

Paper 3603-2018

Getting to Know the No-Show: Predictive Modeling of Missing a Medical Appointment

Joe Lorenz and Kayla Hawkins, Grand Valley State University

ABSTRACT

Patients not showing up for appointments, without cancelling, are still a major loss of efficiency for healthcare clinics, and reduce patient face-time with healthcare providers. This in turn leads to a loss of revenue for healthcare clinics, and overbooking of patient appointments. Using a dataset of 110,000 Brazilian medical appointments collected over a 3-month period, different healthcare populations were identified based on age and frequency of healthcare use. Once healthcare populations were selected, stepwise logistic regression was performed, using PROC LOGISTIC, in order to identify factors within and across populations that are relevant indicators of no-show appointments. The data was split for cross validation, with approximately 80% allocated to build the fitted models, and 20% reserved for validation. The ROC curves for the training and validation datasets were found to be very similar and ranged from 0.68 to 0.75, showing good diagnostic capability for their respective models. The accuracy of the prediction model was also evaluated on the validation dataset, and the model correctly categorized no-show appointments 76-82% of the time. The appointment being on the same day as scheduled, receiving a reminder text, age of the patient, adulthood status and the presence of a chronic condition all contributed to the likelihood of an appointment being a no-show.

INTRODUCTION

With the many demands for personalized universal healthcare, it is surprising that such a large portion of the population is not getting their individual healthcare needs met. In fact, some recent U.S. studies found that one third of adults avoid going to the doctor's all together (Taber, 2015). Having such a large proportion of adults shying away from doctors, it is not a surprise that for some clinics up to 42% of all appointments result in patient no-shows (Lacy, 2005). The prevalence of patient no-shows have become such a standard issue in medical practice that there have been many studies focused on how to optimally schedule appointments when accounting for no-shows (Huang, 2014; Liu, 2016). The loss of revenue, time, and efficiency that clinics incur from holding appointment blocks that are not attended has caused them to look for solutions, most commonly overbooking (Huang, 2014; LaGanga & Lawrence, 2007; Liu, 2016). In fact, instead of trying to prevent no-shows many physicians have adopted a method of overbooking patients to help make up for the loss of revenue (Zeng, 2009). While this does help decrease the loss of revenue, and does increase the overall efficiency of the healthcare clinic, there is some evidence that overbooking negatively affects the patients experience. The increased waiting time that patients experience because of overbooking causes them to become less satisfied with their healthcare experience, and may discourage patients from scheduling appointments in the future (LaGanga & Lawrence, 2007). Instead of treating the adverse effects of patient no-shows through overbooking, a better understanding of the factors that increase the likelihood of a patient no-show are needed. The knowledge of what causes healthcare users to schedule appointments and not show up could be used to develop and implement strategies to reduce the odds of patients not attending their appointments.

DATA

Kaggle Inc. houses a large warehouse of publicly available datasets, including a dataset containing information on scheduled healthcare appointments in Brazil. While multiple versions of the dataset were available, the most recent version, updated May 2016, was used. The dataset contained roughly 110,000 scheduled appointment observations from various clinics, in the state of Espirito Santo, over a 3-month period. There were 15 variables in the dataset including the characteristics of the visit,

such as unique patient and appointment identification numbers, the neighborhood of the participating healthcare clinic, date of scheduling, the appointment date, and the binary response variable of whether or not the patient was a no-show. Being a no-show meant the patient did not attend the appointment without calling ahead to cancel, appointments that were cancelled were not included in the dataset. In addition, demographic variables pertaining to the patient were present including age, gender, and whether the patient has any history of alcoholism, hypertension, diabetes, and disability. Whether or not an SMS text messages was sent, as a reminder of the appointment date, was also included in the dataset. The final variable was whether the patient was a part of the Bolsa Familia social welfare program, where recipients received financial aid to help fight and reduce poverty (Lindert et al., 2007). Using the provided appointment information, new variables were derived including whether the appointment was on the day it was scheduled, whether the patient was an adult, what day of the week the appointment was on, and the time between the appointment scheduling and the actual appointment. This information was then used to see if there were different populations of healthcare users.

