Median: 380000	Median: 8220000		Median: 0.0000	Median: 0.000000	Median: 1.0000
Mean: 1990246	Mean: 9248438		Mean: 0.1007	Mean: 0.004506	Mean: 0.9781
3rd Qu.: 2480000	3rd Qu.: 11500000		3rd Qu.: 0.0000	3rd Qu.: 0.000000	3rd Qu.: 1.0000
Max.: 68150000	Max.: 120000000		Max.: 1.0000	Max.: 1.000000	Max.: 1.0000
HazineLavazemKhanegi	HazineDarmani		video	computer	mobile
Min.: 0	Min.: 0		Min.: 0.0000	Min.: 0.0000	Min.: 0.0000
1st Qu.: 700000	1st Qu.: 200000		1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.: 1.0000
Median: 1090000	Median: 520000		Median: 0.0000	Median: 0.0000	Median: 1.0000
Mean: 1317564	Mean: 2162923		Mean: 0.2111	Mean: 0.2414	Mean: 0.9652
3rd Qu.: 1600000	3rd Qu.: 1600000		3rd Qu.: 0.0000	3rd Qu.: 0.0000	3rd Qu.: 1.0000
Max.: 26800000	Max.: 500300000		Max.: 1.0000	Max.: 1.0000	Max.: 1.0000
HazineHamilonaghl	HazineErtebatat		freezer	yakhchal	yakhchalfreezer
Min.: 0	Min.: 0		Min.: 0.0000	Min.: 0.000	Min.: 0.0000
1st Qu.: 450000	1st Qu.: 400000		1st Qu.: 0.0000	1st Qu.: 0.000	1st Qu.: 0.0000
Median: 1000000	Median: 670000		Median: 0.0000	Median: 0.000	Median: 1.0000
Mean: 1502440	Mean: 803620		Mean: 0.2066	Mean: 0.383	Mean: 0.6292
3rd Qu.: 2000000	3rd Qu.: 1000000		3rd Qu.: 0.0000	3rd Qu.: 1.000	3rd Qu.: 1.0000
Max.: 47800000	Max.: 20100000		Max.: 1.0000	Max.: 1.000	Max.: 1.0000
HazineTafrihatFarhangi	HazineGhazaAmade		ojaghgaz	jarobarghi	lebasshoyi
Min.: 0	Min.: 0		Min.: 0.00	Min.: 0.0000	Min.: 0.0000
1st Qu.: 0	1st Qu.: 0		1st Qu.: 1.00	1st Qu.: 1.0000	1st Qu.: 1.0000
Median: 0	Median: 0		Median: 1.00	Median: 1.0000	Median: 1.0000
Mean: 202387	Mean: 224337		Mean: 0.99	Mean: 0.8999	Mean: 0.7734
3rd Qu.: 0	3rd Qu.: 200000		3rd Qu.: 1.00	3rd Qu.: 1.0000	3rd Qu.: 1.0000
Max.: 180250000	Max.: 100000000		Max.: 1.00	Max.: 1.0000	Max.: 1.0000
HazineKalaMotefaregheh	HazineKalaBadavam		khayati	panke	coolerabimoteharek
Min.: 0	Min.: 0.000e+00		Min.: 0.0000	Min.: 0.0000	Min.: 0.0000
1st Qu.: 450000	1st Qu.: 1.126e+07		1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.: 0.0000
Median: 850000	Median: 5.183e+07		Median: 0.0000	Median: 0.0000	Median: 0.0000
Mean: 1111492	Mean: 1.075e+08		Mean: 0.2716	Mean: 0.2095	Mean: 0.1059
3rd Qu.: 1400000	3rd Qu.: 1.350e+08		3rd Qu.: 1.0000	3rd Qu.: 0.0000	3rd Qu.: 0.0000
Max.: 16790000	Max.: 2.187e+09		Max.: 1.0000	Max.: 1.0000	Max.: 1.0000
Target			coolergazimoteharek	zarfshoyi	
Min.: 0			Min.: 0.00000	Min.: 0.00000	
1st Qu.: 17553333			1st Qu.: 0.00000	1st Qu.: 0.00000	
Median: 28956667			Median: 0.00000	Median: 0.00000	
Mean: 37834783			Mean: 0.01867	Mean: 0.02961	
3rd Qu.: 46404167			3rd Qu.: 0.00000	3rd Qu.: 0.00000	
Max.: 654833333			Max.: 1.00000	Max.: 1.00000	

شکل ۱: خلاصه‌ای از وضعیت هر متغیر مجموعه داده‌ها

در نهایت این مجموعه داده را با نسبت ۱۰/۳۰/۶۰ به ترتیب مجموعه

آموزشی، مجموعه اعتبارسنجی و مجموعه آزمون تقسیم می‌کنیم.

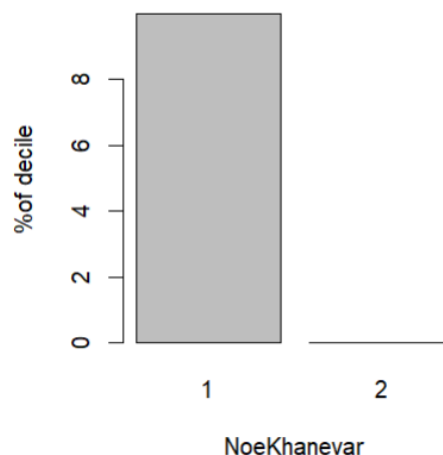
فصل دوم: تصویری سازی

در این فصل، مجموعه‌ای از نمودارها را که می‌توانند برای اکتشاف ماهیت چندمتغیره یک مجموعه داده، مورد استفاده قرار گیرند را تشریح می‌کنیم و بر استفاده از ارائه‌های گرافیکی با هدف اکتشاف داده‌ها، به ویژه در رابطه با تحلیل خودکار پیشگویانه، تمرکز می‌کنیم.^۱ جمله معروف (یک عکس بیش از هزار کلمه می‌ارزد)، به توانایی فشرده‌سازی اطلاعات شفاهی پراکنده در یک تصویر گرافیکی جمع و جور و سریعاً فهمیدنی، باز می‌گردد. در مورد اعداد، تصویری سازی داده‌ها و خلاصه سازی عددی، ما را به یک ابزار قدرتمند برای اکتشاف داده‌ها و یک راه موثر برای ارائه نتایج مجهز می‌سازد.^۲

بسته "ggplot" که هاردلی ویکمن تهیه کرده‌اند، به رایج‌ترین بسته گرافیکی اختصاصی در R برای تصویری سازی در گستره وسیعی از زمینه‌ها بدل شده است.^۳

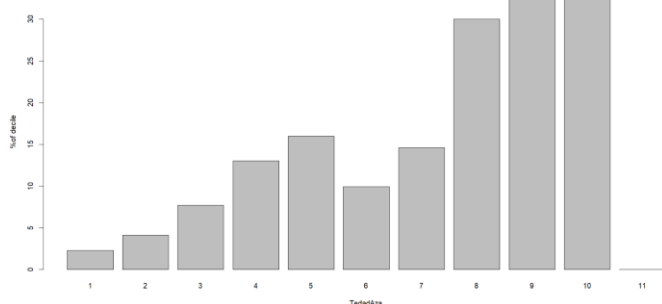
در ابتدا برای درک بهتر، متغیرها را با نمودارهایی نشان می‌دهیم. سه نمودار پایه که موثرترین نمودارها هستند، عبارت‌اند از نمودارهای میله‌ای، نمودارهای خطی و نمودارهای پراکنش.

حال نمودارهایی را براساس متغیر پیشگو و متغیرهای برآمد نشان می‌دهیم.



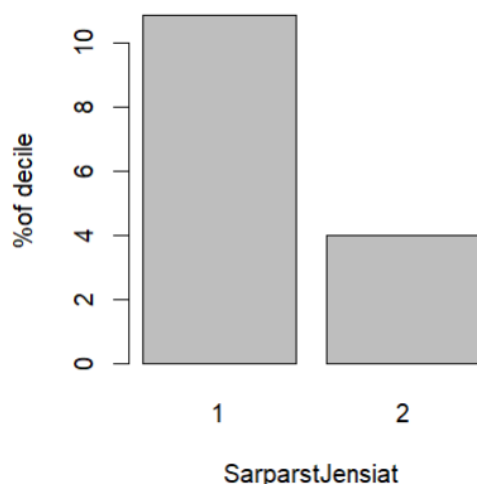
نمودار ۱: نمودار درصد خانوار دهک دهم نسبت به نوع خانوار

در ابتدا درصد خانوار دهک دهم را نسبت به نوع خانوار نشان می‌دهیم. نمودار ۱ نشان دهنده آن است که ۹۹,۹۸ درصد از خانوارهای نوع ۱ که همان معمولی ساکن است، در دهک دهم درآمدی قرار دادند. پس یعنی می‌توان گفت که تقریباً تمام خانوارهای قرار گرفته دهک دهم نمونه، جز این دسته از خانوارها هستند.



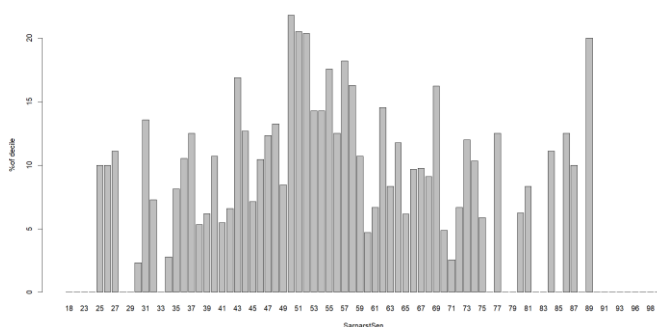
نمودار ۲: نمودار درصد خانوار دهک دهم نسبت به تعداد اعضا خانوار

در نمودار ۲، به راحتی می‌توان مشاهده نمود که خانوارهای دارای ۸، ۹ و ۱۰ عضو با درصد بیشتری در دهک دهم این نمونه را قرار گرفته‌اند.



نمودار ۳: درصد خانوارهای قرار گرفته در دهک دهم نسبت به جنسیت سرپرست.

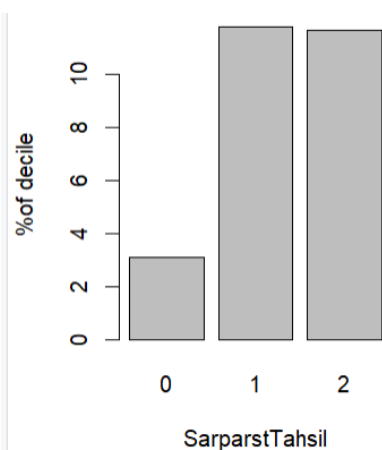
همانطور که در نمودار ۳ مشاهده می‌کنید عمده سرپرستان خانوار قرار گرفته در دهک دهم را مردان تشکیل می‌دهند.



نمودار ۴: درصد خانوارهای دهک دهم نسبت به سن سرپرستان خانوار

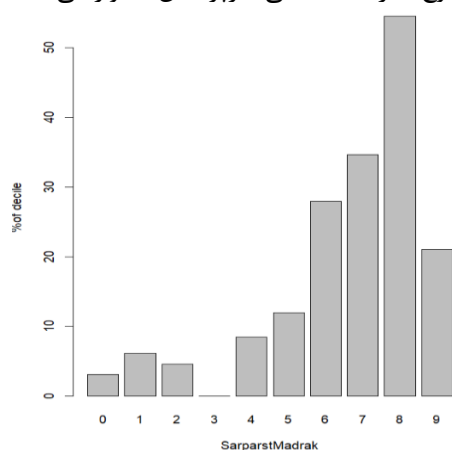
همانطور که در نمودار ۴ مشاهده می‌نمایید، عمده سرپرستان خانوار در رده سنی ۳۴ الی ۶۵ سال می‌باشند. پس می‌توان این رده‌های سنی را در ۴ دسته کلی به صورت زیر تقسیم نمود. (نمودار ۵)

- ۱۸ الی ۳۰ سال به عنوان نوجوان (۲۴)
- ۳۱ الی ۴۰ سال به عنوان جوان (۳۵)
- ۴۱ الی ۶۰ سال به عنوان میانسال (۵۰)
- ۶۰ الی ۱۰۰ سال به عنوان سالمند (۸۰)



نمودار ۷: نسبت دهک درآمدی دهم به وضعیت تحصیل سرپرست خانوار

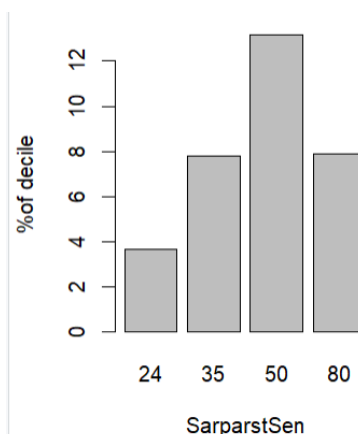
در رابطه با سطح تحصیلات سرپرست خانوار می‌توانیم متغیرهای در این خصوص را ادغام کنیم. در واقع همانطور که ذکر شد می‌توانیم اطلاعات لازم را از متغیر مدرک تحصیلی نیز استخراج کنیم. نمودار ۸، نشان دهنده نوع مدرک تحصیلی سرپرستان خانوار می‌باشد.



