



2022

# پروژه درس مدلسازی و تصمیم‌گیری داده محور

استاد درس: دکتر نفیسه صدقی

معصومه محمودی - محمد زارعی

دانشگاه صنعتی شریف | نیم‌سال دوم ۱۴۰۰ - ۱۴۰۱

## فهرست

ح	چکیده
۱	۱ توضیح مساله و داده‌ها
۱-۱	۱-۱ وام خودرو
۲-۱	۲-۱ معرفی دیتاست (Dataset)
۳-۱	۳-۱ Data Dictionary
۲	۲ تحلیل کاوش گرایانه داده (EDA)
۱-۲	۱-۲ بررسی آماری دیتاست
۲-۲	۲-۲ تحلیل تک متغیره (Univariate Analysis)
۳	۳ آماده‌سازی داده (Data Preparation)
۱-۳	۱-۳ داده‌های پرت (Outliers)
۲-۳	۲-۳ داده‌های گم‌شده یا ناموجود (NA - Missing)
۳-۳	۳-۳ مقداردهی (Imputaion)
۱-۳-۳	۱-۳-۳ Predictive Mean Matching (PMM)
۲-۳-۳	۲-۳-۳ Polytomous (Multinomial) Logistic Regression
۴	۴ مصورسازی داده (Data Visualization)
۱-۴	۱-۴ رابطه میان Default و Client_Income_Type
۲-۴	۲-۴ رابطه میان Default و Client_Education
۳-۴	۳-۴ رابطه میان Default و Gender
۴-۴	۴-۴ رابطه میان Default و Loan_Contract_Type
۵-۴	۵-۴ رابطه میان Default و Client_Housing_Type
۶-۴	۶-۴ رابطه میان Default و Age_Days

۱۹	.....Employed_Days و Default	۷-۴
۲۰	.....Registration_Days و Default	۸-۴
۲۰	.....ID_Days و Default	۹-۴
۲۱	.....Own_House_Age و Default	۱۰-۴
۲۱	.....Client_City_Rating و Default	۱۱-۴
۲۲	.....Client_Permanent_Match_Tag و Default	۱۲-۴
۲۲	.....Score_Source_1 و Default	۱۳-۴
۲۳	.....Score_Source_2 و Default	۱۴-۴
۲۳	.....Score_Source_3 و Default	۱۵-۴
۲۴	.....Social_Circle_Default و Default	۱۶-۴
۲۴	.....Phone_Change و Default	۱۷-۴
۲۵	.....(Descriptive Statistics)	۵ آمار توصیفی
۲۵	.....Client_Education و Default	۱-۵
۲۵	.....Client_Gender و Default	۲-۵
۲۶	.....آزمون میانگین درآمد	۳-۵
۲۶	.....بازه اطمینان درآمد	۴-۵
۲۷	.....Client_Housing_Type و Default	۵-۵
۲۷	.....آزمون سن	۶-۵
۲۸	.....Client_City_Rating و Default	۷-۵
۳۰	.....(Predictive Models)	۶ مدل های پیش بینی
۳۰	.....Logistic Regression	۱-۶
۳۲	.....Decision Tree	۲-۶
۳۳	.....Random Forest	۳-۶

۷ نتیجه گیری..... ۳۶

۸ مراجع..... ۳۷

## فهرست اشکال

شکل ۱-۲ نمودار Bar Chart برای متغیر Default	۹
شکل ۱-۳ تعداد مقادیر NA متغیرها	۱۲
شکل ۱-۴ رابطه میان Default و Client_Income_Type	۱۶
شکل ۲-۴ رابطه میان Default و Client_Education	۱۷
شکل ۳-۴ رابطه میان Default و Gender	۱۷
شکل ۲-۴ رابطه میان Default و Loan_Contract_Type	۱۸
شکل ۵-۴ رابطه میان Default و Client_Housing_Type	۱۸
شکل ۶-۴ رابطه میان Default و Age_Days	۱۹
شکل ۷-۴ رابطه میان Default و Employed_Days	۱۹
شکل ۸-۴ رابطه میان Default و Registration_Days	۲۰
شکل ۹-۴ رابطه میان Default و ID_Days	۲۰
شکل ۱۰-۴ رابطه میان Default و Own_House_Age	۲۱
شکل ۱۱-۴ رابطه میان Default و Client_City_Rating	۲۱
شکل ۱۲-۴ رابطه میان Default و Client_Permanent_Match_Tag	۲۲
شکل ۱۳-۴ رابطه میان Default و Score_Source_1	۲۲
شکل ۱۴-۴ رابطه میان Default و Score_Source_2	۲۳
شکل ۱۵-۴ رابطه میان Default و Score_Source_3	۲۳
شکل ۱۶-۴ رابطه میان Default و Social_Circle_Default	۲۴
شکل ۱۷-۴ رابطه میان Default و Phone_Change	۲۴
شکل ۱-۶ نمودار ROC مدل Logistic Regression برای داده‌های train	۳۱
شکل ۲-۶ نمودار ROC مدل Logistic Regression برای داده‌های test	۳۱
شکل ۳-۶ نمودار ROC مدل Decision Tree برای داده‌های train	۳۲
شکل ۴-۶ نمودار ROC مدل Decision Tree برای داده‌های test	۳۳
شکل ۵-۶ نمودار اهمیت متغیرها در مدل Random Forest	۳۴
شکل ۶-۶ نمودار ROC مدل Random Forest برای داده‌های train	۳۵

شکل ۶-۷ نمودار ROC مدل Random Forest برای داده‌های test ..... ۳۵

شکل ۷-۱ منحنی ROC مدل‌های ارائه شده ..... ۳۶

## چکیده

در ابتدا بر روی داده های این دیتاست عملیات مورد پیش پردازش (Preprocessing) قرار گرفت تا داده هایی که در فرایند بررسی و مصور سازی و پیش بینی و نتیجه گیری کاربرد زیادی ندارند و بیشتر باعث اختلال در عملکرد می شوند تا حد امکان کنار گذاشته شوند. سپس با استفاده از نمودارهای مختلف ارتباط بین متغیرها و تاثیری که بر روی یکدیگر دارند را بررسی شد. با استفاده از نمودارها و Visualization داده ها اطلاعات کلی از مدل ارتباط متغیرها با یکدیگر بدست می آوریم. در ادامه با توجه به این که متغیر پاسخ از نوع دسته ای (Categorical) می باشد، چند مدل طبقه بندی (Classification) ایجاد شده و از میان آن ها یک مدل که عملکرد بهتری داشته برای پیش بینی متغیر پاسخ در دیتاست test استفاده شده است.

واژگان کلیدی: سرمایه گذاری، خودرو، مبتدی، بانکداری، رگرسیون لجستیک

## ۱ توضیح مساله و داده ها

### ۱-۱ وام خودرو

خودروها رایج ترین دارایی های غیرمالی در میان افراد هستند. تقریباً سه چهارم خریدهای خودرو از طریق اعتبار تأمین می شود و وام های خرید خودرو یکی از رایج ترین اشکال استقراض خانوارها است. وام دهندگان در بازار خودرو با خطراتی مواجه هستند. اولین و آشکارترین ریسک، عدم بازپرداخت وام است. یعنی شخصی که برای خرید خودرو وام گرفته است و آن را پس نمی دهد. دومین ریسک مهم برای وام دهندگان در این بازار، ریسک پیش پرداخت است. یعنی خریدار خودرو وام را زودتر پرداخت می کند و جریان پرداخت بهره وام دهنده را کاهش می دهد.

مؤسسه مالی غیر بانکی (NBFI)<sup>۱</sup> یا شرکت مالی غیر بانکی (NBFC)<sup>۲</sup> مؤسسه مالی است که مجوز کامل بانکی ندارد یا توسط یک آژانس نظارتی بانکی ملی یا بین المللی نظارت نمی شود. NBFC خدمات مالی مرتبط با بانک را تسهیل می کند، مانند سرمایه گذاری، تجمع ریسک<sup>۳</sup>، پس انداز قراردادی<sup>۴</sup> و کارگزاری بازار<sup>۵</sup>. یک NBFI به دلیل افزایش عدم پرداخت در رده وام خودرو در تلاش برای نشان دادن سود است.

هدف شرکت تعیین توانایی های بازپرداخت وام مشتری و درک اهمیت نسبی هر پارامتری است که به توانایی وام گیرنده برای بازپرداخت وام کمک می کند.

در این پروژه تلاش بر این است که با توجه به متغیرها و مشخصاتی که از مشتریان (وام گیرندگان) در اختیار داریم و تعدادی از مشتریان که در بخش train وضعیت بازپرداخت آنها مشخص شده، پیش بینی کنیم آیا مشتریان بازپرداخت وام خودرو را انجام می دهند یا خیر؟

---

<sup>1</sup> Non-Banking Financial Institution

<sup>2</sup> Non\_Bank Financial Comapny

<sup>3</sup> Risk pooling

<sup>4</sup> Contractual Savings

<sup>5</sup> Market Brokering



## ۲-۱ معرفی دیتاست (Dataset)

این مسئله شامل دو دیتاست با نام‌های Train\_Dataset و Test\_Dataset است. مدل‌های تحلیلی با کمک بخشی از دیتاست Train\_Dataset ساخته شده و بر روی بخش دیگری از آن، مورد ارزیابی قرار گرفت. سپس، داده‌های دیتاست Test\_Dataset به عنوان ورودی پیش‌بینی به مدل با عملکرد بهتر داده شد. با توجه به این که متغیر پاسخ در دیتاست Test\_Dataset وجود ندارد، لازم است که مدل‌های ارائه‌شده پیش از پیش‌بینی نهایی با استفاده از بخشی از دیتاست Train\_Dataset مورد آزمایش قرار گیرد.

## ۳-۱ Data Dictionary

دیتاست مورد بررسی دارای ۱۲۱۸۵۶ مشاهده (observation) است که تعداد مشتریان بررسی شده می باشد و از ۴۰ متغیر مختلف درباره‌ی اطلاعات مربوط به هر مشتری استفاده شده که در جدول زیر انواع آنها و توضیحات هر کدام قابل مشاهده است. در این جدول، منظور از «Cat.»، دسته‌ای (Categorical) و منظور از «Num.»، عددی (Numerical) است.

جدول ۱-۱

شماره	متغیر	type	توضیح	سطوح داده (Cat.)
۱	ID	Num.	شناسه درخواست وام مشتری	
۲	Client_Income	Num.	درآمد مشتری به واحد دلار	
۳	Car_Owned	Cat.	خودرویی که قبل از درخواست وام برای خودروی دیگر متعلق به مشتری باشد	۰: متعلق به مشتری نباشد ۱: متعلق به مشتری باشد
۴	Bike_Owned	Cat.	دوچرخه متعلق به مشتری (۰ به معنای خیر و ۱ به معنای غیر آن است)	۰: متعلق به مشتری نباشد ۱: متعلق به مشتری باشد
۵	Active_Loan	Cat.	وام فعال دیگری در زمان درخواست وام	۰: وام فعال ندارد ۱: وام فعال دارد
۶	House_Own	Cat.	مشتری دارای خانه می باشد یا خیر	۰: مشتری خانه ندارد

۱: مشتری خانه دارد				
	تعداد فرزندان که مشتری دارد	Cat.	Child_Count	۷
	مبلغ اعتبار وام به واحد دلار	Num.	Credit_Amount	۸
	سالیانه وام به واحد دلار	Num.	Loan_Annuity	۹
	چه کسی مشتری را هنگام درخواست وام همراهی کرد	Cat.	Accompany_Client	۱۰
	نوع درآمد مشتریان	Cat.	Client_Income_Type	۱۱
	سطح تحصیلات مشتری	Cat.	Client_Education	۱۲
D: طلاق گرفته S: مجرد M: متاهل W: بیوه	وضعیت تأهل مشتری	Cat.	Client_Marital_Status	۱۳
	جنسیت مشتری	Cat.	Client_Gender	۱۴
CL: وام نقدی RL: وام گردان	نوع وام	Cat.	Loan_Contract_Type	۱۵
	وضعیت مسکن مشتری	Cat.	Client_Housing_Type	۱۶
ارزش بالاتر یعنی مشتری در جای خوبی زندگی می کند	جمعیت نسبی منطقه ای که مشتری در آن زندگی می کند	Num.	Population_Region_Relative	۱۷
	سن مشتری در زمان ارسال درخواست وام	Num.	Age_Days	۱۸
	چند روز قبل از درخواست وام، مشتری شروع به کسب درآمد کرد	Num.	Employed_Days	۱۹
	چند روز قبل از درخواست وام، مشتری ثبت نام خود را تغییر داد	Num.	Registration_Days	۲۰
	چند روز قبل از درخواست وام، مشتری مدرک هویتی خود را تغییر داد که با آن وام درخواست دهد.	Num.	ID_Days	۲۱
	سن خانه مشتری به سال	Num.	Own_House_Age	۲۲
*: شماره ثبت نشده	شماره موبایل توسط مشتری ارائه شده یا خیر	Cat.	Mobile_Tag	۲۳

۱: شماره ثبت شده				
۰: شماره ثبت نشده ۱: شماره ثبت شده	شماره تلفن خانگی توسط مشتری ارائه شده یا خیر	Cat.	Homephone_Tag	۲۴
۰: شماره در دسترس نبود ۱: شماره در دسترس بود	آیا شماره تلفن کار قابل دسترسی بود یا خیر	Cat.	Workphone_Working	۲۵
	نوع شغل مشتری	Cat.	Client_Occupation	۲۶
	تعداد اعضای خانواده مشتری	Cat.	Client_Family_Members	۲۷
۱: متوسط ۲: خوب ۳: عالی	رتبه شهر مشتری	Cat.	Cleint_City_Rating	۲۸
۰: یکشنبه ۱: دوشنبه ۲: سه شنبه ۳: چهارشنبه ۴: پنجشنبه ۵: جمعه ۶: شنبه	روزی از هفته که مشتری درخواست وام کرده است	Cat.	Application_Process_Day	۲۹
	ساعتی از روز که مشتری درخواست وام کرده است	Num.	Application_Process_Hour	۳۰
Yes: مطابقت دارد No: مطابقت ندارد	آدرس تماس مشتری با آدرس دائمی مطابقت دارد یا خیر	Cat.	Client_Permanent_Match_Tag	۳۱
Yes: مطابقت دارد No: مطابقت ندارد	آدرس کار مشتری با آدرس تماس مطابقت دارد یا خیر	Cat.	Client_Contact_Work_Tag	۳۲
	نوع سازمانی که مشتری در آن کار می کند	Cat.	Type_Organization	۳۳

۳۴	Score_Source_1	Num.	امتیاز از منبع دیگری گرفته شده است. (نرمال شده)	
۳۵	Score_Source_2	Num.	امتیاز از منبع دیگری گرفته شده است. (نرمال شده)	
۳۶	Score_Source_3	Num.	امتیاز از منبع دیگری گرفته شده است. (نرمال شده)	
۳۷	Social_Circle_Default	Num.	چند نفر از دوستان/عضو خانواده مشتری در ۶۰ روز گذشته پرداخت وام را نپذیرفته اند	
۳۸	Phone_Change	Num.	مشتری چند روز قبل از درخواست وام تلفن خود را عوض کرده است	
۳۹	Credit_Bureau	Cat.	تعداد کل درخواست های وام در سال گذشته	
۴۰	Default	Cat.	پرداخت یا عدم پرداخت وام	۰: پرداخت وام ۱: عدم پرداخت وام

## ۲ تحلیل کاوش گرایانه داده (EDA)

به فرآیند کاوش داده‌ها با بکارگیری خلاصه‌های عددی و نمودارها<sup>۶</sup>، برای شناسایی روابط احتمالی میان متغیرها، Exploratory Data Analysis (EDA) گفته می‌شود. به کمک EDA، می‌توان ناهنجاری‌هایی همچون داده‌های پرت (Outliers) و یا مشاهدات غیرعادی را یافت، الگوها را کشف نمود و پرسش‌هایی ایجاد کرد که بعدها توسط روش‌های آماری رسمی‌تر، مورد بررسی قرار گیرد. در اجرای EDA، تحلیلگر یا دانشمند داده همچون کارآگاهی عمل می‌کند که به دنبال یافتن سرنخ‌ها و بینش‌هایی برای شناسایی دلایل ریشه‌ای چالشی است که در تلاش برای حل آن است.