THE PROBLEM

No-show patients have such a large effect on a clinic's day to day practice by causing a loss in efficiency, over-scheduling to compensate, and a decrease in physician-patient face time. Therefore, it is important to understand why no-shows happen, especially in order to implement solutions which could help combat no-show appointments. To better diagnose the nature of patients who no-show, exploratory analysis was done to identify different populations of healthcare users. Once established the separate populations, were then analyzed in order to determine what unique factors exist within these patient groups, and what factors exist across all healthcare users that are relevant in predicting no-shows. Since the response variable, whether or not an appointment was a no-show, is a binary outcome, binary logistic regression was used to identify factors resulting in no-shows.

DATA CLEANING & VALIDATION

While the majority of the data was appropriately coded and reflected the methods described for the variables in the data dictionary, validation checks were done to guarantee the quality of the data. This included ensuring all variables that were described as binary-coded had only two possible responses. 199 observations were deleted, because they were outside of the binary, and considered data entry errors. In addition, five appointments were scheduled to take place prior to the day they were on, therefore they were considered data entry errors and deleted. Lastly the scheduled date of the appointment and the date the appointment was scheduled were used to create a day of the week variable using the INTICK function. This was used to see if there were any cyclical patterns in no-shows, as well as used to derive the length between when an appointment was scheduled and when the appointment occurred on. However, because only 36 of the 110,000 observations were on a Saturday these were not used for analysis. The variables patient id, age, and neighborhood were the first considered when defining different populations of healthcare users.

ANALYSIS

POPULATION IDENTIFICATION

Exploratory analysis of patient age, using PROC MEANS, showed two distinct age groups within the dataset; adults and children. Figure 1 showed that the percentage of no-shows clearly increases as the patient reaches adulthood (see Appendix A for Figure 1). However, when adulthood is reached the percentage of no-shows drop and continue dropping as patients get older. Due to the differing healthcare needs of adults and children this was not surprising, and a variable was created indicating if the patient was an adult. In Brazil legal adulthood is considered 18 years old (*Brazil: Statute of the Child and*

Adolescent). In addition to using age to identify populations, the idea of whether patients occurred multiple times in the 3 month period or only appeared once was of interest. The two age groups were further subset into whether or not the patient only had one scheduled appointment or if the patient had multiple appointments. Since patients never occurred at multiple clinics, the neighborhood variable did not end up being utilized to categorize populations. In total 4 populations were identified; patients under 18 that appeared once, patients over 18 that appeared once, patients under 18 that appeared multiple times, and patients over 18 that appeared multiple times. The percentage of no-shows for the different categorical variables were then produced using PROC FREQ, for each population. The most interesting finding from these frequencies, found in Table 1, was the different proportion of no-show appointments for the four populations when the appointment was on that day (see Appendix A for Table 1). These populations were then analyzed separately to try and identify the most relevant predictors of no-shows.

STEPWISE LOGISTIC REGRESSION

In order to ensure proper data was being allocated to create and validate the model, a simple random sample of the observations was selected, with 80% of each patient population designated to build the model and 20% reserved for validation. A separate model was then fit for each of the four populations, using PROC LOGISTIC to perform logistic regression. The forward-stepwise logistic regression entry criteria was a significance level of 0.15, as this is the asymptotic equivalence of AIC (Shtatland et al., 2001). The models were built to predict the probability of an appointment resulting in a no-show. Table 2 reports the area under the ROC curves of the fitted models (see Appendix A for Table 2). These range from 0.68 to 0.75, meaning the models shows good diagnostic capability. The final models resulting from the stepwise regression can be seen in Table 3 (see Appendix A for Table 3).

ODDS RATIO OF THE SELECTED MODELS

Point estimates of the odds ratios, as well as 95% confidence intervals, were reported in order to understand whether the variables present in the final model contributed to an increased or decreased likelihood of a no-show. For categorical variables the odds ratios can be interpreted as the change in odds of a no-show for a selected category, when compared to the reference category. For continuous variables the odds ratios can be interpreted as the multiplicative change in the odds of a no-show, when the continuous variable has increased by one unit.

CROSS-VALIDATION

Once a model was fit for each analysis populations, the remaining 20% of the data was used to validate it. The fitted models were ran on the validation data, and the area under the ROC curve can be seen in Table 2 (see Appendix A for Table 2). The fitted models performed similar between the data that was used to build it and the data that was withheld. The validation data was also utilized to investigate prediction accuracy of the original model for each subpopulation. Using the model from the training data, PROC LOGISTIC was used to obtain estimated probabilities for each of the appointments in the validation dataset. Examples of the estimated probabilities for a variety of patients can be found in Table 4, where an estimated probability of at least 0.5 was used to classify each test set observation as a no-show (see Appendix A for Table 4). The fraction of correct classification, shown in Table 2, exceeded the ROC value for each model and showed that these models would be able to practically discriminate the no-shows of new patients at a fairly high rate (See Appendix A for Table 2).