نمودار ۸: تفکیک مدرک تحصیلی سرپرستان خانوار در دهک دهم نمونه

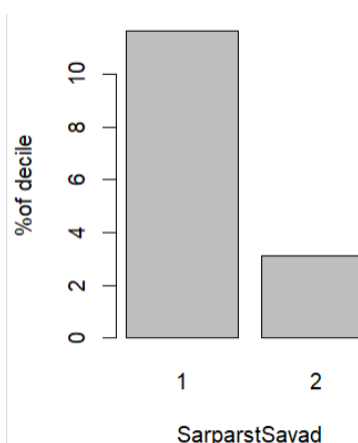
همانطور که در نمودار ۸ مشاهده می‌کنید، درصد بیشتری از سرپرستان خانوار قرار گرفته در دهک دهم درآمدی، مطابق بخش تعریف متغیرها، دارای مدرک تحصیلی کارشناسی/لیسانس، کارشناسی ارشد و دکتری عمومی/ فوق لیسانس، دکترای تخصصی هستند. پس بنابراین می‌توان نتیجه گرفت که مدرک تحصیلی سرپرستان در درآمد آن‌ها تاثیر به سزایی دارد. به عنوان مثال سرپرست خانواری که دارای مدرک تحصیلی دکترای تخصصی می‌باشد حتما در دهک دهم قرار می‌گیرد.

در بررسی متغیر **SarpartFaaliat** در میابیم که سرپرستان شاغل درصد بیشتری از دهک درآمدی دهم را شامل می‌شوند. همچنین تعداد افراد دارای درآمد بدون کار نیز تعداد قابل توجهی دارند. البته لازم به



نمودار ۵: نمودار درصد خانوار دهک دهم نسبت به رده بندی سنی

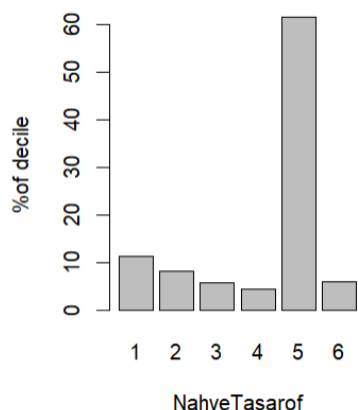
در نمودار ۶، مشهود است که درصد افراد باسواد قرار گرفته در دهک دهم به مراتب بالاتر از درصد افرادی بی‌سوادی است که در این دهک قرار گرفته اند. مطابق نمودار ۶ باسواد بودن سرپرستان در قرارگیری در دهک دهم متاثر است.



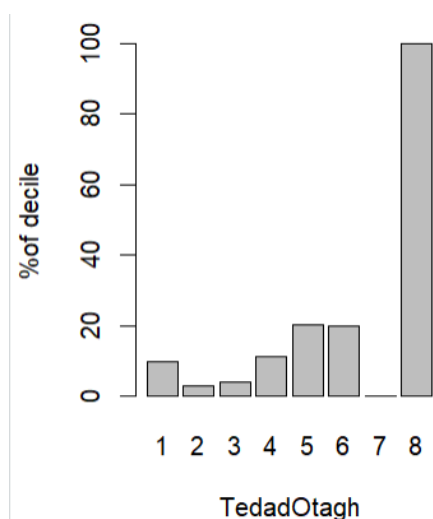
نمودار ۶: باسواد یا بی‌سواد بودن سرپرستان خانوار در دهک دهم

همچنین میتوان درصد سرپرستان در حال تحصیل و فارغ التحصیل را نسبت به دهک درآمدی مشاهده نمود. لازم به ذکر است مقادیر صفر نشان دهنده افراد بی‌سواد است.^۱

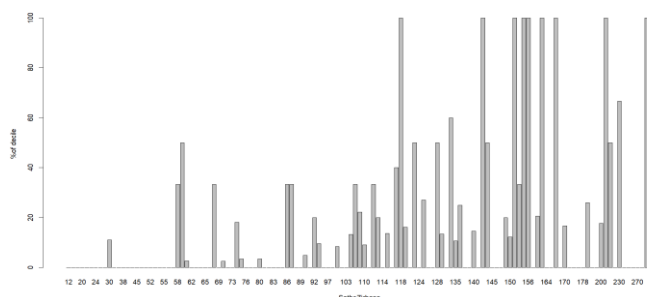
از آنجایی که درصد سرپرست خانوارهای در حال تحصیل و فارغ-التحصیل قرار گرفته در دهک دهم، تفاوت معناداری ندارد پس می‌توان گفت که این متغیر تاثیری در ارائه مدل بهتر ندارد.



نمودار ۱۱: نمودار نسبت دهک دهم خانوار به وضعیت محل سکونت خانوار
از نمودار متغیر **TedadOtagh** در میابیم که اگر خانواری ۸ اتاق داشته-
باشد، ۱۰۰ درصد در دهک دهم قرار میگیرد. البته لازم به ذکر است با
توجه به تعداد اتاق سایر خانوارهای قرار گرفته در دهک دهم بدهی
است که خانوارهای با ۵ اتاق و بالاتر با درصد بیشتری در دهک دهم
قرار می گیرند.

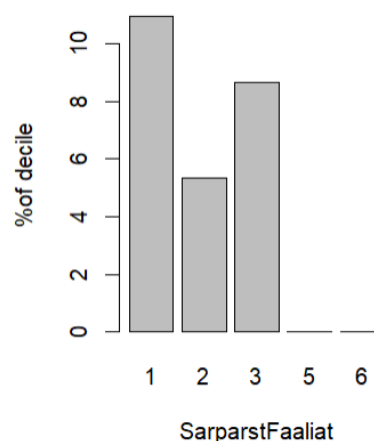


نمودار ۱۲: درصد قرارگیری خانوارها در دهک دهم نسبت به تعداد اتاق محل سکونت
در بررسی متغیر **SatheZirbana** درمی یابیم که همانند متغیر تعداد
اتاق در اختیار خانوار، سطح زیر بنا نیز دقیقاً چنین رفتاری دارد. یعنی
در نمودار ۱۳ خانوارهایی که در خانه هایی با متراژ بالای ۱۱۰ متر مربع
سکونت می کنند با درصد بیشتری در دهک درآمدی دهم قرار گرفته اند.



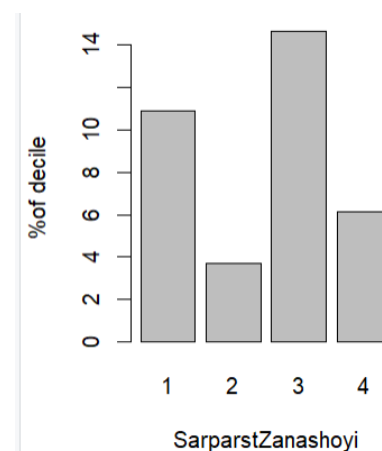
نمودار ۱۳: نمودار درصد خانوارهای قرار گرفته در دهک دهم نسبت به متراژ محل سکونت

ذکر است با توجه به درصد افراد جویای کار قرار گرفته در این دهک
می بایست در خصوص این افراد بررسی های دقیق تری صورت پذیرد.



نمودار ۹: نسبت خانوار قرار گرفته در دهک دهم به وضعیت شغلی سرپرستان خانوار

در بررسی متغیر **SarpartZanashoyi** در میابیم که وضعیت زناشویی
افراد قرار گرفته در دهک دهم در وضعیت طلاق بسیار بیشتر از سایر
وضعیت هاست، و این نشان دهنده آن است که به دلیل طلاق و
مشکلات ناشی از آن به طرز محسوسی درآمد خانوار افزایش یافته است.

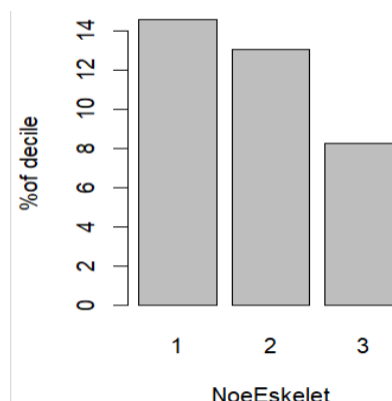
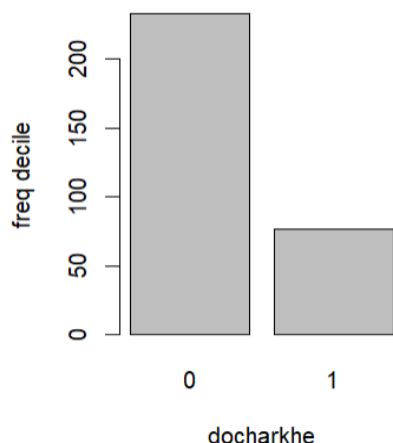


نمودار ۱۰: نمودار نسبت خانوار دهک دهم درآمدی به وضعیت زناشویی سرپرست خانوار

در ادامه به وضعیت محل سکونت خانوار و ارتباط آن با درآمد افراد
می پردازیم.

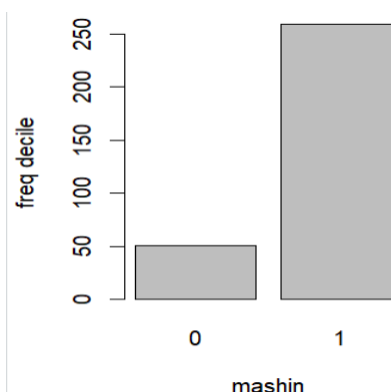
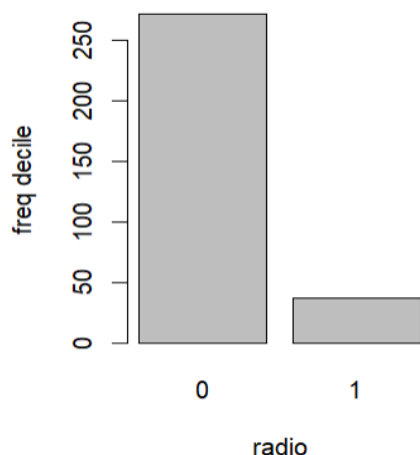
همانطور که در نمودار ۱۱ به وضوح دیده می شود، اکثر خانوارهایی که
نحوه تصرف محل سکونتشان از نوع در برابر خدمت است. یعنی خانوار
به ازای کارکرد در ارگان ها یا سازمان های دولتی و غیردولتی مشمول
سکونت در خانه های به اصطلاح سازمانی گردیده است. درصد
قرارگیری آن ها در دهک مذکور به مراتب بیشتر از سایر وضعیت ها
است. به عبارتی ۶۰ درصد خانوار هایی که نحوه تصرف محل
سکونتشان از نوع مذکور است در دهک دهم قرار گرفته اند.

در نمودار ۱۴، درصد بیشتر نوع اسکلت ساختمان خانوارهای قرار گرفته در دهک دهم، از نوع فلزی است.