اگرچه جداول خلاصه آماری شامل اطلاعاتی همچون میانگین و انحراف استاندارد نیز بخشی از EDA هستند، بیشتر افراد در اجرای EDA بر گراف‌ها تمرکز دارند. تحلیلگر گراف‌ها و دیگر ابزارهای کاوش گرایانه را بکار می‌گیرد و به جایی می‌رود که داده‌ها او را هدایت می‌کنند. هرگاه گراف یا تحلیلی نتواند اطلاعات کافی در اختیار تحلیلگر قرار دهد، او داده‌ها را از جنبه دیگری مورد بررسی قرار می‌دهد.

دو بخش مهم EDA، تحلیل تک متغیره (Univariate Analysis) است که در سطوح اولیه به کمک جداول خلاصه آماری و مصورسازی داده (Data Visualization) انجام می‌شود.

### ۲-۱ بررسی آماری دیتاست

نخستین نکته‌ای که لازم است پس از بارگذاری دیتاست در نرم‌افزار RStudio مورد بررسی قرار گیرد، فرمت ذخیره‌سازی متغیرهای دیتاست است. در این دیتاست نیز برخی متغیرهای عددی به صورت کاراکتر (Character) ذخیره شده بودند که لازم بود پیش از بررسی آماری آن‌ها را به فرمت عددی تبدیل نمود. از میان ۴۰ متغیری که در این دیتاست وجود دارد، یک متغیر ID<sup>۷</sup>، ۱۶ متغیر Num. و ۲۳ متغیر Cat. دیده می‌شود. البته برخی متغیرهای دسته‌ای ترتیبی (Ordinal Cat.) به گونه‌ای هستند که با توجه به نوع تحلیل و نیاز می‌توان آن‌ها را Cat. و یا Num. در نظر گرفت.

برای بررسی اولیه آماری دیتاست در R، معمولاً از دستور summary() استفاده می‌شود اما در این پروژه از دستور describe() از پکیج «Hmisc» استفاده شد که اطلاعات بیشتری همچون پنج عدد بزرگترین و

<sup>۶</sup> Visualizations

<sup>۷</sup> Identifier

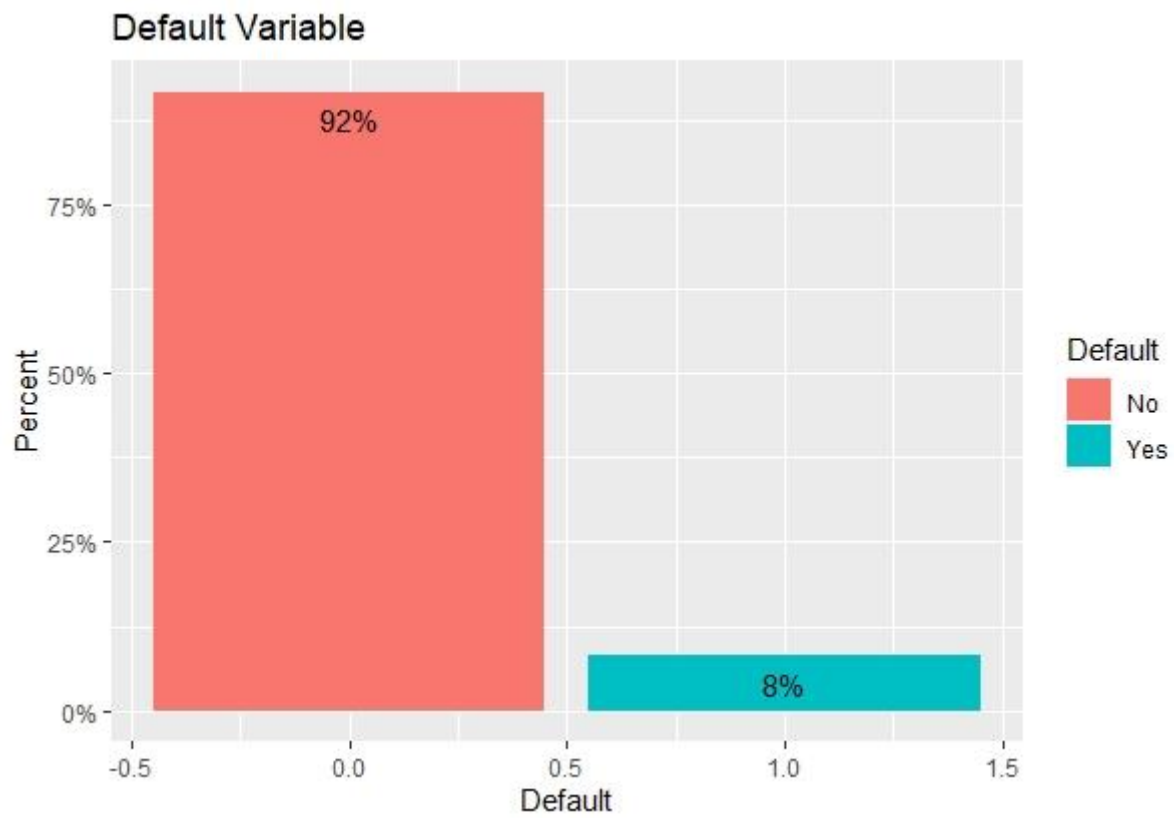
کوچکترین برای متغیرهای Num. و درصد هر یک از سطوح متغیرهای Cat. را ارائه می‌دهد. در ادامه به بیان مهم‌ترین نکات بررسی آماری متغیرهای دیتاست اشاره شده است.

## ۲-۲ تحلیل تک متغیره (Univariate Analysis)

تحلیل تک متغیره ساده‌ترین شکل تحلیل داده است که در آن داده‌های مورد تجزیه و تحلیل، تنها از یک متغیر تشکیل شده است. از آنجایی که تنها حضور یک متغیر مطرح است، با دلایل و روابط سر و کار ندارد. هدف اصلی تحلیل تک متغیره توصیف داده‌ها و یافتن الگوهای موجود در آن‌هاست. با توجه به تعداد بالای متغیرها در این دیتاست، در ادامه، به بیان مهم‌ترین نتایج تحلیل تک متغیره در این دیتاست که به کمک خلاصه‌های آماری و نمودارها بدست آمده، پرداخته شده است:

- **Client\_Income**: میانگین و میانه این متغیر به ترتیب ۱۶۸۶۵ و ۶۷۵۰ بوده است که نشان از چولگی زیاد و نرمال نبودن توزیع داده‌ها دارد. نکته دیگر این است که بزرگترین ۵ داده این متغیر از ۳۸۲۵۰۰ تا ۱۸۰۰۰۰۹ هستند. حضور داده‌ای به بزرگی یک میلیون و ۸۰۰ هزار در کنار داده‌هایی با میانه ۶۷۵۰، صرف نظر از درست بودن یا نبودن داده، می‌تواند بر روی برازش مدل، تاثیر منفی داشته باشد.
- **Child\_Count**: این متغیر از ۰ تا ۱۹ مقدار گرفته است ولی توزیع داده‌ها نامتوازن بوده است، به گونه‌ای که اگر هر مقدار را یک دسته در نظر بگیریم، هر چه از صفر به سمت ۱۹ حرکت می‌کنیم، تعداد افراد کاهش می‌یابد، این کاهش به قدری شدید است که افراد بدون فرزند، ۸۲۸۳۴ نفر بوده‌اند، در حالی که افراد با ۵ فرزند، ۳۴ نفر و افراد دارای بیش از ۷ فرزند، انگشت شمار بوده‌اند.
- **Accompany Client**: این متغیر Cat.، هفت سطح دارد که یکی از آن‌ها تنها دارای ۱۲ داده است و نام آن هم «##» است که نامشخص است.
- **Client\_Income\_Type**: این متغیر Cat.، هشت سطح دارد ولی چهار مورد از آن‌ها به ترتیب یک، دو، شش و هشت داده دارند.
- **Client\_Gender**: برای این متغیر سه سطح تعریف شده است که یکی از آن‌ها «XNA» است که دقیقاً مشخص نیست منظور از این مقدار چه بوده است، با این حال این سطح تنها سه داده دارد.

- **Employed\_Days** : در توضیحات این متغیر آمده که چند روز پیش از درخواست وام، شخص آغاز به کسب درآمد کرده است، با این حال بیش از ۲۱۰۰۰ داده از این متغیر مقداری بیش از ۳۰۰۰۰۰ دارند که تقریباً معادل ۹۰۰ سال است.
- **Mobile\_Tag** : این یک متغیر Boolean است ولی سطح صفر آن تنها یک داده دارد و ۱۲۱۸۵۵ داده مقدار ۱ دارند.
- **Client\_Family\_Members** : این متغیر نیز شرایطی همچون متغیر **Child\_Count** دارد و انتظار می‌رود که رابطه‌ای میان این دو متغیر برقرار باشد.
- **Type\_Organization** : یک متغیر Cat. است که ۵۸ سطح دارد. به نظر می‌رسد برخی از سطوح این متغیر، زیرمجموعه‌ای از یک سطح کلی بزرگتر باشند، برای نمونه ۱۳ سطح این متغیر دارای نام‌های «Industry: type 1» تا «Industry: Type 13» هستند. همچنین یک سطح نامشخص به نام «XNA» وجود دارد.
- **Source\_Score\_1** ، **Source\_Score\_2** و **Source\_Score\_3** : دقیقاً مشخص نیست که این سه متغیر بیانگر چه ویژگی هستند و در راهنمای دیتاست نیز توضیحی در این باره ارائه نشده است، تنها می‌دانیم که این سه متغیر به صورت پیوسته مقادیر صفر تا یک را اختیار می‌کنند.
- **Default** : این متغیر، متغیر پاسخ مساله است و همان‌گونه که پیش از این گفته شد، چون یک متغیر Boolean است، برای پیش‌بینی آن لازم است که از مدل‌های **Classification** استفاده شود. نکته‌ای که درباره این متغیر وجود دارد عدم توازن آن است که می‌تواند بر دقت مدل برازش شده تاثیر منفی داشته باشد. (



شکل ۱-۲ نمودار Bar Chart برای متغیر Default



### ۳ آماده‌سازی داده (Data Preparation)

برای دستیابی به یک مدل پیش‌بینی با دقت مناسب، لازم است که داده‌ها پس از بررسی مورد ویرایش قرار گیرند تا مدل نهایی خطای کمتری داشته باشد. این آماده‌سازی می‌تواند شامل مراحل مختلفی باشد که در این دیتاست، شامل مدیریت داده‌های پرت یا نادرست و بررسی و جایگزینی (Imputation) داده‌های گم شده (Missing Values) است. در طی فرآیند آماده‌سازی داده‌ها ممکن است تصمیم گرفته شود که بخشی از داده‌ها حذف شوند. در این پروژه به منظور مقایسه و دستیابی به مدل مناسب دو دیتاست با دو رویکرد متفاوت ایجاد می‌شود به گونه‌ای که:

۱. **df.norm**: در دیتاست اول هیچ داده‌ای حذف نمی‌گردد و تنها به گونه‌ای که الگوی داده‌ها حفظ

شود، جایگزینی صورت می‌گیرد.

۲. **df.omit**: ترکیبی از حذف و جایگزینی صورت می‌گیرد.

#### ۳-۱ داده‌های پرت (Outliers)

یک تابع برای مشخص نمودن داده‌های پرت نوشته شد و تعداد داده‌های پرت متغیرهای Num. محاسبه شد، در دیتاست دوم تمامی داده‌های پرت حذف شدند ولی در دیتاست اول، موارد زیر اعمال شد:

- **Client\_Income**: تعداد داده‌های پرت، ۵۳۹۲ بدست آمد که تعداد نسبتاً زیادی است، با این حال

در مورد این متغیر، گسترده بودن دامنه داده‌ها تا حدودی طبیعی است، بنابراین سطح ۱۰۰۰۰۰ به عنوان مبنا در نظر گرفته شد و داده‌های بیشتر از این سطح با عدد ۱۰۰۰۰۰ جایگزین شد، بدین ترتیب مقدار ۱۰۱ داده جایگزین شد.

- **Child\_Count**: تعداد داده‌های پرت این متغیر، ۱۶۵۹ بدست آمد ولی با بررسی سطوح داده عدد

۴ به عنوان پایه در نظر گرفته شد، مقادیر بیشتر از ۴ با عدد ۴ جایگزین شد.