RESULTS

The fitted models for each of the populations were unique, indicating that the groups of healthcare users do have different factors which contribute to whether or not an appointment results in a no-show. The largest difference between single appointment and multiple appointment health care users

was the magnitude of the effect of the appointment being on the day it was scheduled. If the visit was not on the day of scheduling, patients with one appointment were 16.6 to 24.5 times more likely to miss the appointment than if visit was on the scheduling day, based on adulthood group. Patients who had multiple appointments had a much smaller odds of no-showing, which were 4.7 to 9.6 times more likely to no-show in the same situation. While it is obvious that a patient is more likely to attend an appointment scheduled on the same day, since it was probably made because of an unforeseen illness, it is surprising that it resulted in a much larger effect for patients that had fewer appointments. The reason for this could be that patients who do not schedule frequent appointments might only go to the physician if something is wrong. Due to this they may not schedule check-up visits, which may be overlooked since they are scheduled months in advance. Another factor that is present in all models, but unique across two populations, is age.

Remembering Figure 1, age was the variable that had the most unique distribution of percentage of no-show appointments (see Appendix A for Figure 1). The odds of no-showing for the youth populations increase multiplicatively by 1.02 to 1.03 times, for each year older they are. This is very different from the adult population, which is 0.98 to 0.99 times less likely to be a no-show, multiplicatively for every additional year. There could be many reasons that the odds switch from being an increased odds to decreased odds when age population changes. For younger patients it is the more likely that they are not scheduling appointments themselves, and that a guardian is reminding or bringing the child to the visit, lowering the odds of the appointment being missed. However, the new found freedom that comes as a child ages, may find the patient being responsible for scheduling and attending their own appointments, and while adjusting to this new responsibility forget that they scheduled an appointment at all. The decrease of odds of missing an appointment when adults get older, could be due to the increase in the number of medical appointments for the elderly, and the seriousness of not regularly monitoring health symptoms. This may cause older adults to keep closer track of their healthcare appointments. While there are many differences in the models for the four populations, there are a few similarities that exist in all analysis models.

Looking at commonalities in the odds ratios across the patient populations, there are many variables that have similar effects in each of the models. As mentioned above, the appointment being on the day it was scheduled increased the likelihood of the visit being attended in every population. This makes sense, because appointments scheduled immediately are more likely to be remembered than check-ups scheduled a year in advance. Not receiving an SMS text message reminder of the appointment also caused patients to be 1.07 to 1.18 times more likely to no-show than if they had received one. This implies that SMS reminders are an effective method of reaching patients in all populations. In the 3 relevant models, patients who did not receive financial aid from the Brazilian government were 0.76 to 0.80 times less likely to no-show than those who do. An explanation of this could be that those who are using this scholarship program might be less able to afford taking off work and no-show an appointment because of it. Patients without chronic conditions, like a history of alcoholism or diabetes, were also 0.18 to 0.88 less likely to be no-shows than those with these conditions. These findings are consistent with current research showing that patients with diabetes are more likely to no-show, and the high recidivism rates of alcoholics may cause patients to no-show (Nuti et al., 2012; Sinha, 2001). Hypertension was the only chronic condition for which patients with the condition were less likely to no-show, but like many variables it was unique to individual patient populations.

While many of the resulting odds ratios were not surprising, there were other variables that had results which were less intuitive. Hypertension performed in the opposite direction of all other chronic conditions in the model it appeared in, finding that hypertensive patients were less likely to no-show than non-hypertensive patients. Disability was a significant predictor in 2 models; adults with one appointment and adolescents with multiple appointments. However, this variable indicated that adults without a disability are 1.34 times more likely to have no-show appointments than those with a disability. Where as, children with disabilities were 0.34 times less likely to have a no-show appointment than those without.

While there is no clear reason for this switch from one population to another, the reason it is not significant in the other two populations could be that there was not a large enough sample of patients with disabilities present. Adults over 18 with one appointment also had gender as a significant variable, where women were 0.92 time less likely to miss an appointment than men. This could be due to adult women having more healthcare appointments than men, and because of that being more likely to remember appointments scheduled in advance (National Center for Health Statistics, 2010). Day of the week was also a variable that was significant in all populations, except minors with one appointment. For most cases Friday appointments were more likely to be no-shows, however Monday and Tuesday appointments were more likely to be skipped by adults with multiple appointments.