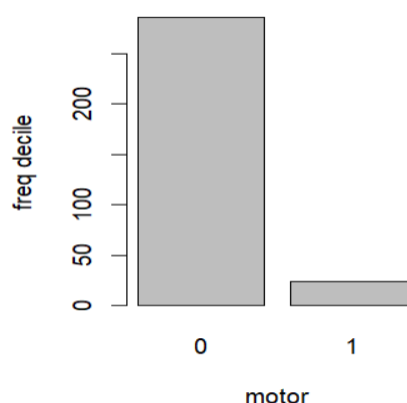
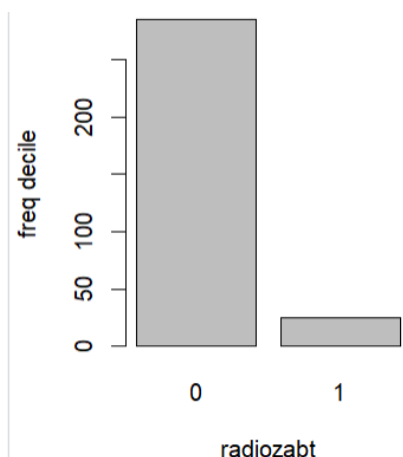
نمودار ۱۴ درصد: نوع اسکلت خانوارهای قرار گرفته در دهک دهم درآمدی

در نمودارهای زیر نیز تعداد خانوارهای دارای وسایل حمل و نقل قرار گرفته در دهک دهم را مشاهده می‌نمایید. لذا از توضیح این نمودار-ها به دلیل سادگی صرف نظر می‌کنیم.



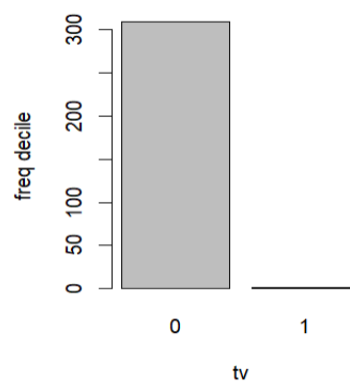
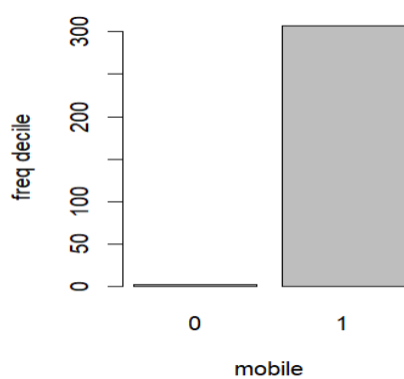
نمودار ۱۸: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن رادیو

نمودار ۱۵: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن خودرو



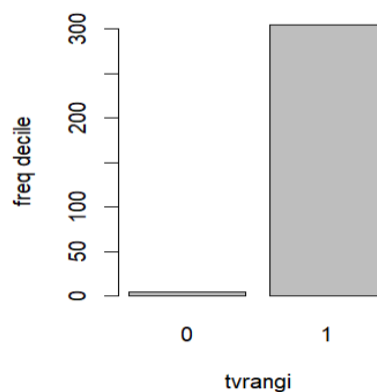
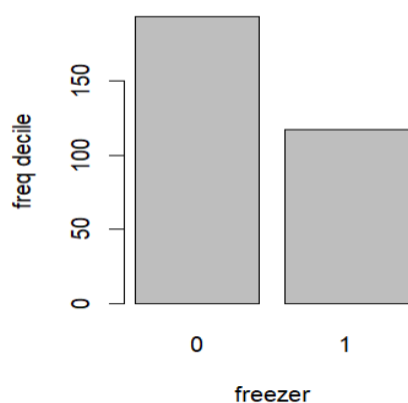
نمودار ۱۹: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن رادیو ضبط

نمودار ۱۶: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن موتورسیکلت



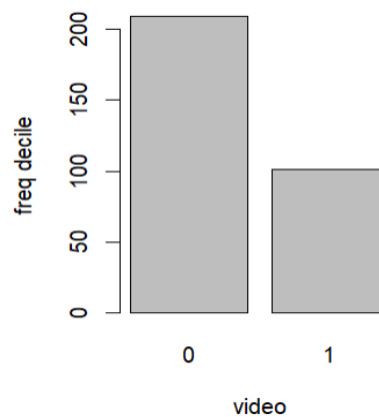
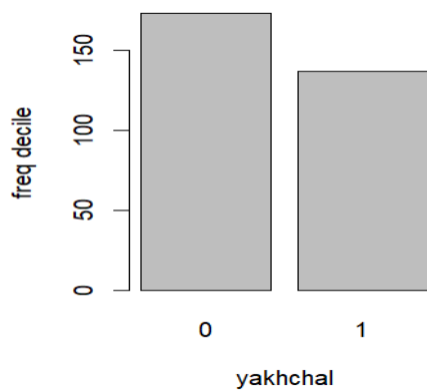
نمودار ۲۴: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن گوشی همراه

نمودار ۲۰: تعداد خانوارهای قرار گرفته در دهک دهم از لحاظ دارا بودن تلویزیون سیاه سفید



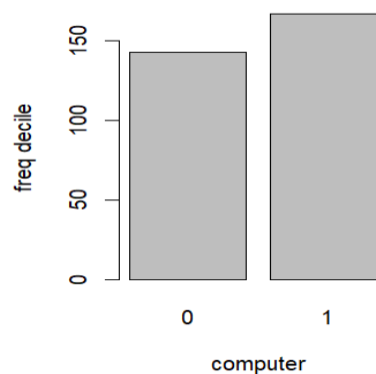
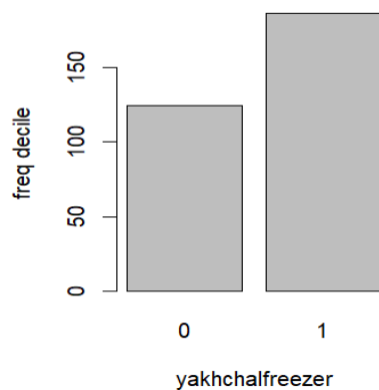
نمودار ۲۵: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن فریزر

نمودار ۲۱: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن تلویزیون رنگی



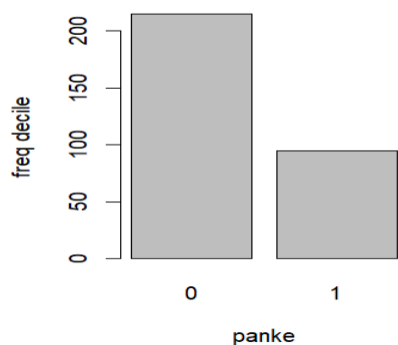
نمودار ۲۶: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن یخچال

نمودار ۲۲: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن دستگاه ویدیو

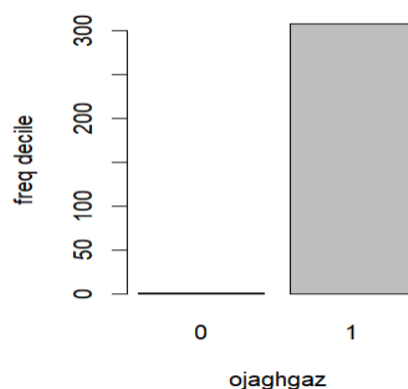


نمودار ۲۷: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن یخچال فریزر

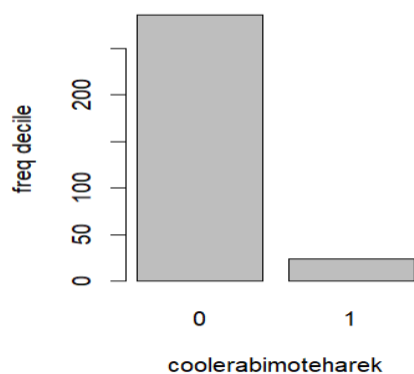
نمودار ۲۳: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن کامپیوتر و تبلت



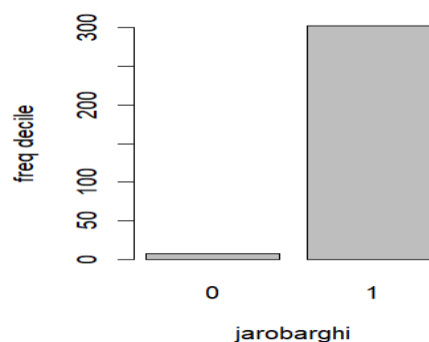
نمودار ۳۲: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن پنکه



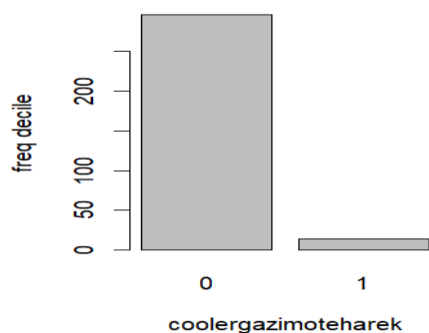
نمودار ۲۸: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن اجاق گاز



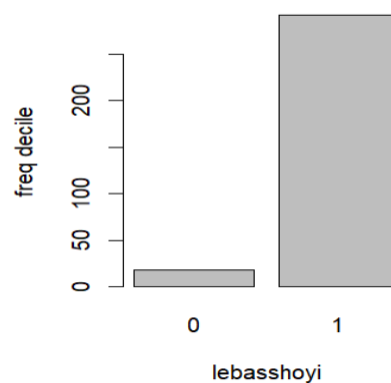
نمودار ۳۳: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن کولر ایبی متحرک



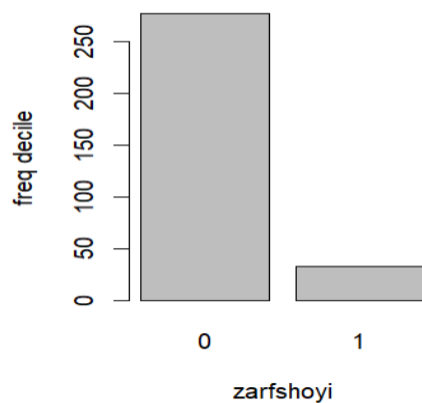
نمودار ۲۹: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن جاروبرقی



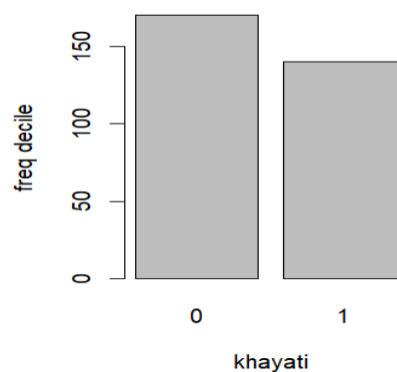
نمودار ۳۴: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن کولرگازی متحرک



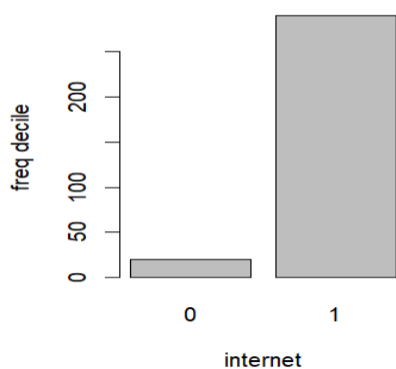
نمودار ۳۰: تعداد خانوارهای قرار گرفته در دهک دهم از لحاظ دارا بودن ماشین لباسشویی



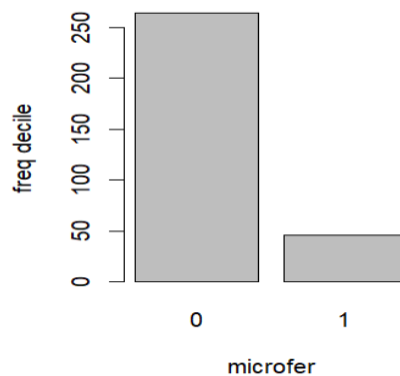
نمودار ۳۵: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن ماشین ظرفشویی