- **Client\_Family\_Members**: همانند متغیر پیشین با این تفاوت که عدد ۶ به عنوان مبنا در نظر گرفته شد.

- **Credit\_Bureau**: همانند متغیر پیشین با این تفاوت که عدد ۱۰ به عنوان مبنا در نظر گرفته شد.

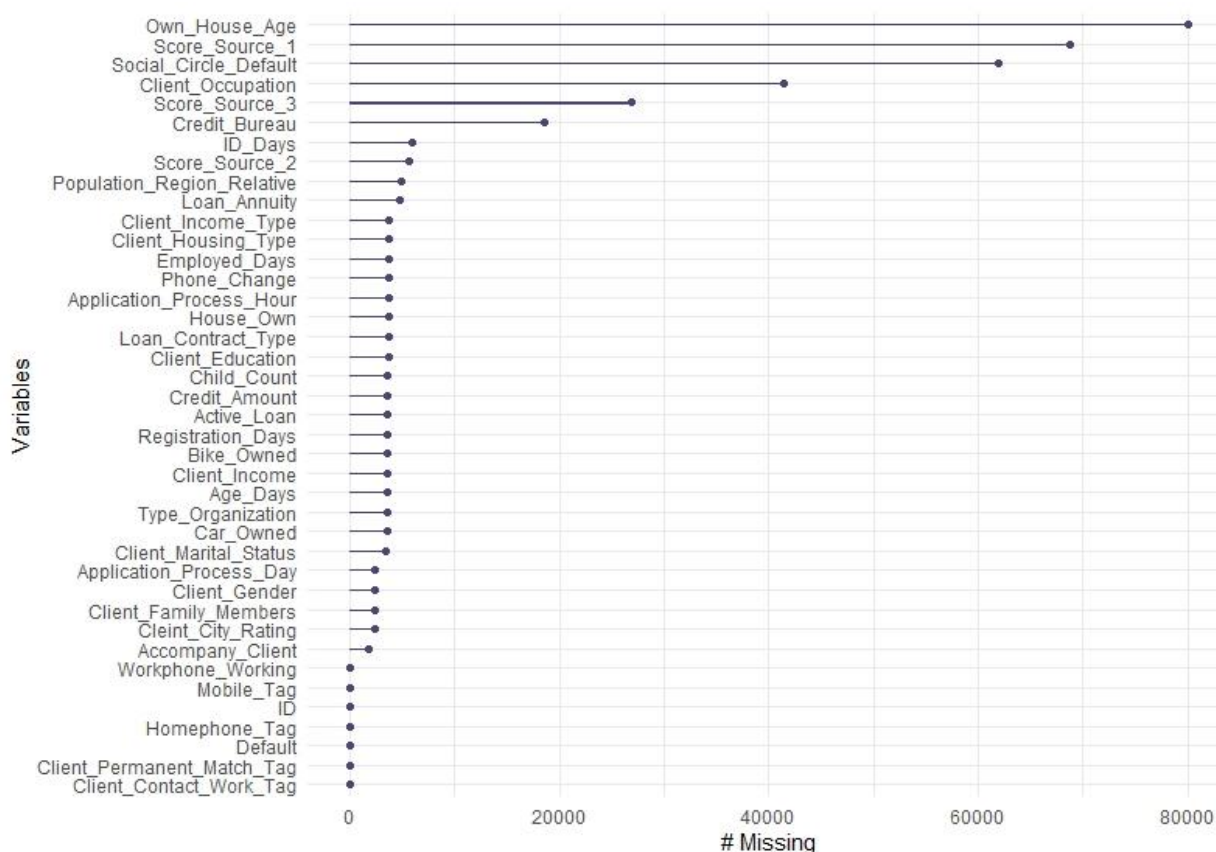
- **Employment\_Days**: تعداد داده های پرت ۲۱۴۸۸ بدست آمد، با توجه به غیر منطقی بودن مقادیر این داده ها، تمامی آن ها به داده گم شده یا ناموجود (NA)<sup>۸</sup> تبدیل شد تا با روش هایی که در ادامه بیان خواهد شد، با مقدار مناسب جایگزین شود.
- **Score\_Source\_2**: تعداد داده های پرت، ۶ بدست آمد که به NA تبدیل شد.  
در مورد متغیرهای Cat.، نیز موارد زیر اعمال شد:
- **Client\_Income\_Type**: چهار سطح این متغیر که تعداد بسیار کمی داده داشتند، در یک سطح جدید به نام «Other»، ادغام شدند.
- **Accompany\_Client**: سطح نامشخص «##» در دیتاست اول به NA تبدیل شد و در دیتاست دوم حذف شد.
- **Client\_Gender**: سطح نامشخص «XNA» در دیتاست اول به NA تبدیل شد و در دیتاست دوم حذف شد.
- **Type\_Organization**: سطوح مشابه این متغیر در یکدیگر ادغام شدند و به این ترتیب تعداد سطوح متمایز از ۵۸ به ۳۵ کاهش یافت. با این حال این تعداد نیز زیاد است و در اجرای الگوریتم های جایگزینی که در ادامه بیان خواهد شد مشکل ایجاد خواهد کرد، از این رو، داده های NA این متغیر به سطح نامشخص «XNA» تبدیل شد که پیش از این حدود ۲۲۰۰۰ داده داشت.

### ۲-۳ داده های گم شده یا ناموجود (NA - Missing)

برخی از داده ها در دیتاست به صورت NA ثبت شده اند ولی برخی دیگر به ویژه در فرمت کاراکتر خالی (NULL) هستند و این مساله سبب می شود که در نگاه سطحی به عنوان NA در نظر گرفته نشوند. بدین منظور در هنگام وارد نمودن دیتاست، لازم است در دستور `read.csv()`، ویژگی `na.strings` را برابر با مقدار TRUE قرار دهیم. در ادامه برای بررسی داده های NA، تعداد کل سطرها که دیتاست شمارش شد که عدد ۲۵۶۸ بدست آمد که معادل ۲ درصد از کل داده ها است، به عبارت دیگر، حدود ۹۸ درصد از سطرها دیتاست، دست کم یک مقدار NA دارند. سپس تعداد NA ها در متغیرهای مختلف مورد بررسی

<sup>8</sup> Not Available

قرار گرفت و مشخص شد که از میان ۴۰ متغیر موجود در داده، شش مورد بیش از پنج درصد مقدار NA دارند که این مقادیر به ترتیب از ۶۵ تا ۱۵ درصد است (شکل ۳-۱).



شکل ۳-۱ تعداد مقادیر NA متغیرها

تغییرات مرحله آماده‌سازی داده که تا کنون بیان شده در قالب یک تابع، پیاده‌سازی شد. بدین ترتیب، در نتیجه اعمال این تغییرات، دیتاست اول بدون حذف هیچ گونه داده‌ای دارای ۱۲۱۸۵۶ مشاهده و دیتاست دوم دارای ۱۱۰۸۸۱ مشاهده که نشان از حذف ۱۰۹۷۵ سطر دارد. همچنین تعداد سطرهای بدون NA از ۲۵۶۸ مورد به ۲۶۳۵ مورد در دیتاست اول و ۲۱۴۲ مورد در دیتاست دوم تغییر کرده است.

### ۳-۳ مقداردهی (Imputaion)

مدیریت داده‌های از دست رفته از اهمیت بسیار بالایی برخوردار است، زیرا که بسیاری از الگوریتم‌های یادگیری آماری (Statistical learning) و یادگیری ماشین (Machine learning) از NA ها پشتیبانی نمی‌کنند، با این حال الگوریتم‌های K-nearest و Naïve Bayes از داده‌های از دست رفته پشتیبانی می‌کنند. دلیل دیگر این است که ممکن است این داده‌ها در نهایت منجر به ایجاد یک مدل جانبدارانه (Biased

(Model) و یا عدم دقت در تجزیه و تحلیل‌ها گردند. به طور کلی دو راه برای مدیریت داده‌های از دست رفته وجود دارد که یکی حذف و دیگری مقداردهی (Imputation) داده‌های از دست رفته است. حذف داده‌های از دست رفته می‌تواند از راه حذف سطرها و یا ستون‌های دیتاست صورت گیرد.

Imputation به معنای جایگزینی یک مقدار از دست رفته (NA) با یک مقدار دیگر بر اساس یک برآورد معقول است. به عبارت دیگر از داده‌های موجود در دیتاست برای مقداردهی داده‌های از دست رفته به منظور دستیابی به یک دیتاست کامل‌تر استفاده می‌شود. روش‌های متفاوتی برای Imputation وجود دارد که از میان آن‌ها می‌توان به انتخاب داده تصادفی، بکارگیری میانه یا میانگین برای داده‌های Num. و بکارگیری مد برای داده‌های Cat. نام برد. با این حال این روش‌ها هرگز توصیه نمی‌شوند زیرا می‌توانند کیفیت داده‌ها را کاهش داده و الگوی داده‌ها را به طور کلی تغییر دهند. برای نمونه در مورد این دیتاست، متغیر Own\_House\_Age حدود ۸۰ هزار داده از دست رفته از ۱۲۱ هزار مشاهده دارد، حال اگر این داده‌ها را با میانه جایگزین کنیم، به یکباره ۶۵ درصد به فراوانی میانه افزوده می‌شود و این روش قطعاً الگوی داده‌ها را به طور کلی تغییر خواهد داد.

یکی دیگر از روش‌های Imputation بکارگیری الگوریتم‌های پیش‌بینی برای داده‌های از دست رفته است که در این پروژه از این روش استفاده شده است. بدین منظور از پکیج «mice»<sup>۹</sup> استفاده شده است که یک پکیج مقداردهی چندگانه (Multiple Imputation) است. روش کار mice به این صورت است که در ابتدا زنجیره‌ای از معادلات بر اساس دیگر متغیرها برای متغیری که قرار است مقادیر از دست رفته آن impute شوند تشکیل می‌دهد. سپس، الگوریتم تعیین شده برای متغیر مورد نظر را به صورت تکرارشونده به تعداد تکرارهای از پیش تعیین شده و بر اساس زنجیره معادله تشکیل شده اجرا می‌نماید. هر بار که تعداد تکرارهای تعیین شده به پایان می‌رسد، یک imputation به صورت یک dataframe بر اساس نتایج تکرارها برای دیتاست مورد نظر ایجاد می‌شود. تعداد imputation ها نیز قابل تنظیم است و پس از پایان اجرای mice می‌توان تعیین کرد که یکی از imputation ها انتخاب شود و یا برآیندی از آن‌ها به کمک دستور pool() محاسبه و سپس جای‌گذاری شود. در پکیج mice، ۲۷ الگوریتم متفاوت برای گونه‌های مختلف داده در نظر گرفته شده است که می‌توان آن‌ها را به صورت دستی یا پیش فرض برای impute کردن داده‌های از دست رفته بکار گرفت. به طور کلی بکارگیری پکیج‌های imputation برای متغیرهای Cat. با تعداد سطوح بالا

<sup>۹</sup> Multivariate Imputation via Chained Equations

آسان نیست و در این پژوهش نیز موفق نشدیم متغیر `Type_Organization` را به وسیله `imputation` مقداردهی کنیم و به همین دلیل داده‌های `NA` این متغیر با سطح نامشخص «`XNA`» جایگزین شد.

در این پروژه از الگوریتم `Predictive Mean Matching` برای متغیرهای `Num.` و `Ordinal Cat.`، الگوریتم `Polytomous Logistic Regression` برای متغیرهای `Cat.` با تعداد سطوح بیشتر از دو و الگوریتم `Logistic Regression` برای متغیرهای `Boolean` استفاده شده است. همچنین به دلیل زمان‌بر بودن فرآیند `imputation`، تنها یک `imputation` حاصل از پنج بار تکرار برای دیتاست‌های اول و دوم که در گام پیشین آماده شده‌اند، در نظر گرفته شده است. در ادامه به توضیح دو الگوریتم بکارگرفته شده در این پروژه پرداخته شده است.

### ۱-۳-۳ Predictive Mean Matching (PMM)

`PMM` مقدار پیش‌بینی شده برای متغیر هدف `Y` را به کمک مدل `imputation` خاص خود محاسبه می‌کند. برای هر مقدار از دست رفته، متد `PMM` مجموعه کوچکی از کاندیداها (معمولاً با سه، پنج یا ده عضو) از تمامی سطرهایی که مقدار پیش‌بینی شده برای آن‌ها نزدیک‌ترین مقدار به مقدار پیش‌بینی شده برای سطر با مقدار از دست رفته است، تشکیل می‌دهد. یک کاندیدا به طور تصادفی از میان کاندیداها انتخاب می‌شود و مقدار مشاهده شده آن برای جایگزینی مقدار از دست رفته در نظر گرفته می‌شود. فرض بر این است که توزیع سلول از دست رفته با داده‌های مشاهده شده کاندیداها یکسان است.

`PMM` یک روش ساده و همه‌کاره است، این روش برای تمامی گونه‌های متغیرها کاربرد دارد، با این حال عملکرد آن برای متغیرهای عددی بهتر است. این متغیر همچنان در برابر تبدیل متغیر هدف نیز انعطاف پذیر است برای نمونه وارد نمودن  $\log(Y)$  می‌تواند نتایجی نزدیک به ورود  $\exp(Y)$  داشته باشد. این روش امکان بکارگیری برای متغیرهای گسسته را نیز فراهم می‌آورد. مقداردهی‌ها بر اساس مقادیر مشاهده شده در سطرها دیگر دیتاست است، بنابراین این مقداردهی‌ها واقع‌بینانه هستند، محاسبات خارج از محدوده داده‌های موجود رخ نخواهد داد و خبری از مشکلات مربوط به مقداردهی‌های نادرست و بی‌معنی همچون منفی شدن وزن و غیره نخواهد بود. همچنین مدل این روش یک مدل ضمنی (`implicit`) است به این معنی که نیاز به یک مدل صریح (`explicit`) برای مقادیر از دست رفته نیست، بنابراین این روش نسبت به بسیاری از روش‌های دیگر آسیب‌پذیری کمتری دارد.