SUGGESTIONS

While there was a lot of useful information that can be used to help make informed decisions about future studies, there are limitations to this analysis. The data was only collected over a three month period. However, a 12 month data collection period would help bring insights into whether or not there are any seasonal patterns in the number of appointments that result in no-shows. This would also introduce more visits that occur on Saturdays, and hopefully allow those visits to be included in future calculations. The variable for date scheduled also included a time that the appointment was scheduled on, but the appointment day variable did not include a time the appointment was supposed to be on. This variable could be helpful in the future for investigating whether or not the appointment was in the morning or the afternoon, and if it is a relevant factor for predicting a no-show. In addition, information on whether a patient is a new patient to a clinic or an already established patient would be useful. A wider inclusion of chronic conditions, including what type of disability a patient has, could help identify populations of patients most in need of frequent health care visits that are at risk of no-showing appointments. Finally, if a reason for the visit, such as emergency, new condition development, or if the appointment was a regularly scheduled check-up, was collected for each appointment this could provide some insight as to why patients no-show. Healthcare providers are often looking to expand their practices, and identifying whether a potential patient has high odds of a no-show would provide insight as how to gain and retain new patients. The most important area for future research would be testing methods to reach patients who are at high-risk of being no-shows, perhaps starting with the most at risk population, which was identified as young adults.

CONCLUSION

While there are many ways to gather more complete data of what causes health care users to skip appointments, there are still many noteworthy insights into no-shows that were revealed in this analysis. Utilizing SAS®, see Table 5 for a full list of SAS® procedures and functions used, relevant predictors were identified for unique populations of healthcare patients (see Appendix A for Table 5). These findings could help healthcare clinics perform more efficiently and reach patients that are at-risk of being no-shows. Some results were intuitive, and showed that current methods to reach patients, such as receiving a SMS reminder, reduces the odds of a no-show appointment occurring. Chronically ill patients being more likely to no-show supports that these patients are not being reached and could be in danger of serious complications due to not monitoring their health. In addition, it was interesting that age plays different roles within the populations of patients that were identified, with youngest and oldest patients being less likely to no-show and those from 18 to middle age having higher odds of being no-shows. Methods designed specifically to reach this young adult age group could not only help to improve clinical efficiency, but may help these patients attend doctor visits more often and take charge of their healthcare. All of these factors are important for clinical decision making, in order to reach out to populations of patients with higher odds of no-shows, and should be further investigated.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Robert Downer for his wisdom and guidance in this endeavor.

REFERENCES

- Brazil: *Statute of the Child and Adolescent [Brazil]*, Law n° 8.069, 13 July 1990, available at: <http://www.refworld.org/docid/4c481bcf2.html>
- Huang, Y., & Zuniga, P. (2014). Effective cancellation policy to reduce the negative impact of patient no-show. *Journal of the Operational Research Society*, 65(5), 605–615. <https://doi.org/10.1057/jors.2013.1>
- Lacy, N. L. (2004). Why We Don't Come: Patient Perceptions on No-Shows. *The Annals of Family Medicine*, 2(6), 541–545. <https://doi.org/10.1370/afm.123>
- LaGanga, L. R., & Lawrence, S. R. (2007). Clinic Overbooking to Improve Patient Access and Increase Provider Productivity. *Decision Sciences*, 38(2), 251–276. <https://doi.org/10.1111/j.1540-5915.2007.00158.x>
- Lindert, K., Linder, A., Hobbs, J., & De la Brière, B. (2007). *The nuts and bolts of Brazil's Bolsa Familia Program: implementing conditional cash transfers in a decentralized context* (Vol. 709). Social Protection Discussion Paper.
- Liu, N. (2016). Optimal Choice for Appointment Scheduling Window under Patient No-Show Behavior. *Production and Operations Management*, 25(1), 128–142. <https://doi.org/10.1111/poms.12401>
- National Center for Health Statistics. (2010, January 22). New Study Profiles Women's Use of Health \Care. Retrieved January 15, 2018, from <https://www.cdc.gov/nchs/pressroom/01news/newstudy.htm>
- Nuti, L. A., Lawley, M., Turkcan, A., Tian, Z., Zhang, L., Chang, K., ... Sands, L. P. (2012). No-shows to primary care appointments: subsequent acute care utilization among diabetic patients. *BMC Health Services Research*, 12(1). <https://doi.org/10.1186/1472-6963-12-304>
- Sinha, R. (2001). How does stress increase risk of drug abuse and relapse? *Psychopharmacology*, 158(4), 343–359. <https://doi.org/10.1007/s002130100917>
- Shtatland, E. S. , Cain, E. , & Barton, M. B. (2001). The perils of stepwise logistic regression and how to escape using them using information criteria and the output delivery system. In *Proceedings of the 13th Annual SAS Users Group International (SUGI), USA* (pp. 222-226). Cary, NC: SAS Institute.
- Taber, J. M., Leyva, B., & Persoskie, A. (2015). Why do People Avoid Medical Care? A Qualitative Study Using National Data. *Journal of General Internal Medicine*, 30(3), 290–297. <https://doi.org/10.1007/s11606-014-3089-1>
- World Bank. (nd.). Life expectancy at birth. Retrieved from <https://data.worldbank.org/indicator/SP.DYN.LE00.IN>
- Zeng, B., Turkcan, A., Lin, J., & Lawley, M. (2010). Clinic scheduling models with overbooking for patients with heterogeneous no-show probabilities. *Annals of Operations Research*, 178(1), 121–144. <https://doi.org/10.1007/s10479-009-0569-5>