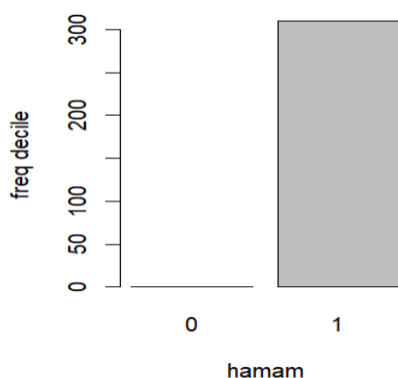
نمودار ۳۱: تعداد خانوار های قرار گرفته در دهک دهم از لحاظ دارا بودن چرخ خیاطی



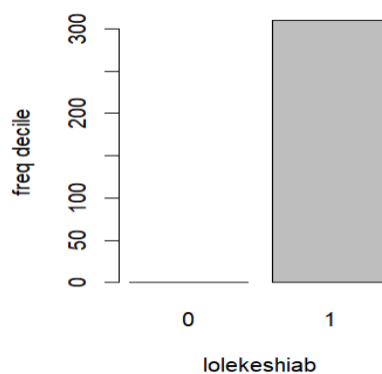
نمودار ۴۰: تعداد خانوارهای دارای دسترسی به شبکه اینترنت قرار گرفته در دهک دهم



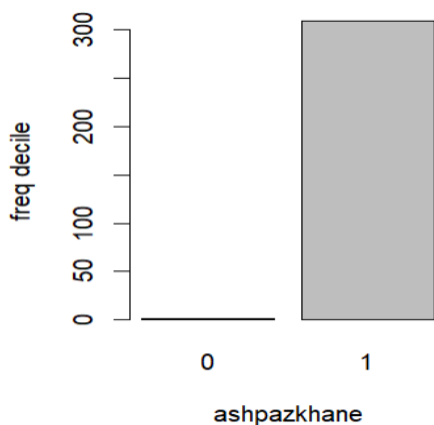
نمودار ۳۶: تعداد خانوارهای قرار گرفته در دهک دهم از لحاظ دارا بودن مایکروفر



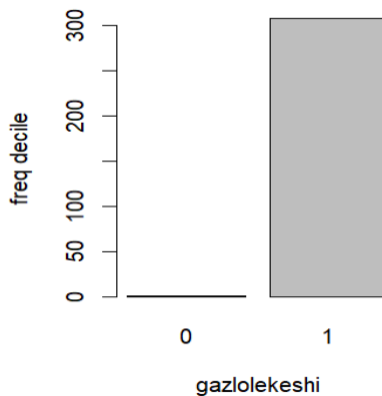
نمودار ۴۱: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن حمام



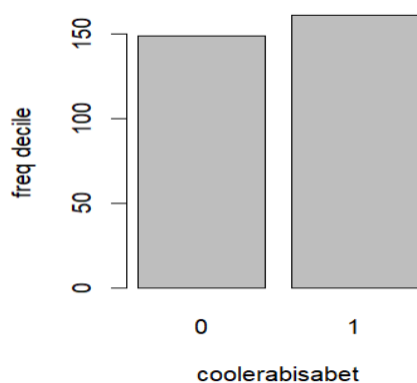
نمودار ۳۷: تعداد خانوارهای قرار گرفته در دهک دهم از لحاظ دارا بودن انشعاب آب لوله کشی



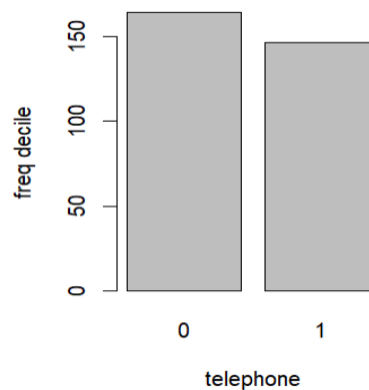
نمودار ۴۲: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن آشپزخانه



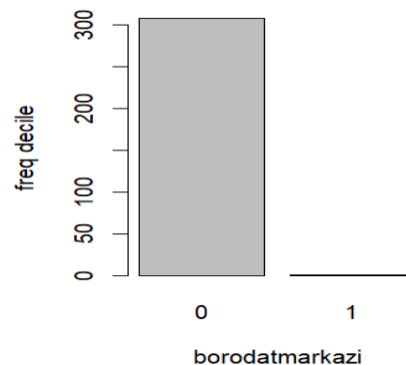
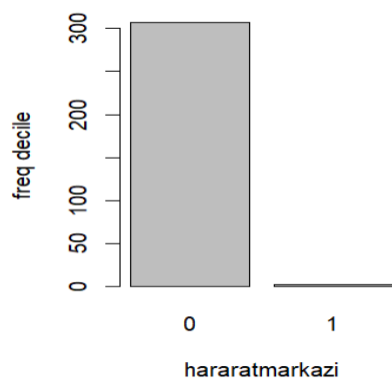
نمودار ۳۸: تعداد خانوارهای قرار گرفته در دهک دهم از لحاظ دارا بودن انشعاب گاز لوله کشی



نمودار ۴۳: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن کولر ایبی ثابت

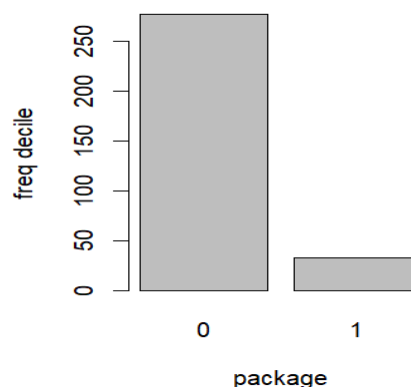
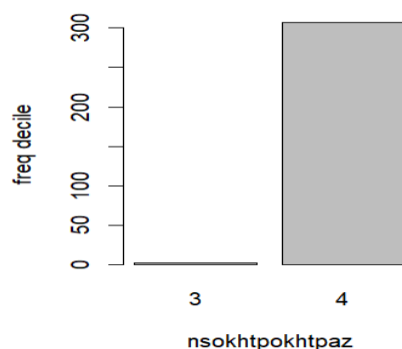


نمودار ۳۹: تعداد خانوارهای قرار گرفته در دهک دهم از لحاظ دارا بودن خط تلفن



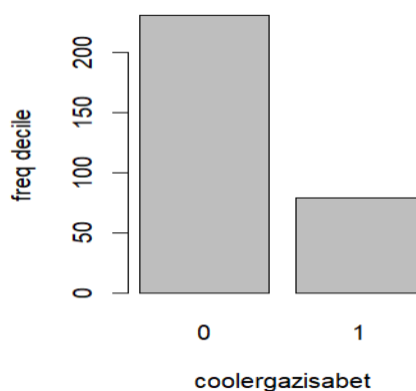
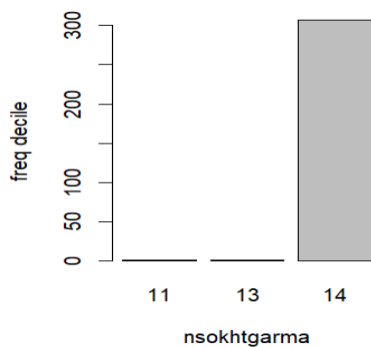
نمودار ۴۴: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن برودت مرکزی

نمودار ۴۸: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن سیستم حرارت مرکزی
در ادامه به بررسی نوع سوخت مورد نیاز برای پخت و پز و ایجاد گرما
برای مصارف گوناگون می پردازیم.



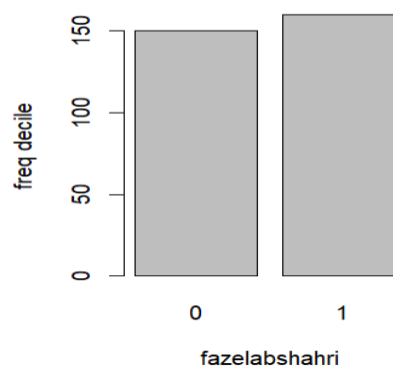
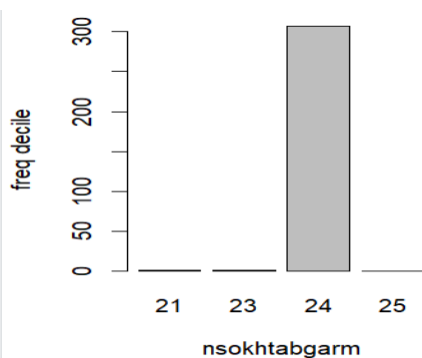
نمودار ۴۵: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن پکیج

نمودار ۴۹: نوع سوخت خانوار های قرار گرفته در دهک دهم برای پخت و پز



نمودار ۴۶: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن کولر گازی ثابت

نمودار ۵۰: نوع سوخت خانوار قرار گرفته در دهک دهم برای گرمایش

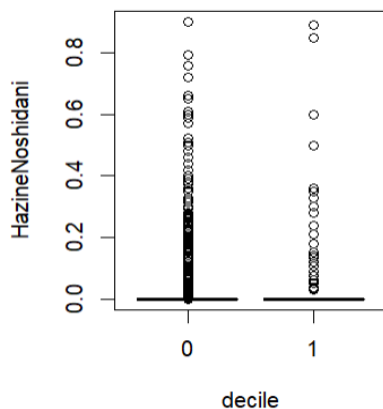


نمودار ۴۷: تعداد خانوار قرار گرفته در دهک دهم از لحاظ دارا بودن انشعاب فاضلاب شهری

نمودار ۵۱: نوع سوخت مورد استفاده خانوارهای قرار گرفته در دهک دهم برای تهیه آبگرم

همانطور که مشهود است در صدک ۵۰ خانوارهای قرار گرفته در دهک دهم میزان هزینه بیش از ۲ میلیون تومان است.

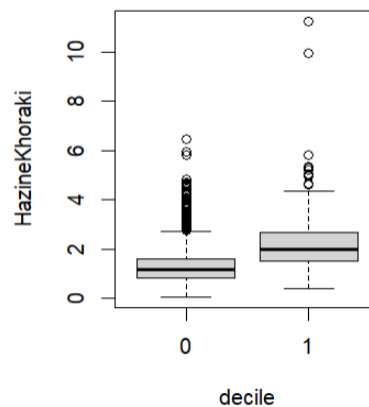
در نمودار ۵۳ نیز می توان به همین نحو نتیجه گیری کرد.



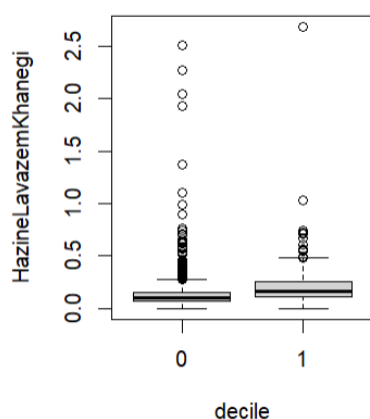
نمودار ۵۳: هزینه نوشیدنی خانوار قرار گرفته در دهک دهم

در ادامه به بررسی متغیرهای هزینه خانوار در زمینه های مختلف را بررسی می کنیم.

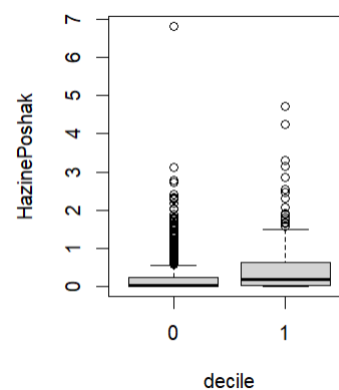
در نمودار ۵۲ ، با استفاده از نمودار پهلوی به پهلوی میزان هزینه ها در قسمت خوراکی خانوار را مشاهده می نمایید.