### ۲-۳-۳ Polytomous (Multinomial) Logistic Regression

رگرسیون لجستیک چندگانه یا چندجمله‌ای یک روش طبقه‌بندی (Classification) است که رگرسیون لجستیک را به مسائل چند سطحی (Multiclass) تعمیم می‌دهد، بدین معنی که بیش از دو نتیجه گسسته احتمالی وجود داشته باشد. به عبارت دیگر، مدلی برای پیش‌بینی احتمال نتایج ممکن یک متغیر وابسته Cat. که بیش از دو سطح دارد، بر اساس مجموعه‌ای از متغیرهای مستقل (که می‌توانند حقیقی، Boolean، Cat. و ... باشند)، ارائه می‌دهد. رگرسیون لجستیک چندگانه با نام‌های مختلفی همچون Polytomous LR، The Maximum، Multinomial logit (mlogit)، Softmax Regression، Multiclass LR، Conditional Maximum entropy و entropy (MaxEnt) Classifier شناخته می‌شود.

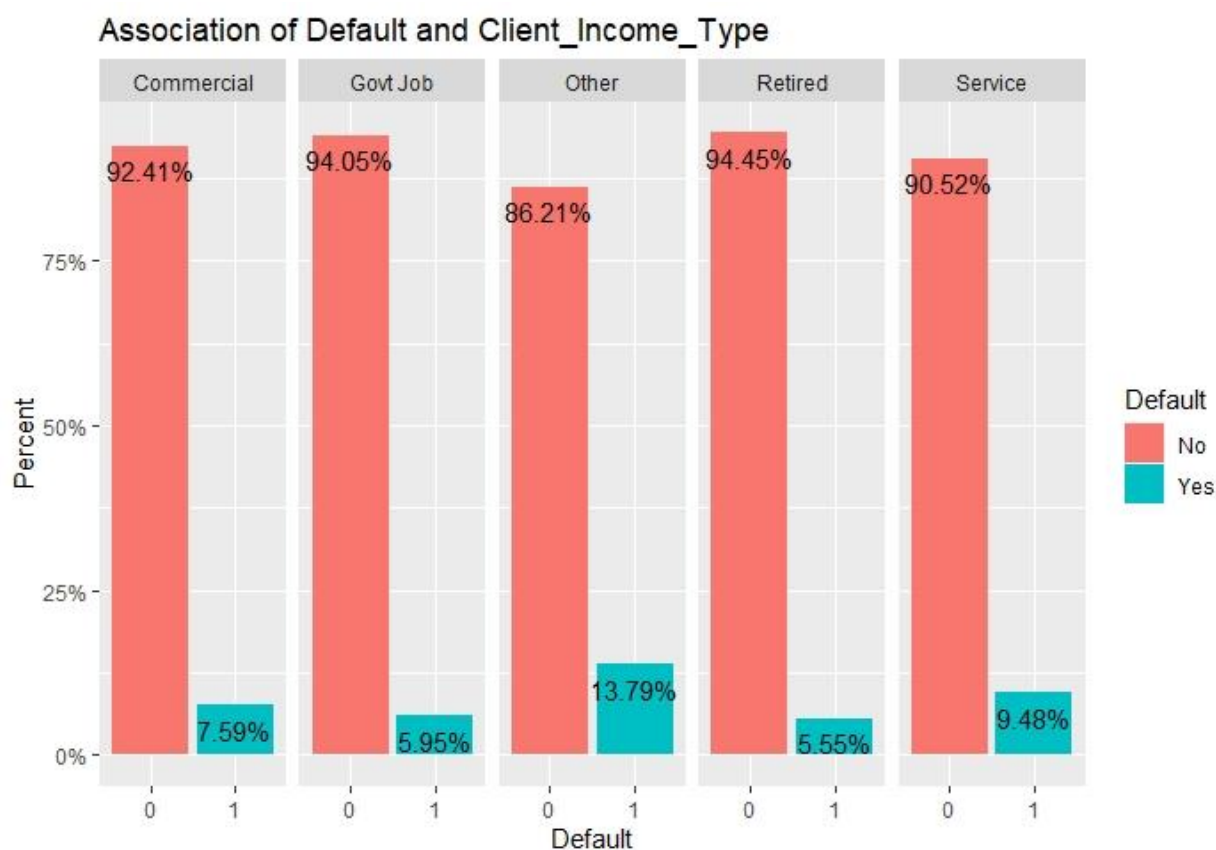
نقطه قوت این روش نسبت به روش رگرسیون لجستیک باینری، بکارگیری اندازه نمونه تمام دسته‌های نتایج در تخمین احتمال پارامترها و واریانس است، زیرا روش باینری تنها از اندازه نمونه دو دسته نتیجه در تخمین احتمال پارامترها و واریانس استفاده می‌کند. یکی از مهم‌ترین فرض‌های این مدل این است که پاسخ‌های نهایی از یکدیگر مستقل‌اند. در تجزیه و تحلیل رگرسیون چندگانه، بیش از یک مدل logit بر داده‌ها برازش می‌گردد، سپس برآیند تمامی مدل‌های logit، مدل رگرسیون چندگانه را تشکیل می‌دهد و همه با هم برای پیش‌بینی احتمال هر نتیجه بکار گرفته می‌شوند.

## ۴ مصورسازی داده (Data Visualization)

در این بخش برای رسم تمامی نمودارها از دیتاست اول استفاده شده است. در نمودارهای Box plot برای نمایش بهتر نمودار، پیش از رسم، داده‌های پرت حذف شده‌اند.

### ۴-۱ رابطه میان Client\_Income\_Type و Default

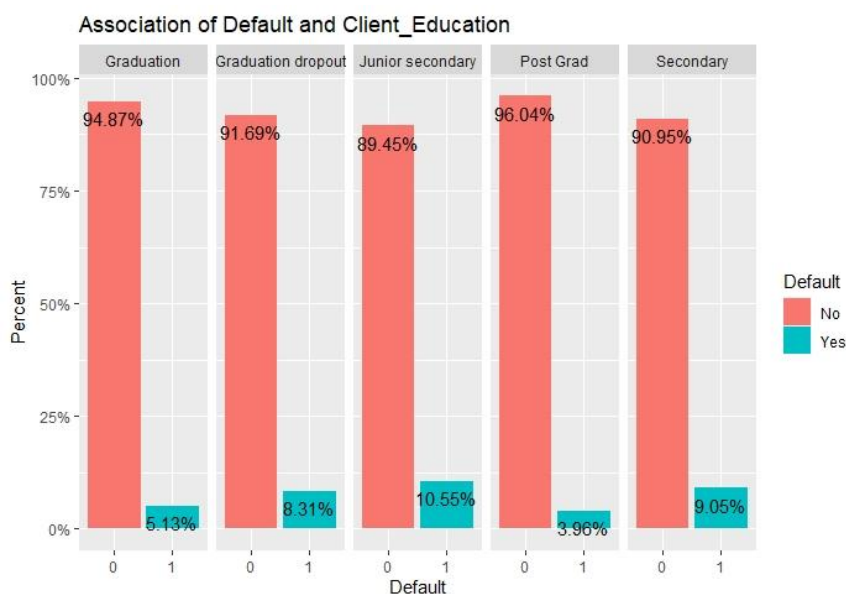
به نظر می‌رسد افراد دارای منبع درآمد خدماتی و یا تجاری درصد default بیشتری داشته‌اند.



شکل ۴-۱ رابطه میان Client\_Income\_Type و Default

## ۲-۴ رابطه میان Client\_Education و Default

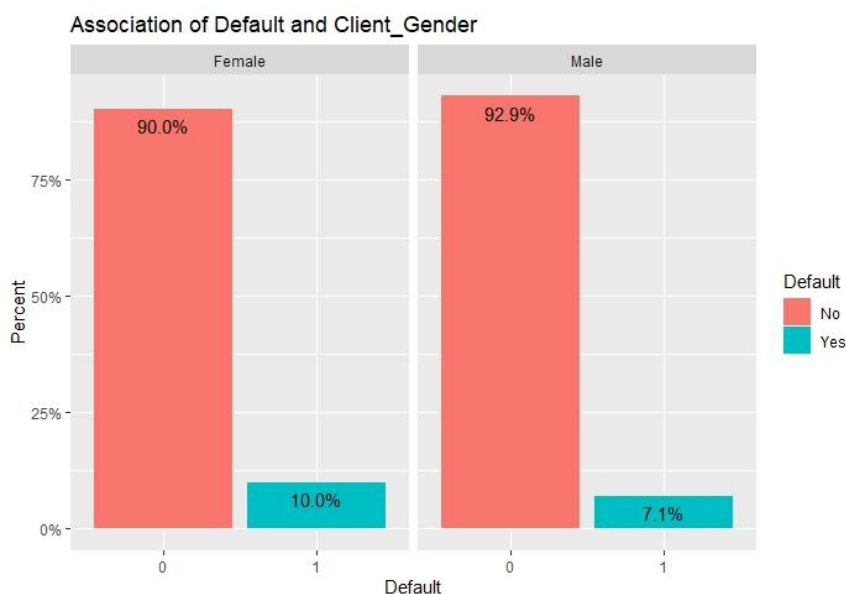
به نظر می رسد افراد دارای سطح تحصیلات Junior Secondary بیشترین درصد Default و افراد دارای سطح تحصیلات Post Grad دارای کمتری درصد default بوده اند.



شکل ۲-۴ رابطه میان Client\_Education و Default

## ۳-۴ رابطه میان Gender و Default

به نظر می رسد درصد default مردان اندکی بیشتر از زنان بوده است.

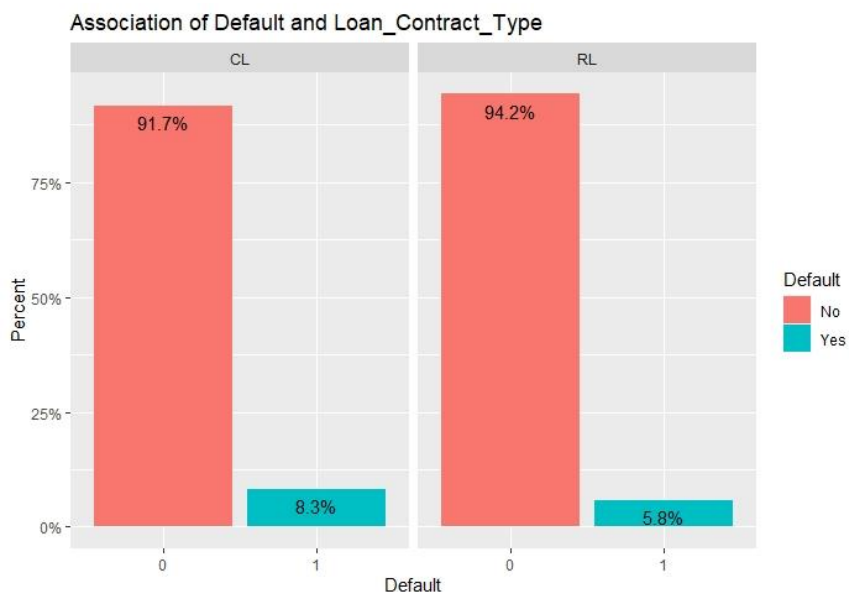


شکل ۳-۴ رابطه میان Gender و Default



#### ۴-۴ رابطه میان Loan\_Contract\_Type و Default

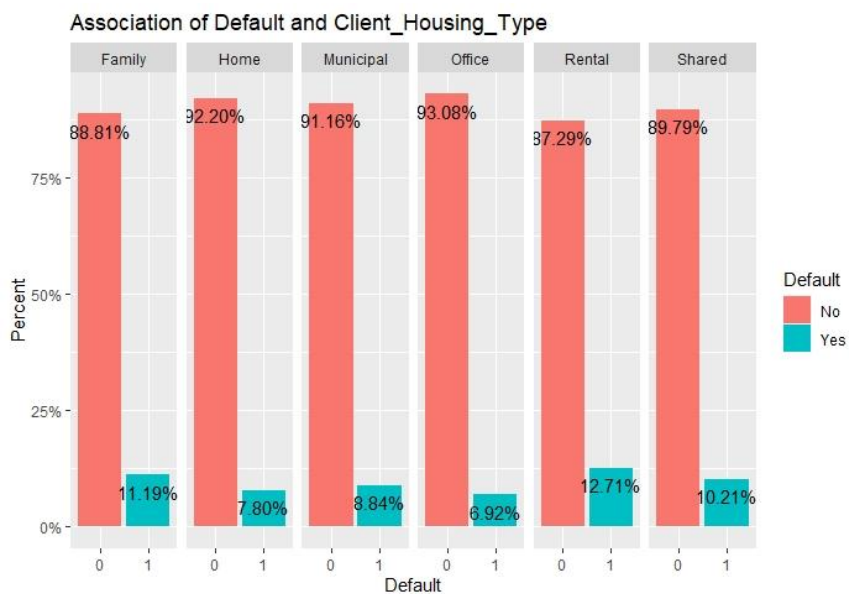
به نظر می رسد افراد دارای وام از نوع چرخشی درصد default کمتری داشته اند.



شکل ۴-۴ رابطه میان Loan\_Contract\_Type و Default

#### ۴-۵ رابطه میان Client\_Housing\_Type و Default

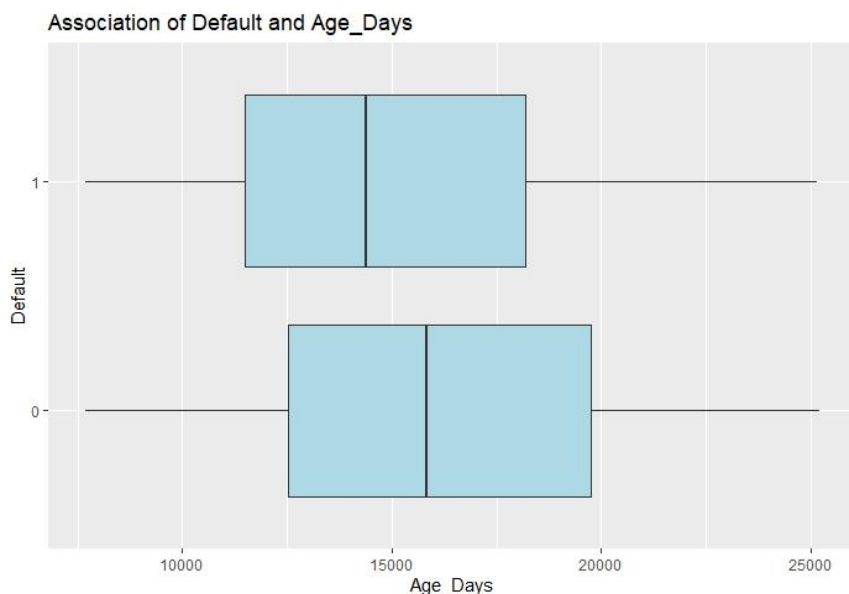
به نظر می رسد افراد اجاره نشین درصد default بیشتری داشته اند.