APPENDIX A. TABLES & FIGURES

Table 1. Percentage of No-Show by Population

Population	One Appointment	One Appointment	Multiple Appointments	Multiple Appointments
	< 18 years old	≥ 18 years old	< 18 years old	≥ 18 years old
Alcoholism (None)	20.00 (28.57)	18.18 (23.11)	23.00 (58.52)	20.34 (18.87)
Diabetes (None)	44.44 (19.99)	15.67 (18.60)	23.02 (28.57)	19.02 (20.40)
Hypertension (None)	30.77 (20.00)	14.77 (19.66)	36 (23.02)	18.54 (20.87)
Disability (None)	22.94 (20.00)	12.82 (18.44)	38.19 (22.91)	17.52 (20.34)
SMS Reminder (None)	15.70 (30.05)	14.99 (24.59)	20.06 (30.37)	16.60 (28.01)
Same Day Visit (Future Visit)	2.21 (32.55)	2.60 (25.82)	9.82 (32.11)	4.23 (28.24)
Scholarship (None)	20.00 (20.01)	22.27 (17.99)	26.68 (22.57)	24.09 (19.84)
Male	20.24	18.45	22.68	19.33
Female	19.79	18.26	23.38	20.65
Monday	21.00	19.04	21.73	21.00
Tuesday	18.05	16.70	22.84	21.28
Wednesday	19.34	17.81	24.42	19.17
Thursday	21.46	17.59	20.99	19.42
Friday	20.98	21.15	24.71	20.26

Note: Values in parenthesis represent the no-show proportion without this condition. E.g. adolescents with one appointment who have a history of alcoholism had no-show appointments for 20% of all appointments, where as those who did not have a history of alcoholism in the same population had no-shows for 28.57% of all appointments.

Table 2. Cross Validation Using ROC Curves

Population	Age Group	Training ROC	Validation ROC	Validation Prediction Accuracy*
Appear Once	< 18	0.7518	0.7564	79.05%
Appear Once	≥ 18	0.7450	0.7503	81.88%
Appear Multiple Times	< 18	0.6822	0.6893	76.35%
Appear Multiple Times	≥ 18	0.7072	0.7030	79.54%

*Validation prediction accuracy is the proportion of training data predicted accurately with the reported model.