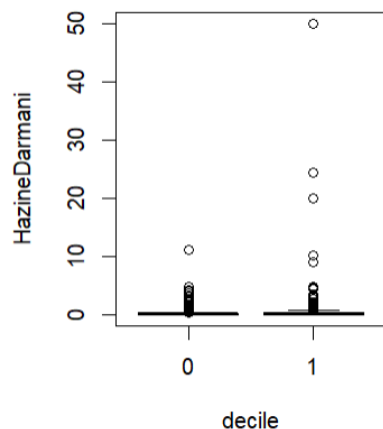
نمودار ۵۲: هزینه خوراکی خانوار قرار گرفته در دهک دهم



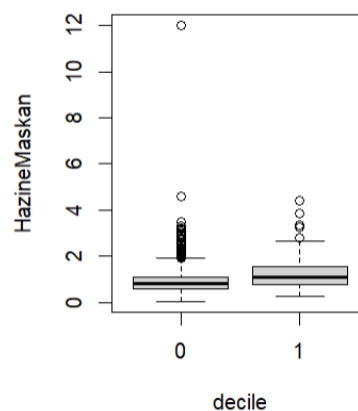
نمودار ۵۶: هزینه لوازم خانگی خانوار قرار گرفته در دهک دهم



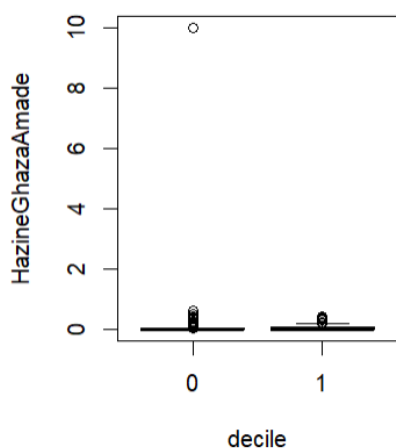
نمودار ۵۴: هزینه پوشاک در خانوارهای قرار گرفته در دهک دهم



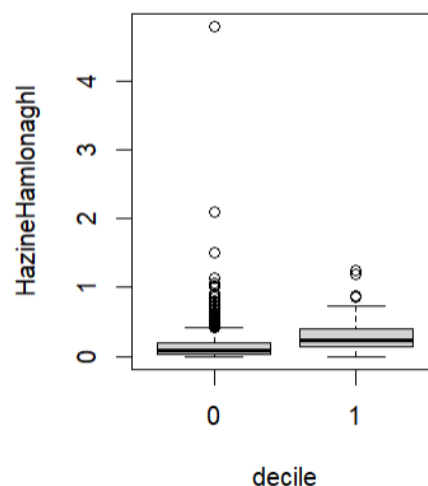
نمودار ۵۷: هزینه درمانی در خانوارهای قرار گرفته در دهک دهم



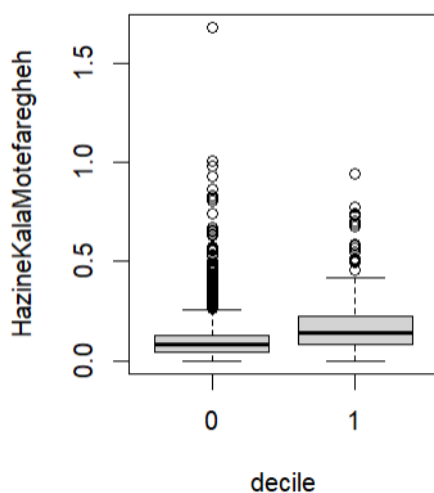
نمودار ۵۵: هزینه مربوط به خانه در خانوارهای قرار گرفته در دهک دهم



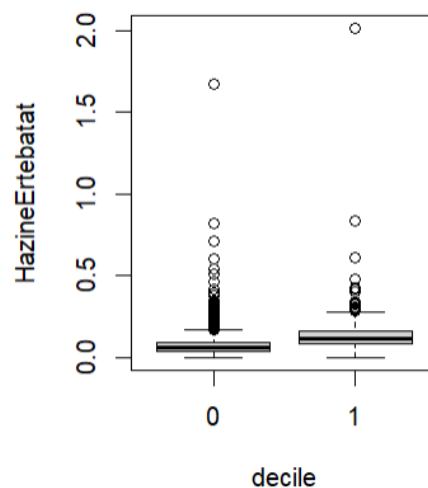
نمودار ۶۱: هزینه غذا آماده و رستوران خانوارهای قرار گرفته در دهک دهم



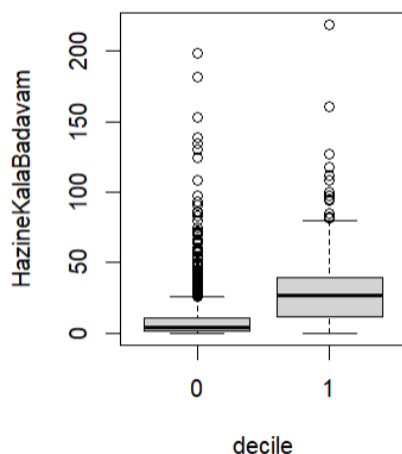
نمودار ۵۸: هزینه حمل و نقل خانوارهای قرار گرفته در دهک دهم



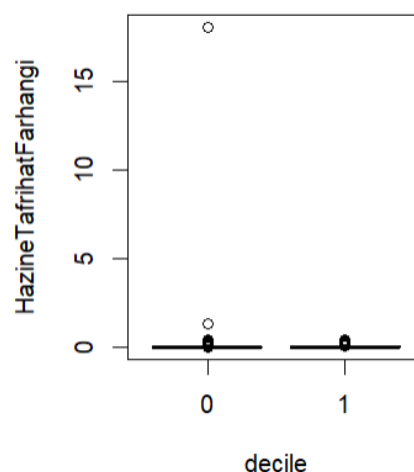
نمودار ۶۲: هزینه کالای متفرقه خانوارهای قرار گرفته در دهک دهم



نمودار ۵۹: هزینه ارتباطات خانوارهای قرار گرفته در دهک دهم



نمودار ۶۳: هزینه کالای بادوام خانوارهای قرار گرفته در دهک دهم در یکسال



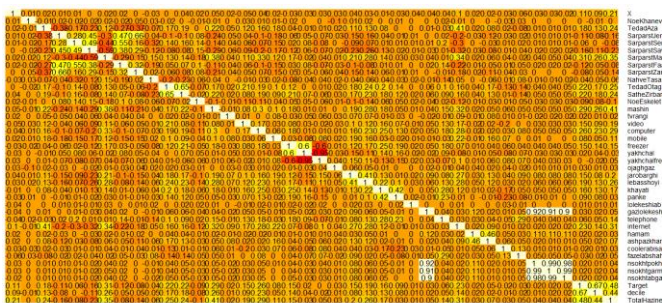
نمودار ۶۰: هزینه تفریحات فرهنگی خانوارهای قرار گرفته در دهک دهم

توضیحات نمودارها به صورت کلی

در نمودارهای ۱۵، ۱۶ و ۱۷ میزان استفاده و علاقه مندی دهک دهم درآمدی نسبت به وسایل نقلیه را مشاهده می کنید. همانطور که مشهود است اکثر خانوارهای قرار گرفته در دهک دهم دارای خودرو می باشند.

$$\begin{aligned} \text{TotalHazine} = & \frac{1}{12} \text{HazineKalaBadavam} \\ & + \text{HazineKalaMotefaregheh} \\ & + \text{HazineGhazaAmade} \\ & + \text{HazineTafrihatFarhangi} \\ & + \text{HazineErtebatat} \\ & + \text{HazineHamlonaghl} \\ & + \text{HazineDarmani} \\ & + \text{HazineLavazemKhanegi} \\ & + \text{HazineMaskan} \\ & + \text{HazinePoshak} \\ & + \text{HazineNoshidani} \\ & + \text{HazineKhoraki} \end{aligned}$$

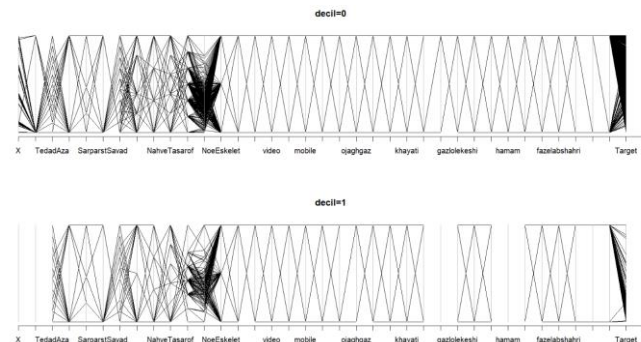
حال نمودار حرارتی این متغیر ها را مشاهده میکنیم.



نمودار ۶۴: نمودار حرارتی همبستگی بین متغیرها مجموعه داده

همانطور که در نمودار ۶۴ مشهود است متغیرهای **freezer**، **yakhchal** و **yakhchalfreezer** همبستگی زیادی مشاهده می شود و همچنین همبستگی بسیار کمی بین **sarparsfaaliat** و **sarparastsen** مشاهده می شود.

پس می توانیم یکی از متغیرهای یخچال یا یخچال فریزر را حذف کنیم. در نهایت ۳۱۰۷ ثبت و ۳۹ متغیر داریم.



نمودار ۶۵: نمودار محور های موازی متغیرها

و از آنجایی که درصد خانوارهای دارای موتورسیکلت و دوچرخه در این دهک بسیار کم است، پس این دو متغیر را حذف میکنیم.

در نمودارهای ۱۸، ۱۹ و ۲۰ نیز همانطور که مشاهده می کنید، متغیر-های رادیو، رادیو ضبط و تلویزیون سیاه سفید نیز معیار مناسبی جهت تشخیص دهک درآمدی نمیباشد. به دلیل آنکه اکثریت خانوارهای قرار گرفته در این دهک این وسایل را ندارند. پس این متغیر ها را نیز می بایست حذف کنیم.

در نمودارهای ۲۲، ۲۳ نیز با توجه به اینکه تفاوت معناداری بین داشتن یا نداشتن دستگاه ویدیو و DVD، کامپیوتر و تبلت در این دهک نیست. پس بنابراین این متغیرها را نمیتوانیم حذف کنیم.

در نمودارهای ۲۱، ۲۴، ۲۸، ۲۹، ۳۰، ۳۷، ۳۸، ۴۰، ۴۱، ۴۲ با توجه به اینکه ۹۹/۹۵ درصد از خانوارهای قرار گرفته در دهک دهم از آن برخوردارند، پس این متغیرها میتواند ما را در تشخیص واقعی این دهک در امدی راهنمایی کند

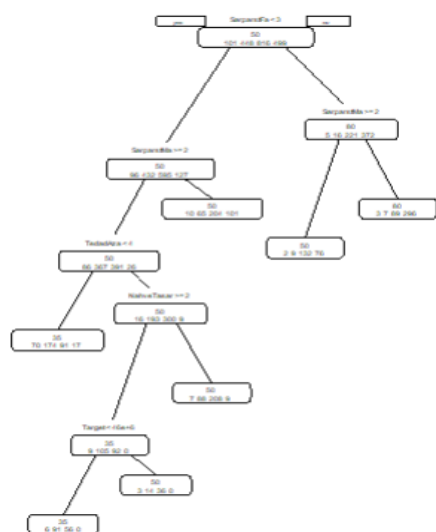
در نمودارهای ۲۵، ۲۶، ۲۷ نیز تفاوتی انچنان بین داشتن و نداشتن وسیله مذکور وجود ندارد پس میبایست حفظ شوند.

در نمودارهای ۳۳، ۳۴، ۳۵، ۳۶، ۴۴، ۴۵، ۴۶، ۴۸ همانطور که مشهود است، این امکانات و تجهیزات در خانوارهای دهک دهم نیز به ندرت دیده می شود، و از آنجایی که هدف پروژه تشخیص دهک مذکور است، پس این متغیر ها هیچ کمکی به ارائه مدل بهتر نمیکند. پس می بایست حذف گردند.

در بررسی نوع سوخت خانوارها در میابیم که ۹۹/۹ درصد از خانوارهای قرار گرفته در این دهک درآمدی از شبکه گاز طبیعی به عنوان منبع گرما در زندگی خود استفاده می کنند. پس می توانیم ان متغیرها را به دو گروه کلی شبکه گاز طبیعی و سایر تقسیم بندی کنیم.