شکل ۴-۵ رابطه میان Client\_Housing\_Type و Default

#### ۴-۶ رابطه میان Age\_Days و Default

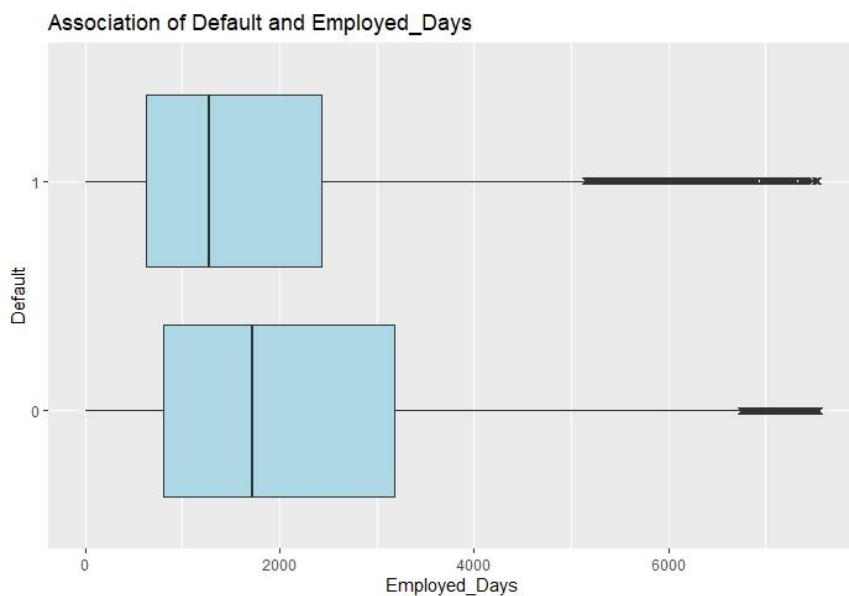
به نظر می رسد افرادی که default کرده اند، به طور متوسط جوان تر بوده اند.



شکل ۴-۶ رابطه میان Age\_Days و Default

#### ۴-۷ رابطه میان Employed\_Days و Default

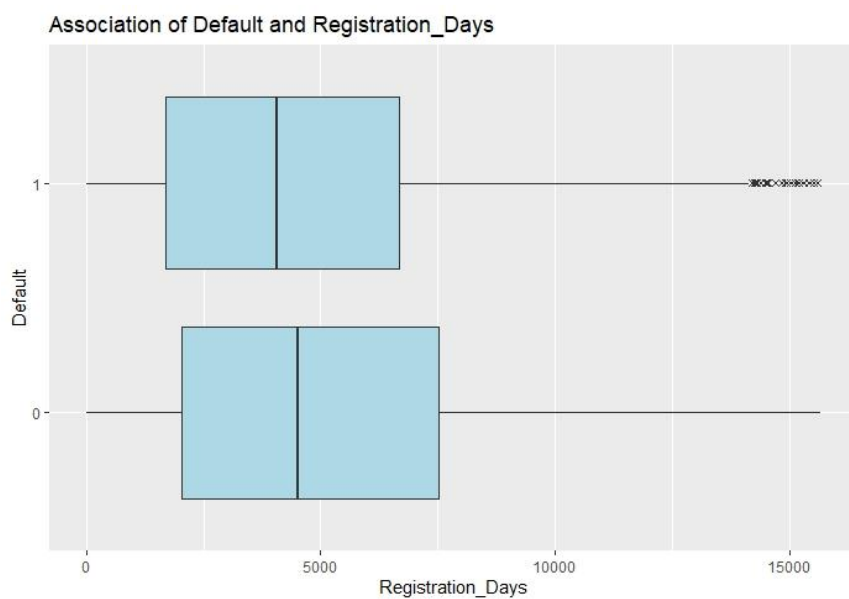
به نظر می رسد افرادی که default کرده اند، به طور متوسط سابقه کار کمتری داشته اند.



شکل ۴-۷ رابطه میان Employed\_Days و Default

#### ۸-۴ رابطه میان Default و Registration\_Days

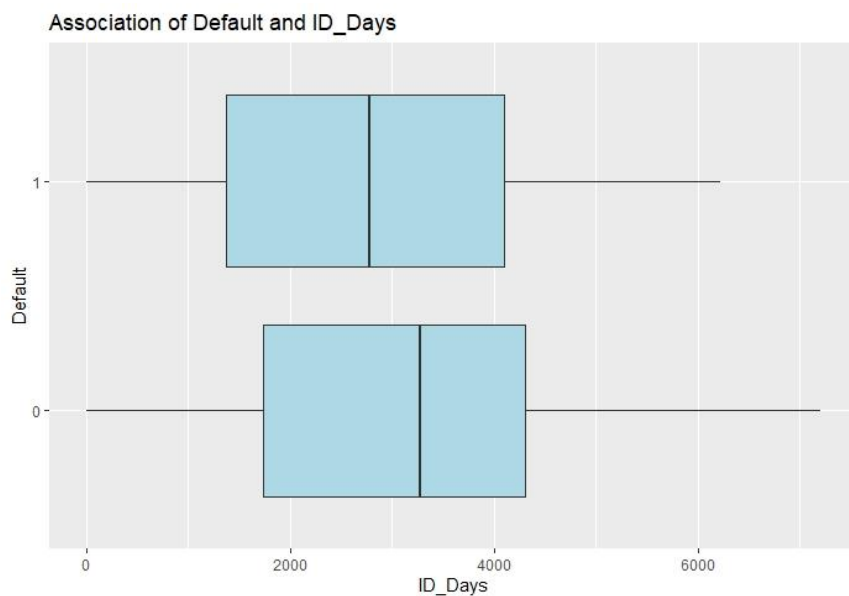
به نظر می‌رسد افرادی که default کرده اند، به طور متوسط عمر حساب بانکی آن‌ها کمتر بوده است.



شکل ۸-۴ رابطه میان Default و Registration\_Days

#### ۹-۴ رابطه میان Default و ID\_Days

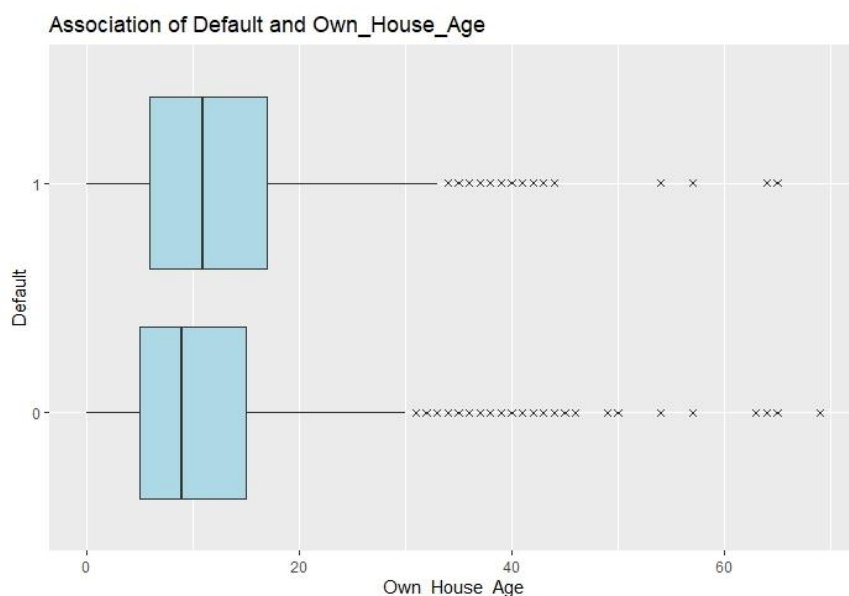
به نظر می‌رسد افرادی که default کرده اند، به طور متوسط مدارک هویتی خود را زودتر تغییر داده اند.



شکل ۹-۴ رابطه میان Default و ID\_Days

#### ۱۰-۴ رابطه میان Own\_House\_Age و Default

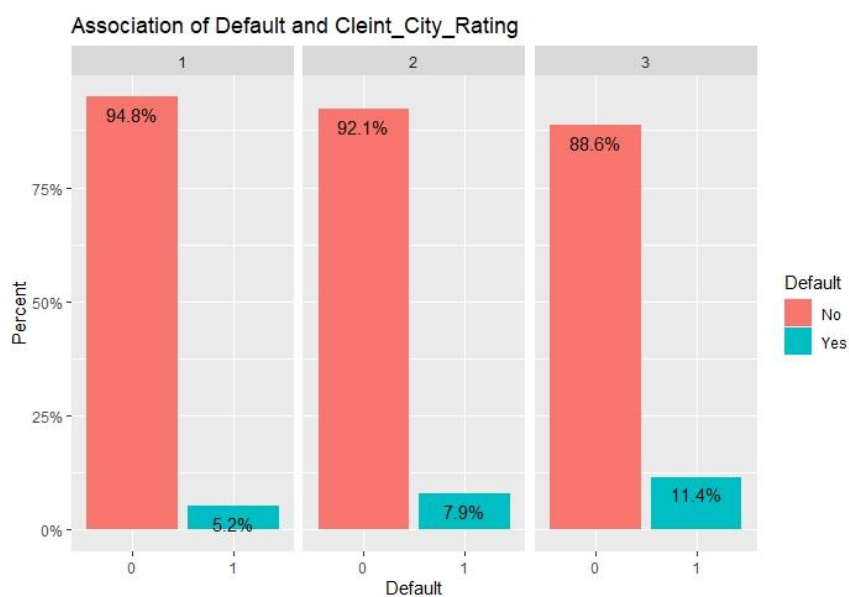
به نظر می رسد افرادی که default کرده اند، به طور متوسط سابقه مالکیت خانه بیشتری داشته اند.



شکل ۱۰-۴ رابطه میان Own\_House\_Age و Default

#### ۱۱-۴ رابطه میان Client\_City\_Rating و Default

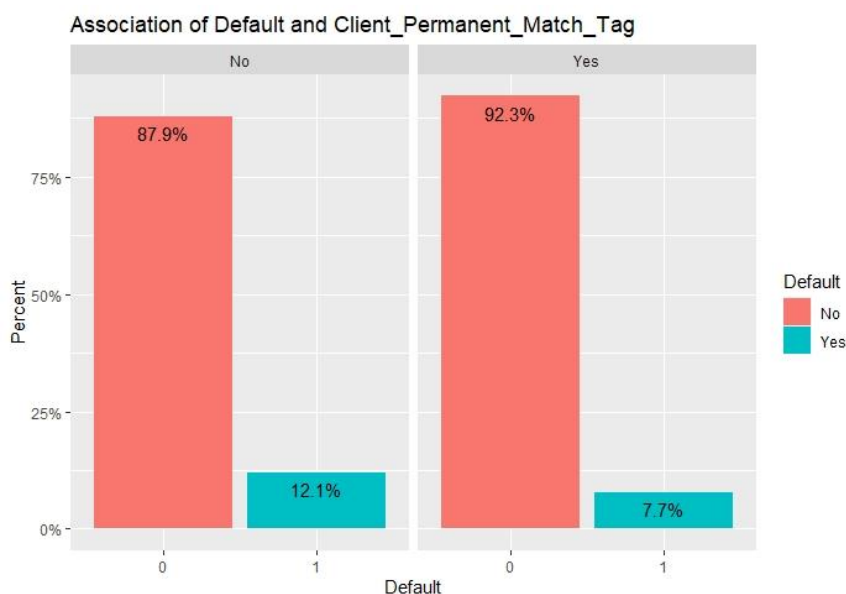
به نظر می رسد بالاتر بودن رتبه شهری رابطه عکس با default کردن داشته باشد.



شکل ۱۱-۴ رابطه میان Client\_City\_Rating و Default

## ۱۲-۴ رابطه میان Default و Client\_Permanent\_Match\_Tag

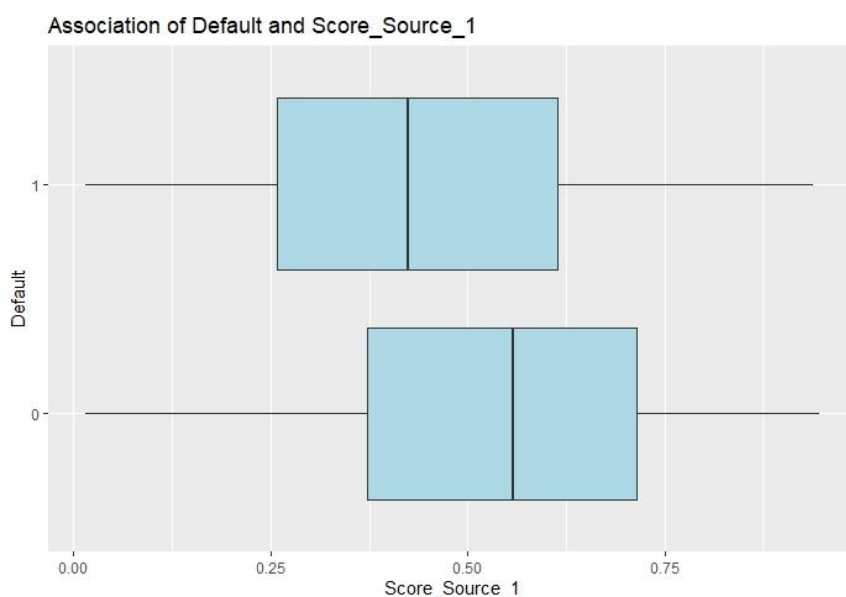
به نظر می‌رسد افرادی که مطابقت آدرس نداشته اند، بیشتر default کرده اند.