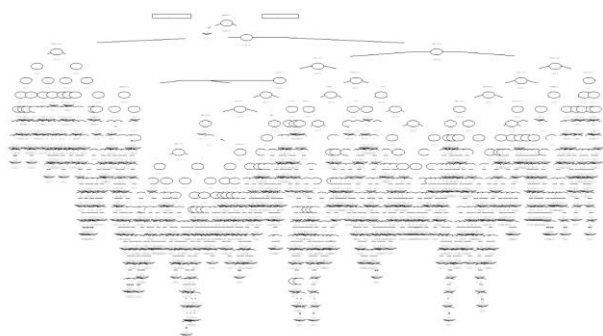
همچنین در نمودارهای ۵۳، ۵۴، ۵۵، ۵۶، ۵۷، ۵۸، ۵۹، ۶۰، ۶۱، ۶۲، ۶۳، که مرتبط با هزینه های خانوار در ماه گذشته می باشد. لذا می-توانیم، تمام هزینه های مربوطه را جمع کرده و در متغیر جدید طبق معادله ۲ جایگزین کنیم.

می‌توانیم درخت رده‌بندی را با توجه به سن سرپرست به صورت زیر مشاهده کنیم.



نمودار ۶۸: درخت رده‌بندی بر اساس سن سرپرست

همچنین با اجرا درخت رده‌بندی عمیق به روی مجموعه آموزشی به شرح زیر داریم.



نمودار ۶۹: درخت رده‌بندی عمیق متغیر سن سرپرست

سپس با اجرا درخت به روی داده‌های اعتبارسنجی و آزمون از صحت عملکرد درخت رده‌بندی مطلع می‌شویم.

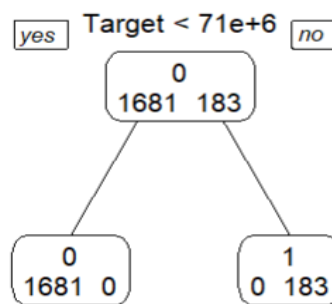
فصل سوم: درخت

درخت رده‌بندی

این فصل یک روش انعطاف پذیر داده رهنمون را تشریح می‌کند که می‌توان برای رده بندی و پیشگویی از آن استفاده کرد. در میان روش‌های داده رهنمون، درختان، شفاف ترین و ساده ترین روش تفسیرند، که بر جداسازی ثبتهای زیرگروهها، از طریق ایجاد بخشهایی در پیشگوییها مبتنی‌اند.

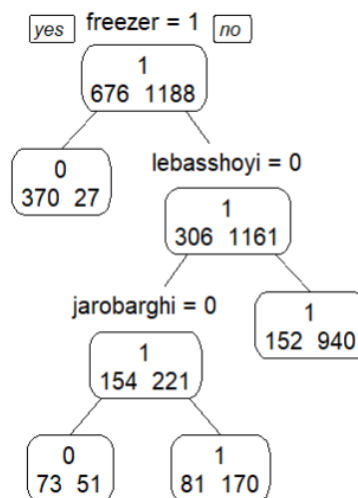
اگر یک فن رده بندی وجود داشته باشد که بدون نیاز به تلاش چندان تحلیلگر، در طیف وسیعی از موقعیتها خوب عمل کند و به راحتی توسط مصرف کننده تحلیل درک پذیر باشد، همانا مدعی قوی روانشناسی درخت است که بری من و همکاران توسعه اش داده اند.

در ابتدا با اجرا درخت به روی مجموعه آموزشی برای رده بندی دهک درآمدی داریم:



نمودار ۶۶: درخت رده‌بندی بر اساس دهک درآمدی دهم

سپس درخت رده‌بندی زیر را به صورت دقیق تر داریم.



نمودار ۶۷: درخت رده‌بندی بر اساس متغیر yakhchalfreezer

```
> confusionMatrix(default.ct.point.pred.train,train.data$SarparstSen)
```

Confusion Matrix and Statistics

	Reference			
Prediction	24	35	50	80
24	0	0	0	0
35	76	265	147	17
50	22	176	580	186
80	3	7	89	296

Overall Statistics

Accuracy : 0.6121
95% CI : (0.5896, 0.6343)
No Information Rate : 0.4378
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.4049

Mcnemar's Test P-Value : < 2.2e-16

Statistics by Class:

	Class: 24	Class: 35	Class: 50	Class: 80
Sensitivity	0.00000	0.5915	0.7108	0.5932
Specificity	1.00000	0.8305	0.6336	0.9275
Pos Pred Value	NaN	0.5248	0.6017	0.7494
Neg Pred Value	0.94582	0.8653	0.7378	0.8618
Prevalence	0.05418	0.2403	0.4378	0.2677
Detection Rate	0.00000	0.1422	0.3112	0.1588
Detection Prevalence	0.00000	0.2709	0.5172	0.2119
Balanced Accuracy	0.50000	0.7110	0.6722	0.7603

خروجی ۳: ارزیابی عملکرد درخت با ماتریس درهم‌ریختگی به روی مجموعه آموزشی

همانطور که مشاهده میکنید، رده‌بندی مدل به روی داده‌های آموزشی

۶۱ درصد صحت دارد، از پیشگویی درخت به روی داده اعتبار سنجی

خروجی زیر حاصل می‌شود.

```
> confusionMatrix(default.ct.point.pred.valid,valid.data$SarparstSen)
```

Confusion Matrix and Statistics

	Reference			
Prediction	24	35	50	80
24	0	0	0	0
35	38	127	82	8
50	7	88	275	91
80	0	2	49	165

Overall Statistics

Accuracy : 0.6084
95% CI : (0.5762, 0.6399)
No Information Rate : 0.4356
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.4022

Mcnemar's Test P-Value : NA

Statistics by Class:

	Class: 24	Class: 35	Class: 50	Class: 80
Sensitivity	0.00000	0.5853	0.6773	0.6250
Specificity	1.00000	0.8210	0.6464	0.9237
Pos Pred Value	NaN	0.4980	0.5965	0.7639
Neg Pred Value	0.95172	0.8671	0.7219	0.8617
Prevalence	0.04828	0.2328	0.4356	0.2833
Detection Rate	0.00000	0.1363	0.2951	0.1770
Detection Prevalence	0.00000	0.2736	0.4946	0.2318
Balanced Accuracy	0.50000	0.7031	0.6619	0.7743

خروجی ۴: ارزیابی درخت به روی مجموعه اعتبار سنجی

همچنین با توجه به اجرای درخت عمیق به روی این متغیر داریم.

همانطور که در خروجی زیر مشاهده میکنید ارزیابی عملکرد به روی مجموعه آموزشی ۱۰۰ درصد است.

```
> confusionMatrix(default.ct.point.pred.train,train.data$decile)
```

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	1681	0
1	0	183

Accuracy : 1
95% CI : (0.998, 1)
No Information Rate : 0.9018
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 1

Mcnemar's Test P-Value : NA

Sensitivity : 1.0000
Specificity : 1.0000
Pos Pred Value : 1.0000
Neg Pred Value : 1.0000
Prevalence : 0.9018
Detection Rate : 0.9018
Detection Prevalence : 0.9018
Balanced Accuracy : 1.0000

'Positive' Class : 0

خروجی ۱: ارزیابی عملکرد درخت با ماتریس درهم‌ریختگی به روی داده آموزشی

همچنین ارزیابی آن به روی داده اعتبار سنجی نیز ۱۰۰ درصد است.

```
> confusionMatrix(default.ct.point.pred.valid,valid.data$decile)
```

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	833	0
1	0	99

Accuracy : 1
95% CI : (0.996, 1)
No Information Rate : 0.8938
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 1

Mcnemar's Test P-Value : NA

Sensitivity : 1.0000
Specificity : 1.0000
Pos Pred Value : 1.0000
Neg Pred Value : 1.0000
Prevalence : 0.8938
Detection Rate : 0.8938
Detection Prevalence : 0.8938
Balanced Accuracy : 1.0000

'Positive' Class : 0

خروجی ۲: ارزیابی عملکرد درخت به روی داده اعتبارسنجی

همچنین دقیقاً همین مقادیر برای درخت عمیق متغیر دهک درآمدی

ثبت شد.

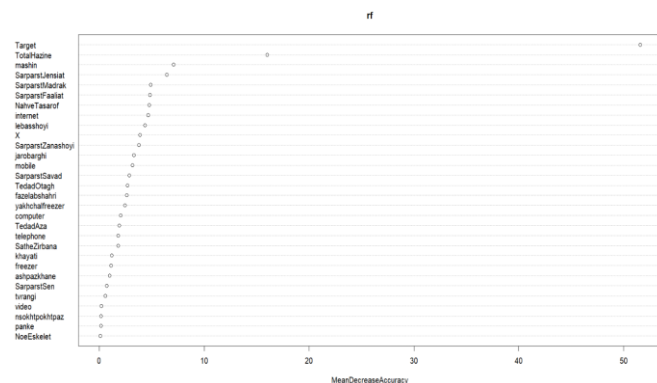
در ارزیابی عملکرد درخت رده‌بندی متغیر سن سرپرست به روی

مجموعه آموزشی به شرح زیر داریم.

جنگل تصادفی

یکی دیگر از روش‌های رده‌بندی درخت‌های تصادفی است. بر خلاف یک درخت منفرد، نتایج یک جنگل تصادفی را نمی‌توان در یک نمودار درخت گونه نمایش داد، بدین ترتیب قابلیت تفسیری که یک درخت منفرد ارائه می‌کند، از دست می‌رود. در عین حال، جنگل‌های تصادفی می‌توانند امتیازهای اهمیت متغیر را که مشارکت نسبی پیشگوهایی مختلف را اندازه می‌گیرد، تولید کنند.

در نمودار زیر جنگل تصادفی اجرا شده به روی مجموعه آموزشی را مشاهده می‌نمایید. همچنین ارزیابی عملکرد این جنگل تصادفی را نیز به شرح خروجی زیر داریم.



نمودار ۷۰: جنگل تصادفی اجرا شده به روی مجموعه آموزشی و اهمیت متغیرها

```
> rf.pred<-predict(rf,valid.data)
> confusionMatrix(rf.pred,valid.data$decile)
Confusion Matrix and Statistics
```

	Reference	0	1
Prediction	0	833	4
	1	0	95

Accuracy : 0.9957
95% CI : (0.989, 0.9988)
No Information Rate : 0.8938
P-Value [Acc > NIR] : <2e-16

Kappa : 0.977

Mcnemar's Test P-Value : 0.1336

Sensitivity : 1.0000
Specificity : 0.9596
Pos Pred Value : 0.9952
Neg Pred Value : 1.0000
Prevalence : 0.8938
Detection Rate : 0.8938
Detection Prevalence : 0.8981
Balanced Accuracy : 0.9798

'Positive' Class : 0

خروجی ۷: ارزیابی عملکرد جنگل تصادفی به روی مجموعه اعتبار سنجی

همانطور که در نمودار ۷۰ مشهود است متغیرهای **target** و **totalAzine** خانوار به شدت در دهک بندی تاثیر گذارند، و همچنین عملکرد جنگل تصادفی در دهک‌بندی مجموعه اعتبارسنجی نیز با دقت ۰.۹۹۵۷ درصد ثبت شده است.

```
> confusionMatrix(deeper.ct.point.pred.train,train.data$SarparsSen)
Confusion Matrix and Statistics
```

	Reference	24	35	50	80
Prediction	24	101	0	0	0
	35	0	448	0	0
	50	0	0	816	0
	80	0	0	0	499

Overall Statistics

Accuracy : 1
95% CI : (0.998, 1)
No Information Rate : 0.4378
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 1

Mcnemar's Test P-Value : NA

Statistics by Class:

	Class: 24	Class: 35	Class: 50	Class: 80
Sensitivity	1.00000	1.0000	1.0000	1.0000
Specificity	1.00000	1.0000	1.0000	1.0000
Pos Pred Value	1.00000	1.0000	1.0000	1.0000
Neg Pred Value	1.00000	1.0000	1.0000	1.0000
Prevalence	0.05418	0.2403	0.4378	0.2677
Detection Rate	0.05418	0.2403	0.4378	0.2677
Detection Prevalence	0.05418	0.2403	0.4378	0.2677
Balanced Accuracy	1.00000	1.0000	1.0000	1.0000

خروجی ۵: ارزیابی عملکرد درخت عمیق به روی مجموعه آموزشی با ماتریس درهم‌ریختگی همانطور که مشاهده میکنید میزان یادگیری درخت عمیق به روی مجموعه داده آموزشی ۱۰۰ درصد است همچنین میزان دقت آن برای رده‌بندی به روی مجموعه آموزشی ۱ است. ولی با توجه به خروجی ۶، این دقت به عدد ۵۳ درصد رسیده است و به طور جدی میزان دقت رده‌بندی به روی مجموعه اعتبار سنجی با مشکل بیش برآزش رو به رو شده است.

یک خطر رشد دادن درخت‌های عمیق‌تر روی داده‌های آموزشی، بیش برآزش است. بیش برآزش به عملکرد ضعیف در داده‌های جدید منجر خواهد شد.^۱

```
> deeper.ct.point.pred.valid<-predict(deeper.ct,valid.data,type="class")
> confusionMatrix(deeper.ct.point.pred.valid,valid.data$SarparsSen)
Confusion Matrix and Statistics
```

	Reference	24	35	50	80
Prediction	24	8	30	11	2
	35	28	105	94	11
	50	8	75	212	73
	80	1	7	89	178

Overall Statistics

Accuracy : 0.5397
95% CI : (0.5071, 0.5721)
No Information Rate : 0.4356
P-Value [Acc > NIR] : 1.136e-10

Kappa : 0.3254

Mcnemar's Test P-Value : 0.4837

Statistics by Class:

	Class: 24	Class: 35	Class: 50	Class: 80
Sensitivity	0.177778	0.4839	0.5222	0.6742
Specificity	0.951522	0.8140	0.7034	0.8548
Pos Pred Value	0.156863	0.4412	0.5761	0.6473
Neg Pred Value	0.958002	0.8386	0.6560	0.8691
Prevalence	0.048283	0.2328	0.4356	0.2833
Detection Rate	0.008584	0.1127	0.2275	0.1910
Detection Prevalence	0.054721	0.2554	0.3948	0.2951
Balanced Accuracy	0.564650	0.6489	0.6128	0.7645

خروجی ۶: ارزیابی عملکرد درخت عمیق به روی مجموعه اعتبار سنجی

پس بنابراین برای دسترسی به مدلی بهتر روش‌های دیگر را بررسی می‌کنیم.

همانطور که مشهود است، معادله زیر حاصل می‌شود.

فصل چهارم: رگرسیون لوژیستیک

در این فصل روش رده‌بندی بسیار قدرتمند و رایجی را که رگرسیون لوژیستیک نامیده می‌شود را اجرا می‌کنیم. رگرسیون لوژیستیک، ایده رگرسیون خطی را برای حالتی که متغیر برآمد، y ، رسته‌ای است، توسعه می‌دهد. در خروجی زیر برازش یک مدل رگرسیون لوژیستیک را مشاهده می‌کنید.

```
Call:
glm(formula = decile ~ ., family = "binomial", data = train.data)
```

```
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-0.00051881 -0.00000002 -0.00000002 -0.00000002  0.00064826
```

```
Coefficients: (1 not defined because of singularities)
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  758.788247872 884130.2108969315  0.001
X             -0.0000992591    0.2929280706  0.000
NoeKhanevar   429.9595118529 234256.0990300766  0.002
TedadAza      -2.9082348184   1266.0892132882 -0.002
SarparstJensiat -11.1710913315  73103.3810471379  0.000
SarparstSen35  -16.9275916730  22872.6623968570 -0.001
SarparstSen50  -34.8997947540  22399.9050100665 -0.002
SarparstSen80  -40.3812627503  22524.6154149806 -0.002
SarparstSavad  -21.9862789346   6834.1358264012 -0.003
SarparstMadrak  -3.3360699503   1017.1716430431 -0.003
SarparstFaaliat -2.0751053683   1830.2599089943 -0.001
SarparstZanashoyi 19.6671128711   72101.6081296941  0.000
NahveTasarof    1.4216216481   1232.7154499149  0.001
TedadOtagh     -9.6517822722   3224.5497703946 -0.003
SatheZirbana    0.1322979777     64.7661783808  0.002
NoeEskelet     -1.2243954255   1637.2034759108 -0.001
mashin        -28.4706560342   4654.6820093825 -0.006
tvrangi        7.8608282168   5049.1886963285  0.002
video         19.3149174731   4214.9497308810  0.005
computer       13.1255480291   2599.7357716833  0.005
mobile        -26.0998125729   29653.2462874807 -0.001
freezer        6.4487722675    7994.8573918092  0.001
yakhchalfreezer 9.7279607335    6417.3876831679  0.002
ojaghgaz     -58.5566938175   60083.9650348164 -0.001
jarobarghi   -106.4682040714   18536.8798478428 -0.006
lebasshoyi    -3.7277233736    6414.9289758174 -0.001
khayati       -1.9343738770    3767.9213785165 -0.001
panke        -28.9811450592   5658.0973016281 -0.005
lolekeshiab  -304.4675853024  163525.2626220236 -0.002
gazlolekeshi  154.2812865702   149836.9266871985  0.001
telephone     8.6423784205    3473.8590844062  0.002
internet      3.4373887329    7738.7511258033  0.000
```

```
SarparstSen80      0.999
SarparstSavad      0.997
SarparstMadrak     0.997
SarparstFaaliat    0.999
SarparstZanashoyi  1.000
NahveTasarof       0.999
TedadOtagh         0.998
SatheZirbana       0.998
NoeEskelet         0.999
mashin             0.995
tvrangi            0.999
video              0.996
computer           0.996
mobile             0.999
freezer            0.999
yakhchalfreezer    0.999
ojaghgaz           0.999
jarobarghi         0.995
lebasshoyi         1.000
khayati            1.000
panke              0.996
lolekeshiab        0.999
gazlolekeshi       0.999
telephone          0.998
internet           1.000
hamam              1.000
ashpazkhane        0.999
coolerabisabet     1.000
fazelabshahri      1.000
nsokhtpokhtpaz     1.000
nsokhtgarma        NA
nsokhtbgarm        0.999
Target             0.976
TotalHazine        0.996
```

(Dispersion parameter for binomial family taken to be 1)

```
Null deviance: 1196.8989138156 on 1863 degrees of freedom
Residual deviance: 0.0000026023 on 1824 degrees of freedom
AIC: 80
```

Number of Fisher Scoring iterations: 25

خروجی ۸: برازش مدل رگرسیون لوژیستیک به روی مجموعه آموزشی

Logit(decile = 1)

```
= 758.78 + 429.95 NoeKhanevar
- 2.90 TedadAZA
- 11.17 SarparstJensiat
- 16.92 SarparstSen35
- 34.89 SarparstSen50
- 40.38 SarparstSen80
- 21.98 SarparstSavad
- 3.33 SarparstMadrak
- 2.07 SarparstFaaliat
+ 19.66 SarparstZanashoyi
+ 1.42 NahveTasarof
- 9.65 TedadOtagh
+ 0.13 SatheZirbana
- 1.22 NoeEskelet
- 28.47 mashin + 7.86 tvrangi
+ 19.31 video + 13.12 computer
- 26.09 mobile + 6.44 freezer
+ 9.72 yakhchalfreezer
- 58.55 ojaghgaz
- 106.46 jarobarghi
- 3.72 lebasshoyi - 1.93 khayati
- 28.98 panke
- 304.46 lolekeshiab
+ 154.28 gazlolekeshi
+ 8.64 telephone + 3.43 internet
+ 150.88 hamam
- 221.26 ashpazkhane
+ 3.63 coolerabisabet
+ 0.58 fazelabshahri
+ 200.72 nsokhtpokhtpaz
+ 50.80 nsokhtgarma
- 98.5 nsokhtbgarm
+ 0.000012 Target
+ 0.00000029 TotalHazine
```

همانطور که در خروجی ۸ مشاهده می‌نمایید **AIC=80** می‌باشد. و

Fisher Scoring iterations به تعداد ۲۵ عدد در این مجموعه می‌باشد.

ارزیابی عملکرد رده بندی

رایج ترین روش ارزیابی عملکرد، استفاده از ماتریس در هم ریختگی است. در خروجی زیر نسبت به ارزیابی عملکرد مدل می پردازیم.

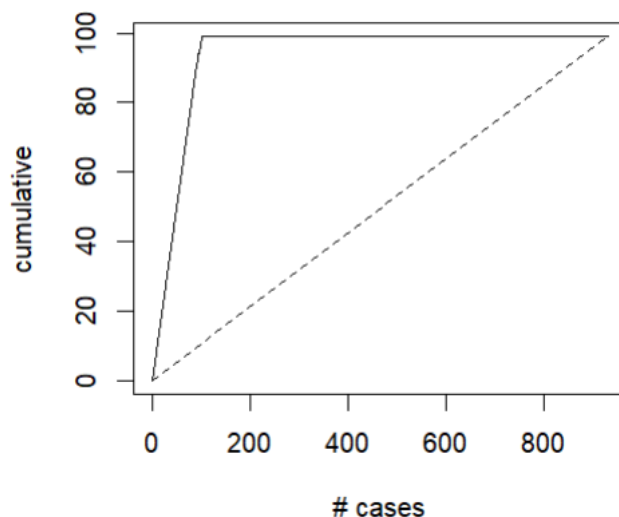
```
> data.frame(actual=valid.data$decile[1:10], predicted=logit.reg.pred[1:10])
```

	actual	predicted
2168	0	0.0000000000000000220446
1248	0	0.0000000000000000220446
1353	0	0.0000000000000000220446
1393	0	0.0000000000000000220446
1367	1	0.9999999999999999779554
1925	0	0.0000000000000000220446
2070	1	0.9999999999999999779554
2699	0	0.0000000000000000220446
24	0	0.0000000000000000220446
541	0	0.0000000000000000220446

خروجی ۹: ارزیابی عملکرد مدل رگرسیون لوژیستیک

همانطور که مشاهده می کنید به نتایج مطلوبی در ارزیابی مدل به روی داده های اعتبارسنجی رسیدیم.

دیگر ابزار مفید برای ارزیابی عملکرد رده بندی مدل، نمودارهای بالابردن و بالا بردن دهکی هستند. در زیر نمودار بالا بردن به دست آمده برای مدل لوژیستیک روی مجموعه اعتبارسنجی را مشاهده می کنید.



نمودار ۷۱: نمودار بالابردن برای ارزیابی مدل لوژیستیک

همانطور که در نمودار ۷۱ نیز مشهود است هرچه فاصله بالا بردن از خط قطری میان صفحه بیشتر باشد، نشان دهنده عملکرد خوب مدل است. پس نتیجه می گیریم که مدل رگرسیون لوژیستیک بسیار عملکرد خوبی در این بخش دارد.

فصل پنجم: k-نزدیک ترین همسایه (k-NN)

در این فصل الگوریتم k-نزدیک ترین همسایه را که می تواند در رده بندی (برای یک برآمد رسته ای) یا پیشگویی (برای یک برآمد عددی) استفاده شود، توصیف می کنیم.

الگوریتم k-نزدیک ترین همسایه، یک روش رده بندی است که در آن، هیچ فرضی درباره شکل رابطه بین عضویت رده Y و پیشگوهای X_1, X_2, \dots, X_p وجود ندارد. پس این روش ناپارامتری است.