شکل ۱۲-۴ رابطه میان Default و Client\_Permanent\_Match\_Tag

## ۱۳-۴ رابطه میان Default و Score\_Source\_1

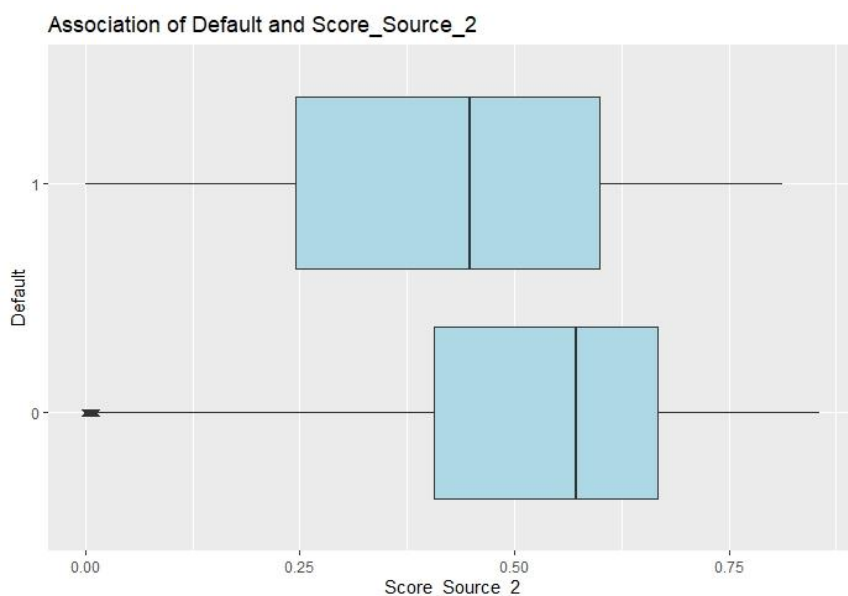
به نظر می‌رسد افرادی که default کرده اند، به طور متوسط مقدار Score\_Source\_1 کمتری داشته اند.



شکل ۱۳-۴ رابطه میان Default و Score\_Source\_1

#### ۱۴-۴ رابطه میان Default و Score\_Source\_2

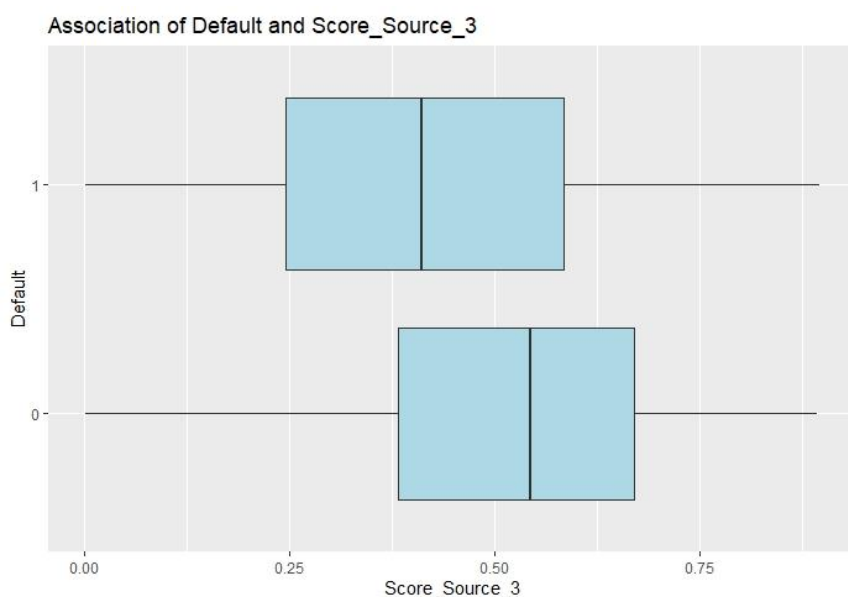
به نظر می‌رسد افرادی که default کرده اند، به طور متوسط مقدار Score\_Source\_2 کمتری داشته اند.



شکل ۱۴-۴ رابطه میان Default و Score\_Source\_2

#### ۱۵-۴ رابطه میان Default و Score\_Source\_3

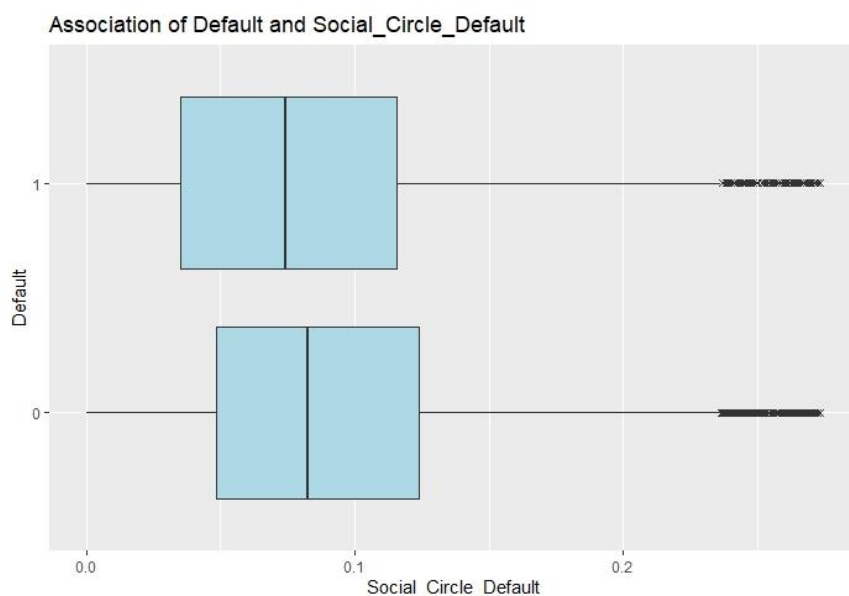
به نظر می‌رسد افرادی که default کرده اند، به طور متوسط مقدار Score\_Source\_3 کمتری داشته اند.



شکل ۱۵-۴ رابطه میان Default و Score\_Source\_3

#### ۱۶-۴ رابطه میان Default و Social\_Circle\_Default

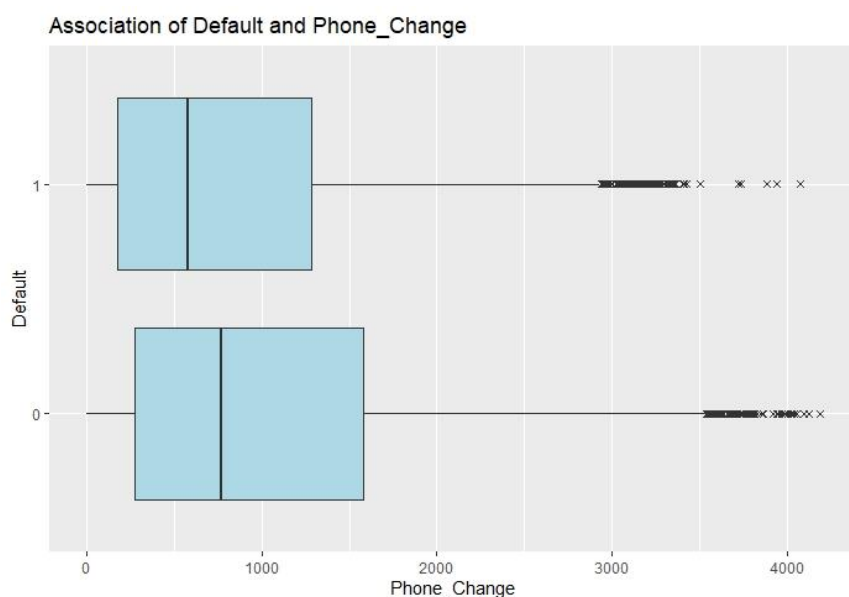
به نظر می‌رسد افرادی که default کرده اند، سابقه خانوادگی default کمتری داشته اند.



شکل ۱۶-۴ رابطه میان Default و Social\_Circle\_Default

#### ۱۷-۴ رابطه میان Default و Phone\_Change

به نظر می‌رسد افرادی که default کرده اند، به طور متوسط شماره تلفن خود را زودتر تغییر داده اند.



شکل ۱۷-۴ رابطه میان Default و Phone\_Change

## ۵ آمار توصیفی (Descriptive Statistics)

در این بخش یا آزمون تاثیرگذار بودن یک عامل مورد بررسی قرار گرفته و یا این که یک آزمون ایجاد شده است.

### ۱-۵ رابطه میان Client\_Education و Default

افرادی که سطح تحصیلات post grad دارند با احتمال بیشتری وام های خود را پرداخت می کنند نسبت به افرادی که سطح تحصیلات graduation دارند.

از تست proportion test استفاده می کنیم:

$$H_0: P_1 > P_2$$

$$H_1: P_1 < P_2$$

N1: تعداد افراد با سطح تحصیلات post grad

N2: تعداد افراد با سطح تحصیلات graduation

Win1: مجموع عدم بازپرداخت وام توسط افراد با سطح تحصیلات graduation

Win2: مجموع عدم بازپرداخت وام توسط افراد با سطح تحصیلات post grad

با انجام تست به نتیجه زیر می رسیدیم:

$$P_1: 0.03960396$$

$$P_2: 0.05127774$$

در نتیجه با توجه به تست انجام شده فرض صفر رد می شود و افرادی که سطح تحصیلات post grad دارند با احتمال کمتری وام های خود را پرداخت نمی کنند.

### ۲-۵ رابطه میان Client\_Gender و Default

ادعا شده که جنسیت مشتریان بر روی پرداخت و عدم پرداخت قسط های وام موثر است. برای تست این فرضیه از آزمون مربع کای دو (استقلال) استفاده می کنیم:



ابتدا یک جدول با متغیرهای Gender و Defult ایجاد میکنیم و سپس از تست  $\chi^2$  استفاده می‌کنیم.

چون مقدار p-value برابر  $e^{-162.2}$  است و از ۰.۰۵ بیشتر شده آزمون معنادار می‌شود.

### ۳-۵ آزمون میانگین درآمد

معاون بانک ادعا کرده که میانگین درآمد مشتریان درخواست کننده وام بیشتر از ۱۶۰۰۰ دلار می‌باشد.

$$H_0: M \leq 20000$$

$$H_1: M > 20000$$

برای تست این فرض از آزمون نرمال z-test با فرض انحراف معیار برابر ۲ استفاده کنیم.

مقدار بحرانی برای تست درآمد را با ۰.۰۵ در نظر میگیریم و تست را با این فرضیات انجام می‌دهیم.

پس از آن با توجه به اینکه خود مسئله مقدار انحراف معیار را نداده و مجهول است از t-test استفاده کردیم و نتایج را مجدداً بررسی کردیم.

همچنین فرض کردیم اگر به اندازه ۰.۲ اختلاف داشتیم میزان قدرت در تشخیص را اندازه میگیریم.

میزان قدرت به اندازه ۰.۹۷ در تشخیص می‌باشد که میزان قدرت مناسبی می‌باشد.

### ۴-۵ بازه اطمینان درآمد

در این بخش برای درآمد مشتریان می‌خواهیم بازه اطمینان محاسبه کنیم.

برای اینکار از نمونه‌های ۳۰ تایی استفاده میکنیم که ۱۰۰۰ بار نمونه‌گیری شود و میانگین مربوط به هر نمونه

ذخیره شود. سپس مقدار Z استاندارد را بدست آورده و در فرمول بازه اطمینان جایگذاری میکنیم که فاصله

اطمینان بصورت زیر می‌باشد:

$$(20142.19, 13460.89)$$

## ۵-۵ رابطه میان Default و Client\_Housing\_Type

ادعا شده افرادی که دارای خانه‌ی از نوع rental می باشند احتمال عدم بازپرداخت وام بیشتری در مقایسه با سایر مشتریان دارند.

$$H_0: P_1 < P_2$$

$$H_1: P_1 > P_2$$

$P_1$ : احتمال عدم بازپرداخت وام توسط افراد با خانه ی rental

$P_2$ : احتمال عدم بازپرداخت وام توسط افراد با خانه ی غیر از rental

$N_1$ : تعداد افراد با خانه ی rental

$N_2$ : تعداد افراد با خانه ی غیر از rental

$Win_1$ : مجموع عدم بازپرداخت وام توسط افراد با با خانه ی rental

$Win_2$ : مجموع عدم بازپرداخت وام توسط افراد با با خانه ی غیر از rental

برای انجام این آزمون فرض از proportion test استفاده میکنیم.

پس از انجام آزمون مقدارهای زیر بدست می آید:

$$P_1: 0.12707775$$

$$P_2: 0.08007267$$

با توجه به نتایج بدست آمده فرض صفر رد می شود و افرادی که دارای خانه ی اجاره ای هستند توان بازپرداخت وام کمتری نسبت به سایر افراد دارند.

## ۵-۶ آزمون سن

ادعا شده افرادی که دارای سن بیشتر از ۴۱ سال (۱۵۰۰۰ روز) می باشند با احتمال بیشتری بازپرداخت وام های خود را در مقایسه با افراد دارای سن کمتر از ۴۱ سال انجام می دهند.

$$H_0: P_1 < P_2$$

$$H_1: P_1 > P_2$$

P1: احتمال بازپرداخت وام توسط افراد با سن بیشتر از ۴۱ سال

P2: احتمال بازپرداخت وام توسط افراد با سن کمتر از ۴۱ سال

N1: تعداد افراد با سن بیشتر از ۴۱ سال

N2: تعداد افراد با سن کمتر از ۴۱ سال

Win1: مجموع بازپرداخت وام توسط افراد با سن بیشتر از ۴۱ سال

Win2: مجموع بازپرداخت وام توسط افراد با سن کمتر از ۴۱ سال

برای انجام این آزمون فرض از  $\text{proportion test}$  استفاده میکنیم.

مقدار احتمال بازپرداخت برا یافراد با بیشتر از ۴۱ سال سن برابر ۰.۹۲۸۱۱۵۲ و برای افراد با کمتر از ۴۱ سال سن برابر ۰.۹۰۳۸۲۰۸ می باشد. در نتیجه فرض صفر رد می شود. اما همانطور که از عدد احتمال ها مشخص است تفاوت چشمگیری با یکدیگر ندارند.

## ۷-۵ رابطه میان Client\_City\_Rating و Default

طبق بررسی هایی که تیم این موسسات انجام دادند به این نتیجه رسیدند که اگر سطح شهر زندگی مشتریان متوسط باشد توان بازپرداخت وام آنها کاهش می یابد.

برای آزمون این فرض یکبار سطح شهر متوسط با خوب و یکبار سطح شهر متوسط با عالی را بررسی میکنیم.  
سطح های متوسط و خوب:

$$H_0: P_1 > P_2$$

$$H_1: P_1 < P_2$$

P1: احتمال عدم بازپرداخت وام توسط افراد با سطح شهر متوسط

P2: احتمال بازپرداخت وام توسط افراد با سطح شهر خوب

N1: تعداد افراد با سطح شهر متوسط

N2: تعداد افراد با سطح شهر خوب

Win1: مجموع بازپرداخت وام توسط افراد با سطح شهر متوسط

Win2: مجموع بازپرداخت وام توسط افراد با سطح شهر خوب

برای انجام این آزمون فرض از  $\text{proportion test}$  استفاده میکنیم.