حال به پیاده سازی این روش به روی مجموعه آموزشی می پردازیم.

```
#K-NN
#initialize normalized training,valid,complete data frames to originals
#use preprocess() from the caret package to normalize decile and samparstsen
norm.values<-preProcess(train.data[,1:2],method=c("center","scale"))
train.data[,1:2]<-predict(norm.values,train.data[,1:2])
valid.data[,1:2]<-predict(norm.values,valid.data[,1:2])
data.df[,1:2]<-predict(norm.values,data.df[,1:2])
#initialize a dataframe with two columns: k, and accuracy.
new.norm.df<-predict(norm.values,valid.data[,1:2])
library(class)
library(FNN)
#use knn()
nn<-knn(train=train.data[,1:2],test=new.norm.df,cl=train.data[,3],k=3)
row.names((train.data)[attr(nn,"nn.index")])
accuracy.df<-data.frame(k=seq(1,14,1),accuracy=rep(0,14))
#compute knn for different k on validation
for(i in 1:14){knn.pred<-knn(train.data[,1:2],valid.data[,1:2],
                             cl=train.data[, 3],k=i)
accuracy.df[i,2]<-confusionMatrix(knn.pred,as.factor(valid.data[, 3]))$overall[1]}
```

همانطور که در تصویر فوق مشاهده می کنید، الگوریتم k-نزدیک ترین همسایه به شرح فوق پیاده سازی می گردد. تا بهترین k را بیابد. خروجی زیر را داریم.

```
> #use knn()
> nn<-knn(train=train.data[,1:2],test=new.norm.df,cl=train.data[,3],k=3)
> row.names(train.data)[attr(nn,"nn.index")]
[1] "10" "18" "18" "32" "56" "83" "90" "97"
[9] "122" "133" "144" "144" "180" "184" "191" "191"
[17] "194" "196" "200" "203" "247" "259" "272" "286"
[25] "298" "300" "304" "313" "323" "324" "329" "330"
> nn
[1] 4 3 3 1 3 4 5 1 4 2 2 2 2 1 3 3 3 2 2 2 1 1 4 2 2 5 5 3
[29] 4 4 4 2 2 4 3 4 3 3 3 3 2 2 2 2 3 4 2 3 2 5 5 2 2 2 3 3
[57] 3 4 6 2 3 3 3 3 4 4 1 1 4 4 2 2 2 2 3 2 3 4 2 3 5 2 4 3
[85] 3 3 3 4 4 5 3 4 2 2 2 3 2 5 3 2 1 5 2 4 1 1 1 3 2 4 2 3
attr(,"nn.index")
[1,] [1,] [2,] [3,]
[1,] 816 339 29
[2,] 96 29 618
[3,] 96 618 1079
[4,] 807 732 1217
[1,] [2,] [3,]
[1,] 0.0001859634 0.0001859634 5.578901e-04
[2,] 0.0003719267 0.0003719267 5.578901e-04
[3,] 0.0001859634 0.0003719267 5.578901e-04
[4,] 0.0001859634 0.0001859634 3.719267e-04
[5,] 0.0003719267 0.0005578901 5.578901e-04
[6,] 0.0001859634 0.0003719267 5.578901e-04
```

خروجی ۱۰: روش پیدا کردن بهترین k در الگوریتم k-NN

در نهایت به خروجی زیر دست میابیم.

k	accuracy
1	0.6
2	0.6
3	0.7
4	0.7
5	0.8
6	0.9
7	0.9
8	0.9
9	1.0
10	0.9
11	0.8
12	0.8
13	0.4
14	0.5

خروجی ۱۱: بهترین K همسایه در الگوریتم K-NN

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	1615	111
1	66	72

Accuracy : 0.905
95% CI : (0.8908, 0.918)

No Information Rate : 0.9018
P-Value [Acc > NIR] : 0.3374011

Kappa : 0.3978

McNemar's Test P-Value : 0.0009422

Sensitivity : 0.9607
Specificity : 0.3934
Pos Pred Value : 0.9357
Neg Pred Value : 0.5217
Prevalence : 0.9018
Detection Rate : 0.8664
Detection Prevalence : 0.9260
Balanced Accuracy : 0.6771

'Positive' Class : 0

خروجی ۱۲: ارزیابی عملکرد شبکه عصبی به روی مجموعه آموزشی

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	794	59
1	39	40

Accuracy : 0.8948
95% CI : (0.8734, 0.9138)

No Information Rate : 0.8938
P-Value [Acc > NIR] : 0.48436

Kappa : 0.3921

McNemar's Test P-Value : 0.05495

Sensitivity : 0.9532
Specificity : 0.4040
Pos Pred Value : 0.9308
Neg Pred Value : 0.5063
Prevalence : 0.8938
Detection Rate : 0.8519
Detection Prevalence : 0.9152
Balanced Accuracy : 0.6786

'Positive' Class : 0

خروجی ۱۳: ارزیابی عملکرد مدل شبکه عصبی به روی مجموعه اعتبار سنجی

همانطور که مشهود است عملکرد بی نظیر شبکه عصبی با ۲ لایه پنهان با ۱۰ و ۵ نورون به گونه‌ای شگفت آوری توجه ما را جلب می‌کند. ولی مشکلی نیز وجود دارد. و آن میزان **Sensitivity** و میزان **Specificity** است. ولی در حالی است که میزان دقت آن در دو مجموعه آموزشی و اعتبار سنجی مطلوب است. ولی با توجه به هدف پروژه این دو مقدار قابل قبول است.

همانطور که مشاهده می‌کنید، بهترین $K=5$ است به دلیل آنکه اگر $K=9$ انتخاب شود، مطمئناً با مشکل بیش برآزش روبه رو خواهیم شد. لذا بهترین و سریع‌ترین را انتخاب می‌کنیم. سپس بار دیگر الگوریتم را با $K=5$ اجرا می‌کنیم.

مزیت عمده روش K -NN سادگی آن و نداشتن فروض پارامتری است. در صورت وجود مجموعه آموزشی به اندازه کافی بزرگ، این روش‌ها به طرز شگفت‌آوری خوب عمل می‌کنند. زمان یافتن نزدیک‌ترین همسایه‌ها در یک مجموعه آموزشی بزرگ می‌تواند کمرشکن باشد، این یکی از معایب این روش است، از راهکار-های حل این مشکل می‌توان به کاهش بعد اشاره کرد.

فصل ششم: شبکه عصبی

در این فصل، شبکه‌های عصبی را توصیف می‌کنیم، روشی داده‌رهنمون و انعطاف‌پذیر که می‌تواند برای رده‌بندی و پیشگویی به کار گرفته شود. شبکه‌های عصبی اگرچه از نقطه نظر تفسیرپذیری، جعبه سیاه تلقی می‌شود، از نظر درستی پیشگوییانه قویا موفق است. شبکه عصبی بر یک مدل فعالیت بیولوژیکی در مغز مبتنی است که در آن سلول‌های عصبی در اتصال با یکدیگرند و از تجربه یاد می‌گیرند. تابع فعال ساز لوژستیک مورد استفاده قرار می‌گیرد. تابع لوژستیک رایج‌ترین تابع در شبکه‌های عصبی است. ارزش عملی آن از این واقعیت بر می‌خیزد که روی مقادیر بسیار کوچک یا مقادیر بسیار بزرگ اثر تخریبی دارد، اما در دامنه‌ای که مقادیر تابع بین ۰،۹ و ۰،۱ باشد، تقریباً خطی است. پس در نتیجه می‌توانیم خروجی راس j در لایه پنهان را به صورت زیر بنویسیم.

$$Output_j = g\left(\theta_j + \sum_{i=1}^p w_{ij}x_i\right) = \frac{1}{1 + e^{-(\theta_j + \sum_{i=1}^p w_{ij}x_i)}}$$

θ اربیبی راس j را نشان می‌دهد.

همانطور که مشاهده می‌کنید، این معادله فرمول رگرسیون لوژستیک است!!

حال به اجرای مدل شبکه عصبی می‌پردازیم.

برای درک بهتر ساختار شبکه عصبی اجرا شده، می‌توانید به شکل زیر توجه کنید.

نتیجه گیری کلی

همانطور که مشاهده نمودید بهترین عملکرد مربوط به شبکه‌های عصبی است پس در نهایت این مدل را به روی مجموعه آزمون با تابع فعال ساز لوژستیک برازش می‌دهیم و می‌توانیم در خروجی زیر نتیجه آن را مشاهده کنیم.

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	271	19
1	12	9

Accuracy : 0.9003
95% CI : (0.8615, 0.9313)
No Information Rate : 0.91
P-Value [Acc > NIR] : 0.7604

Kappa : 0.3144

McNemar's Test P-Value : 0.2812

Sensitivity : 0.9576
Specificity : 0.3214
Pos Pred Value : 0.9345
Neg Pred Value : 0.4286
Prevalence : 0.9100
Detection Rate : 0.8714
Detection Prevalence : 0.9325
Balanced Accuracy : 0.6395

'Positive' Class : 0

خروجی ۱۵: ارزیابی عملکرد مدل نهایی به روی مجموعه آزمون با تابع فعال ساز لوژستیک

همانطور که مشاهده می‌کنید تعادل دقت مدل ۰,۹۰ رسیده است.

و این یعنی می‌توانیم با استفاده از مدل شبکه‌های عصبی با تابع فعال-ساز لوژستیک به بهترین نتیجه در خصوص داده‌های جدید برسیم.^۱

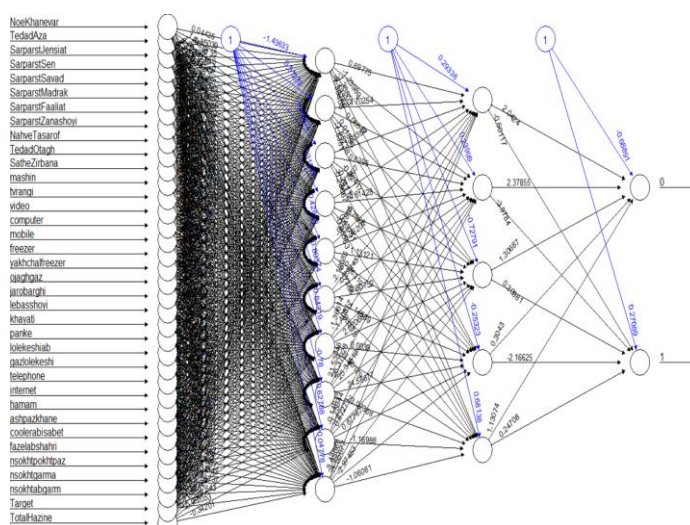
سخن پایانی

در پایان ضمن تقدیر و تشکر از زحمات استاد گرانقدر درس داده‌کاوی جناب آقای دکتر محمدرضا فقیهی حبیب آبادی و همچنین انجمن ریاضی دانشگاه شهید بهشتی، از تمامی خوانندگان کمال تشکر را اعلام می‌دارم.

تهیه‌کننده: سجاد صحت بخش

دانشجوی کارشناسی ارشد دانشگاه شهید بهشتی

رشته ریاضی کاربردی-گرایش علوم داده



نمودار ۷۲: ساختار شبکه عصبی اجرا شده به روی مجموعه آموزشی

همچنین برای اجرای این شبکه عصبی بدین نحو عمل کردیم.

```
#####
library(neuralnet)
Nn<-neuralnet(as.factor(decile)~NoeKhanevar+TedadAza+SarparsTensiat+SarparsSen+
SarparsSavad+SarparsMadrak+SarparsFaaliat+SarparsZanashoyi+
NahveTasarof+TedadOtagh+SatheZirbana+mashin+tvrang+video+
computer+mobile+freezer+yakhchal+freezer+ojaghgaz+jarobarghi+
lebasshoyi+khayati+panke+lolekeshiab+gazlolekeshi+telephone+
internet+hamam+ashpazkhane+coolerabisabet+fazelabshahri+
nsokhtpokhtpaz+nsokhtgarma+nsokhtabgarm+Target+TotalHazine,
data = train.data,linear.output = F,hidden =c(10,5))

plot(Nn)
train.data<-select(train.data,-decile)
training.prediction=compute(Nn,train.data[,1:38])
training.class=apply(training.prediction$net.result,1,which.max)-1
confusionMatrix(as.factor(training.class),as.factor(train.data$decile))

valid.prediction=compute(Nn,valid.data[,1:39])
valid.class=apply(valid.prediction$net.result,1,which.max)-1
confusionMatrix(as.factor(valid.class),as.factor(valid.data$decile))
```

خروجی ۱۴: اجرای شبکه عصبی به روی مجموعه آموزشی و اعتبار سنجی

برجسته‌ترین مزیت شبکه‌های عصبی، عملکرد پیشگویانه خوب آنهاست. آنها به تحمل بالا برای داده‌های نوفه‌ای و توانایی شکار کردن ارتباط‌های بسیار پیچیده بین پیشگوها و یک برآمد مشهورند. مهمترین نقطه ضعف آن‌ها به فراهم آوردن بینش درباره ساختار ارتباط و در نتیجه سوء شهرت به جعبه سیاه مربوط است.