پس از انجام تست احتمال عدم بازپرداخت وام برای افراد در سطح شهر متوسط برابر  $0.05653385$  و برای افراد در سطح شهر خوب برابر  $0.06854134$  می باشد و فرض صفر رد می شود.

سطح های متوسط و عالی:

P1: احتمال عدم بازپرداخت وام توسط افراد با سطح شهر متوسط

P2: احتمال بازپرداخت وام توسط افراد با سطح شهر عالی

N1: تعداد افراد با سطح شهر متوسط

N2: تعداد افراد با سطح شهر عالی

Win1: مجموع بازپرداخت وام توسط افراد با سطح شهر متوسط

Win2: مجموع بازپرداخت وام توسط افراد با سطح شهر عالی

پس از انجام تست احتمال عدم بازپرداخت وام برای افراد در سطح شهر متوسط برابر  $0.05164627$  و برای افراد در سطح شهر خوب برابر  $0.11418387$  می باشد و فرض صفر رد می شود.

## ۶ مدل‌های پیش‌بینی (Predictive Models)

سه مدل Logistic Regression، Decision Tree و Random Forest برای پیش‌بینی در نظر گرفته شد. با توجه به این که دیتاست Test\_Dataset متغیر Default را ندارد، لازم است برای بررسی عملکرد مدل‌های پیش‌بینی از دیتاست Train\_Dataset استفاده شود. بنابراین، دیتاست آماده‌سازی شده به دو بخش train و test تقسیم شد. برای این منظور، نسبت تقسیم (Split Ratio) برابر با ۰.۸ در نظر گرفته شد که در نتیجه آن، ۹۷۴۸۵ داده به دیتاست train و ۲۴۳۷۱ داده به مجموعه تست اختصاص یافت. همچنین در تمامی مدل‌ها متغیرهای ID و Mobile\_Tag از مدل‌سازی حذف شدند. تلاش شده برای آستانه مقداری انتخاب شود که تا حد امکان، مقدار Sensitivity و Specificity با هم برابر شوند، بدین ترتیب شرایط مقایسه مدل‌ها تقریباً برابر خواهد بود.

### ۱-۶ Logistic Regression

این مدل در ابتدا با بکارگیری همه متغیرها به جز ID و Mobile\_Tag برازش شد، سپس متغیرهایی که موثر تشخیص داده نشده بودند، به ترتیب بزرگترین مقدار p-value از مدل حذف شدند. در خروجی مدل نهایی، مقدار  $AIC^{10}$  که می‌توان گفت بیانگر خطای مدل است برابر با ۴۹۷۸۹ مشاهده شد. همچنین، با در نظر گرفتن آستانه برابر با ۰.۰۸، مقادیر accuracy برابر با ۶۶.۷۹٪، sensitivity برابر با ۶۶.۸۱٪ و specificity برابر با ۶۶.۵۳٪ برای داده‌های train بدست آمد.

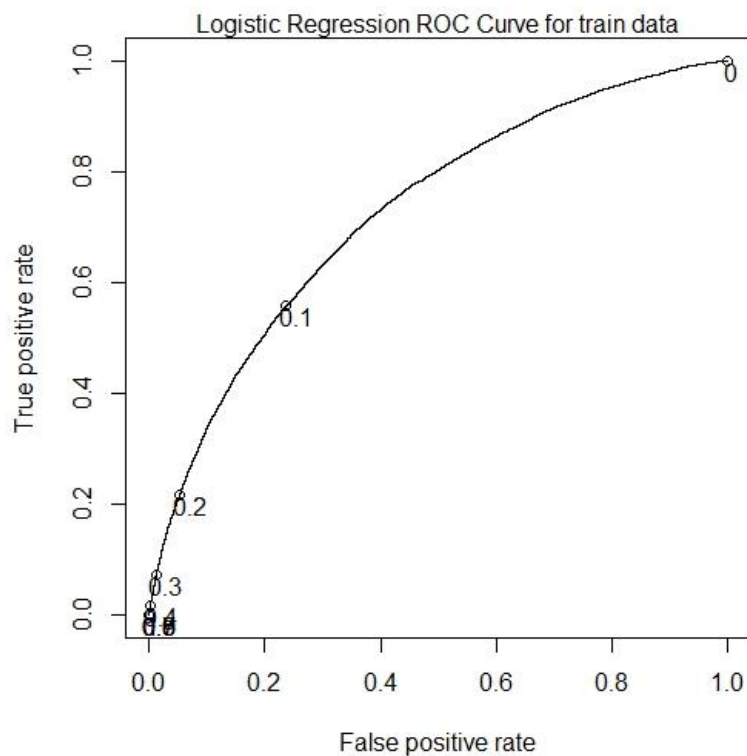
در ادامه، نمودار ROC<sup>۱۱</sup> مدل نهایی برای داده‌های train، رسم شد که در شکل ۶-۱ نشان داده شده است.

با بکارگیری مدل برای داده‌های test و در نظر گرفتن آستانه ۰.۰۸، مقادیر accuracy برابر با ۶۶.۸۶٪، sensitivity برابر با ۶۷.۰۳٪ و specificity برابر با ۶۵.۰۱٪ بدست آمد. نمودار ROC مدل برای داده‌های test در شکل ۶-۲ نشان داده شده است و سطح زیر نمودار ROC (AUC)<sup>۱۲</sup> نیز برابر با ۷۱.۹۳٪ بدست آمده است.

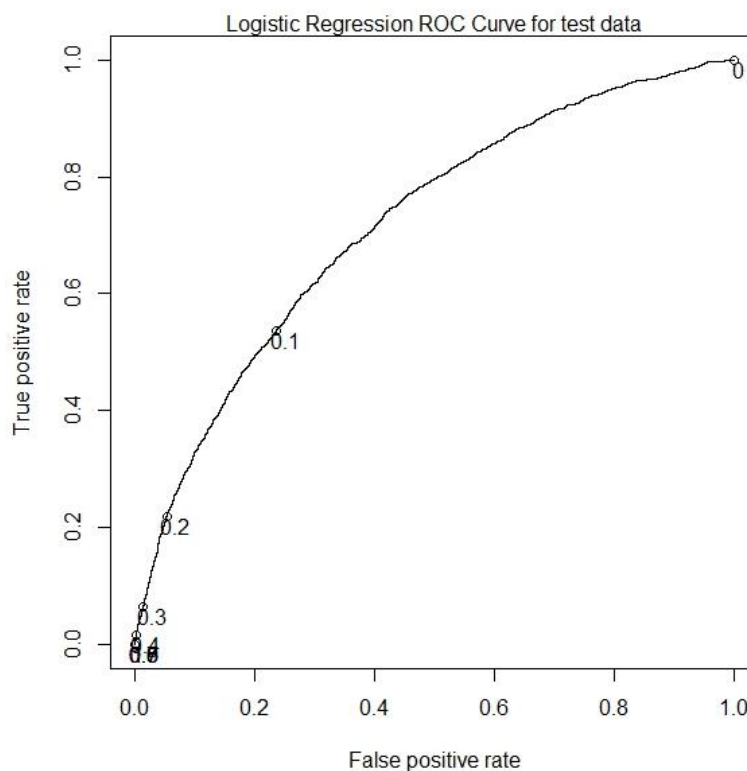
<sup>10</sup> Akaike Information Criterion

<sup>11</sup> Receiver Operating Characteristics

<sup>12</sup> Area Under the Curve



شکل ۶-۱ نمودار ROC مدل Logistic Regression برای داده‌های train

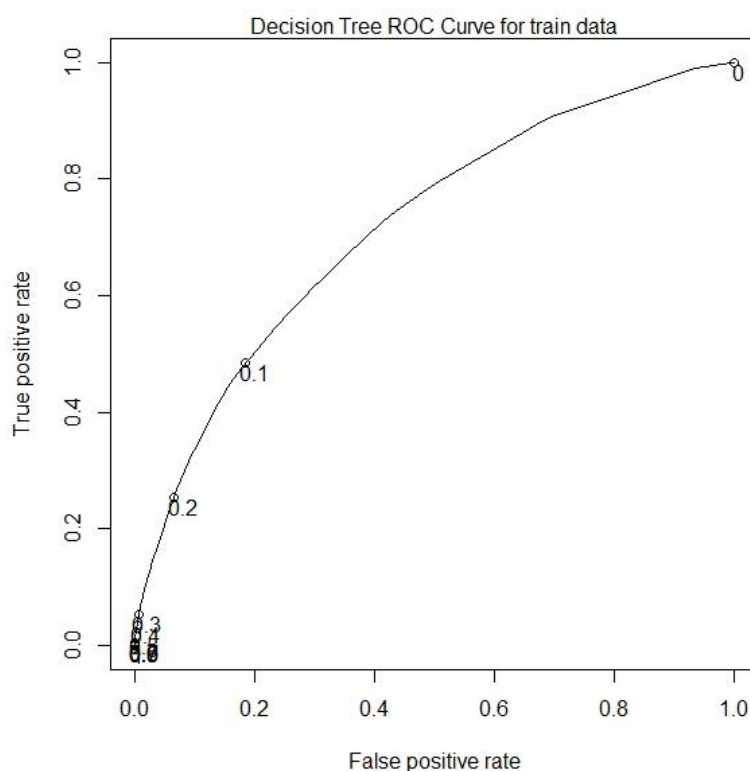


شکل ۶-۲ نمودار ROC مدل Logistic Regression برای داده‌های test

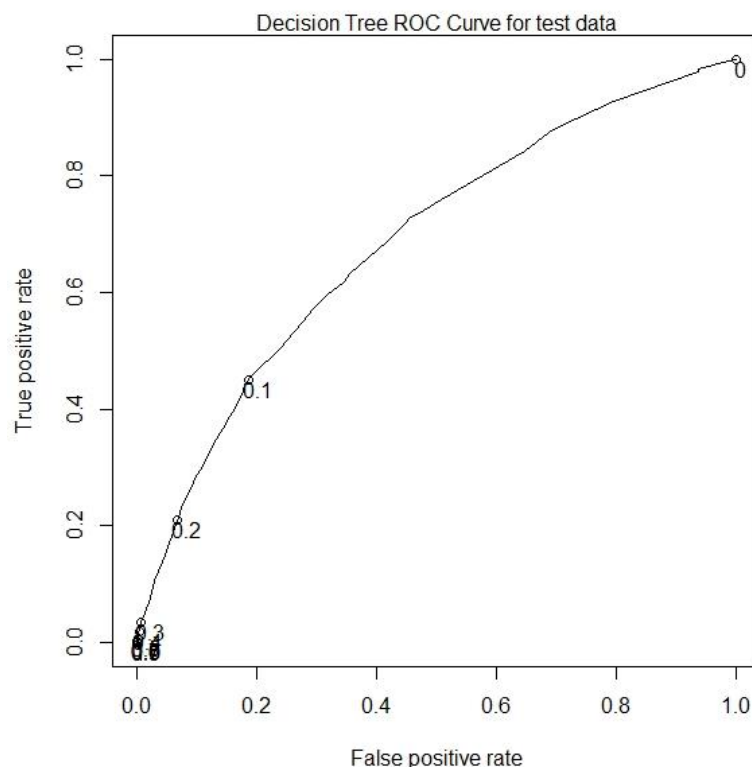
## ۲-۶ Decision Tree

پیش از بکارگیری این مدل داده های عددی استاندارد سازی شدند. در ادامه به دلیل تعداد بالای متغیرها مقدار پارامتر  $cp$  برابر با ۱- قرار داده شد تا مدل بر خلاف پیچیدگی بالا ایجاد شود، در ادامه تعداد شاخه زدن های مدل محدود شد تا یک نتیجه معقول حاصل شود. نتایج در جدول آمده است. همچنین نمودار این مدل با توجه به تعداد زیاد شاخه ها نمایش قابل فهمی ندارد.. با در نظر گرفتن مقدار آستانه برابر ۰.۰۸، معیارهای  $accuracy$  برابر با ۶۵.۷۴٪،  $sensitivity$  برابر با ۶۵.۷۱٪ و  $specificity$  برابر با ۶۶.۰۶٪ برای داده های  $train$  بدست آمد. نمودار ROC مدل برای داده های  $train$  در شکل ۳-۶ رسم شده است.

با بکارگیری مدل برای داده های  $test$ ، معیارهای  $accuracy$  برابر با ۶۵.۰۶٪،  $sensitivity$  برابر با ۶۵.۳۴٪ و  $specificity$  برابر با ۶۱.۸۶٪ بدست آمد. به نظر می رسد، این مدل کارایی معادل با مدل Logistic Regression داشته است. در ادامه، نمودار ROC مدل برای داده های  $test$  رسم شد که در شکل ۴-۶ نشان داده شده است، همچنین مقدار  $AUC$  برابر با ۶۸.۷۰٪ بدست آمد که نسبت به مدل Logistic Regression بدتر است.



شکل ۳-۶ نمودار ROC مدل Decision Tree برای داده های  $train$



شکل ۶-۴ نمودار ROC مدل Decision Tree برای داده‌های test

### ۳-۶ Random Forest

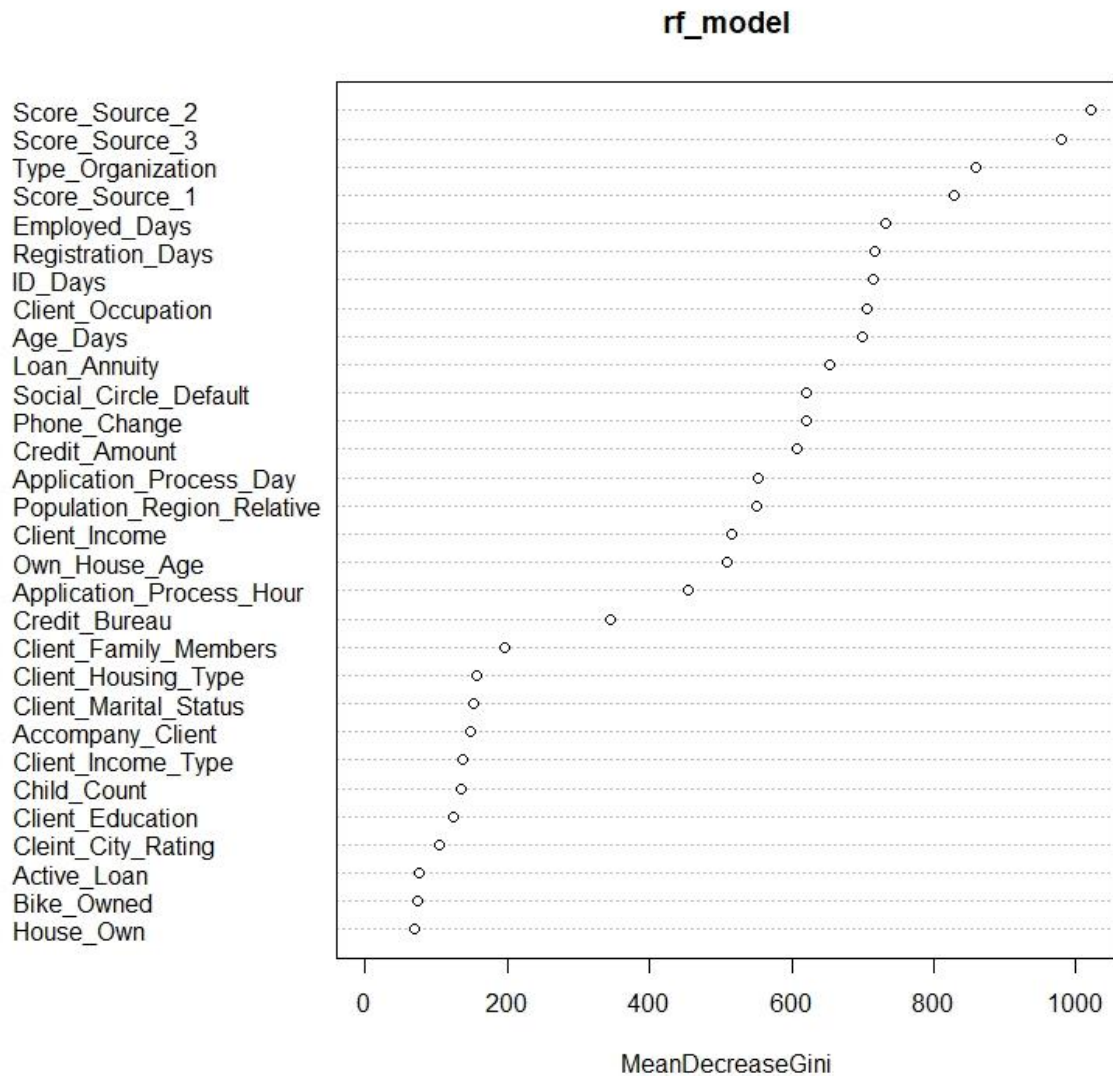
این مدل در ابتدا با بخش بزرگی از متغیرها تشکیل شد. با تشکیل مدل، نمودار اهمیت متغیرها مورد بررسی قرار گرفت. برای بهبود عملکرد مدل، کم اهمیت ترین متغیرها در نمودار اهمیت که در شکل ۶-۵ نشان داده شده است، به ترتیب حذف شدند تا جایی که دیگر پارامتر AUC بهبود نیابد. پس از حذف پنج متغیر روند بهبود مقدار AUC به پایان رسید و مدل نهایی با مقدار خطای OOB<sup>۱۳</sup> برابر با ۷۰.۷۱٪ ایجاد شد. با در نظر گرفتن مقدار آستانه برابر با ۰.۱۱ برای دیتاست train، معیارهای accuracy برابر با ۷۰.۵۷٪، sensitivity برابر با ۷۰.۸۴٪ و specificity برابر با ۶۷.۴۲٪ بدست آمد. همچنین نمودار ROC مدل برای دیتاست train در شکل ۶-۶ رسم شده است.

با بکارگیری مدل برای داده‌های test، معیارهای accuracy برابر با ۷۰.۴۶٪، sensitivity برابر با ۷۰.۷۱٪ و specificity برابر با ۶۷.۷۰٪ بدست آمد. به نظر می‌رسد، کارایی این مدل بر خلاف دو مدل پیشین، بر

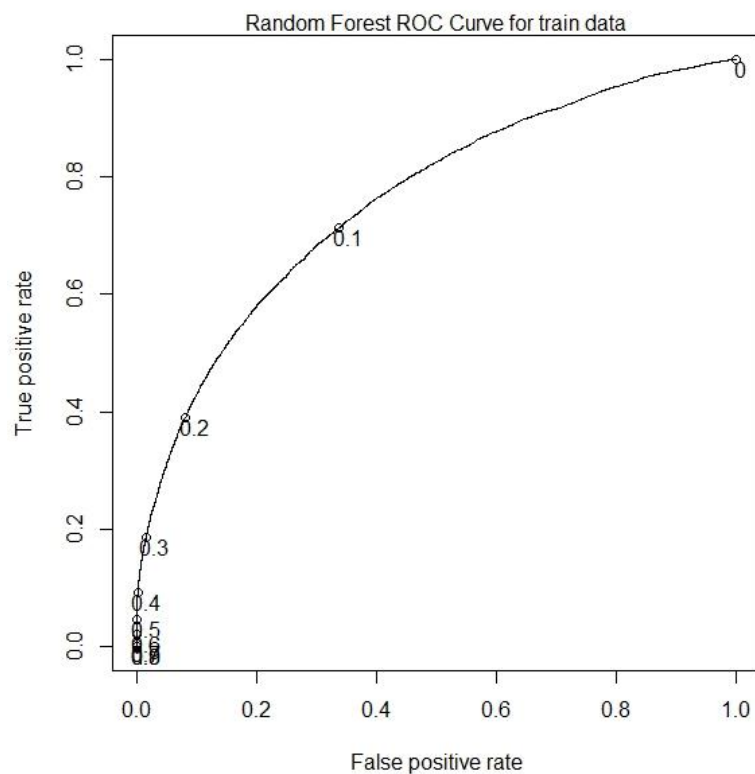
<sup>۱۳</sup> Out Of Bag Error



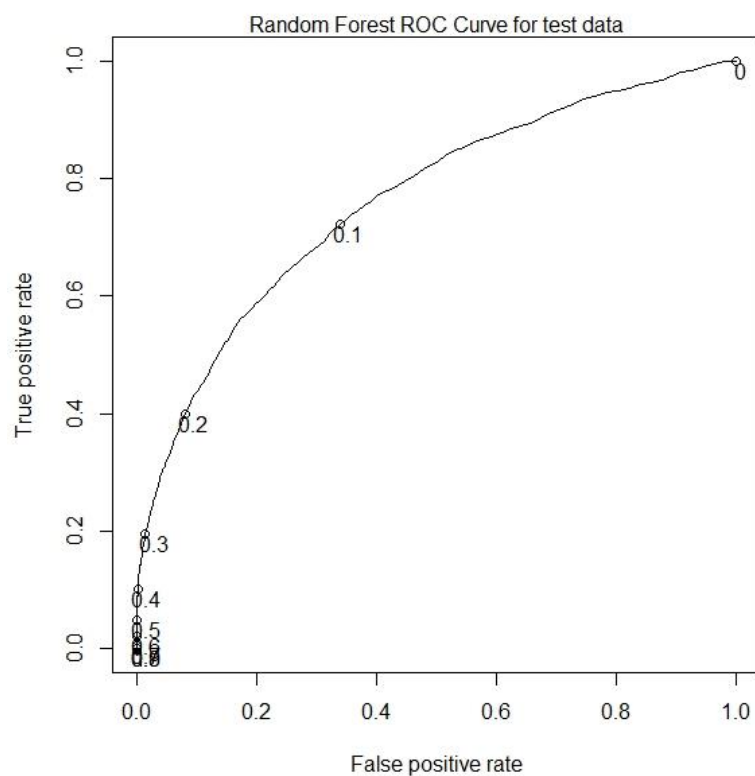
روی داده‌های test بهتر بوده است. در ادامه، نمودار ROC مدل برای داده‌های test رسم شد که در شکل ۶-۷ نشان داده شده است، همچنین مقدار AUC، برابر با ۷۶.۱۸٪ بدست آمده که به وضوح عملکرد بهتری را نسبت به دو مدل پیشین نشان می‌دهد.



شکل ۶-۵ نمودار اهمیت متغیرها در مدل Random Forest



شکل ۶-۶ نمودار ROC مدل Random Forest برای داده‌های train



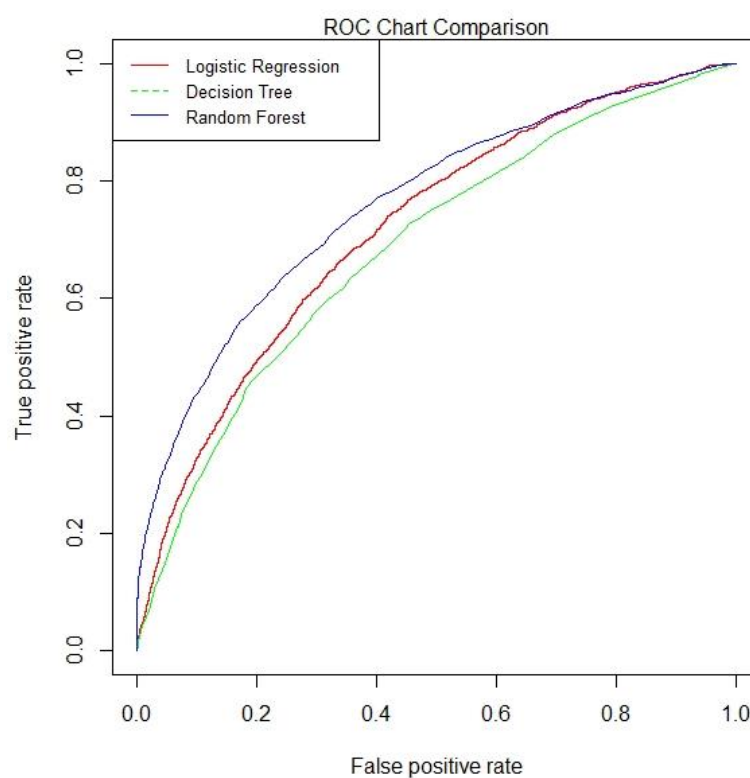
شکل ۶-۷ نمودار ROC مدل Random Forest برای داده‌های test

## ۷ نتیجه گیری

مهم ترین معیارهای مربوط به سه مدل پیش بینی در جدول ۲ آورده شد هاست. همچنین منحنی مربوط به مدل های ارائه شده در مقایسه شده است. در نتیجه می توان گفت مدل Random Forest عملکرد بهتری داشته است.

جدول ۲ - مقایسه معیارهای مدل های پیش بینی ارائه شده

مدل	Accuracy	Sensitivity	Specificity	AUC
Logistic Regression	۶۶.۸۶٪	۶۷.۰۳٪	۶۵.۰۱٪	۷۱.۱۹٪
Decision Tree	۶۵.۰۶٪	۶۵.۳۴٪	۶۱.۸۶٪	۶۸.۷۰٪
<b>Random Forest</b>	<b>۷۰.۴۶٪</b>	<b>۷۰.۷۱٪</b>	<b>۶۷.۷۰٪</b>	<b>۷۶.۱۸٪</b>



شکل ۷-۱ منحنی ROC مدل های ارائه شده

- [1] O. MONK, "Automobile Loan Default Dataset," Kaggle, 2021. [Online]. Available: <https://www.kaggle.com/datasets/saurabhbagchi/dish-network-hackathon>. [Accessed 28 06 2022].
- [2] Anonymous, "Exploratory Data Analysis," JMP Statistical Discovery LLC, [Online]. Available: [https://www.jmp.com/en\\_hk/statistics-knowledge-portal/exploratory-data-analysis.html](https://www.jmp.com/en_hk/statistics-knowledge-portal/exploratory-data-analysis.html). [Accessed 31 07 2022].
- [3] M. Restori, "What is Exploratory Data Analysis," Charito, [Online]. Available: <https://chartio.com/learn/data-analytics/what-is-exploratory-data-analysis/>. [Accessed 31 07 2022].
- [4] N. Tamboli, "All You Need To Know About Different Types Of Missing Data Values And How To Handle It," Analytics Vidhya, 25 07 2022. [Online]. Available: <https://www.analyticsvidhya.com/blog/2021/10/handling-missing-value/>. [Accessed 31 07 2022].
- [5] P. Bhandari, "Missing Data | Types, Explanation, & Imputation," Scribbr, 08 12 2021. [Online]. Available: <https://www.scribbr.com/statistics/missing-data/>. [Accessed 31 07 2022].
- [6] Anonymous, "Tutorial on 5 Powerful R Packages used for imputing missing values," Analytics Vidhya, 05 07 2020. [Online]. Available: <https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#:~:text=MICE%20Package,of%20uncertainty%20in%20missing%20v alues..> [Accessed 31 07 2022].
- [7] Anonymous, "mice: mice: Multivariate Imputation by Chained Equations," RDocumentation, [Online]. Available: <https://www.rdocumentation.org/packages/mice/versions/3.14.0/topics/mice>. [Accessed 31 07 2022].
- [8] J. Josse, "Handling missing values with R," Julie Josse, [Online]. Available: <http://juliejosse.com/wp-content/uploads/2018/06/DataAnalysisMissingR.html>. [Accessed 31 07 2022].

- [9] S. Buuren, " Predictive mean matching," Stevan Buuren, [Online]. Available: <https://stefvanbuuren.name/fimd/sec-pmm.html>. [Accessed 31 07 2022].
- [10] Anonymous, "Extensions to Multinomial Regression," Columbia Public Health, 07 07 2022. [Online]. Available: [https://www.publichealth.columbia.edu/research/population-health-methods/extensions-multinomial-regression#:~:text=Multinomial%20\(Polytomous\)%20Logistic%20Regression&text=In%20polytomous%20logistic%20regression%20analysis,outcome%20is%20compared%20to%20it](https://www.publichealth.columbia.edu/research/population-health-methods/extensions-multinomial-regression#:~:text=Multinomial%20(Polytomous)%20Logistic%20Regression&text=In%20polytomous%20logistic%20regression%20analysis,outcome%20is%20compared%20to%20it). [Accessed 31 07 2022].