Table 3. Odds Ratio Table for Significant Effects

Variable	Patients Appear Only Once		Patients Appear Multiple Time	
	Patient Age < 18	Patient Age ≥ 18	Patient Age < 18	Patient Age ≥ 18
Age	1.02 (1.01, 1.03)	0.98 (0.97, 0.98)	1.03 (1.03, 1.04)	0.99 (0.98, 0.99)
No Alcoholism	NA	0.60 (0.50, 0.73)	0.18 (0.06, 0.53)	0.80 (0.71, 0.92)
No Diabetes	NA	0.85 (0.74, 0.98)	NA	0.88 (0.80, 0.96)
Female vs Male	NA	0.92 (0.85, 1.00)	NA	NA
Not Handicap	NA	1.34 (0.99, 1.83)	0.34 (0.23, 0.51)	NA
No Hypertension	NA	1.11 (1.00, 1.23)	NA	NA
Same Day Appointment	24.53 (19.07, 31.56)	16.62 (14.19, 19.47)	4.67 (4.18, 5.22)	9.55 (8.74, 10.43)
SMS Received	1.21 (1.07, 1.37)	1.18 (1.10, 1.28)	1.16 (1.06, 1.28)	1.07 (1.02, 1.13)
No Scholarship	NA	0.76 (0.67, 0.87)	0.80 (0.70, 0.91)	0.79 (0.73, 0.86)
Mon vs Fri	NA	0.96 (0.85, 1.07)	0.88 (0.77, 1.02)	1.04 (0.96, 1.13)
Tues vs Fri	NA	0.87 (0.77, 0.97)	0.94 (0.92, 1.20)	1.08 (1.00, 1.17)
Weds vs Fri	NA	0.84 (0.75, 0.94)	0.94 (0.82, 1.07)	0.89 (0.82, 0.96)
Thurs vs Fri	NA	0.80 (0.71, 0.91)	0.82 (0.71, 0.95)	0.92 (0.85, 1.00)

Note: NA indicates the variable was not included in the model for this population. The category included in the model is given in parentheses, with 0 indicating an absence of the condition.

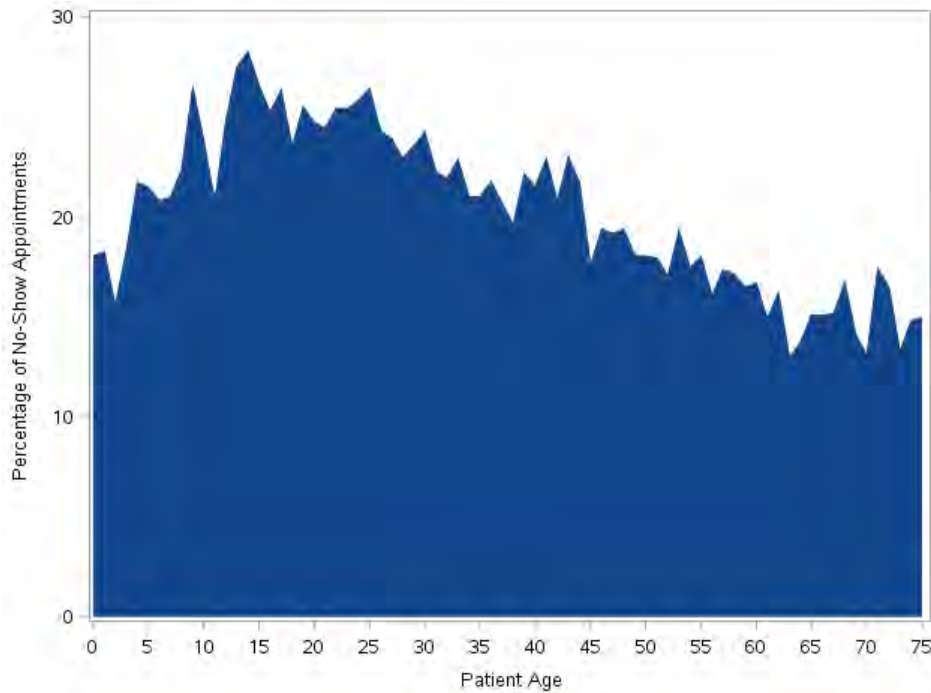
Table 4. Predicted Probabilities for a Cross Section of Patients

Number of Visits	Age	Gender	Same Day Appointment	Scholarship	SMS Reminder	Medical History	Day of the Week	Probability of No-Show
1	6	Male	Yes	No	Yes	None	Tuesday	0.0203
4	7	Female	Yes	No	Yes	None	Wednesday	0.1011
3	16	Female	No	Yes	No	Alcoholism	Friday	0.8223
1	25	Male	Yes	No	No	Hypertension	Monday	0.0531
1	31	Male	No	Yes	No	Diabetes	Friday	0.5792
3	55	Female	Yes	No	Yes	Hypertension	Thursday	0.2250

Table 5. Summary of SAS® Functions and Procedures

Analysis Step	Procedure/Function
Data Import	PROC IMPORT
Data Cleaning	Done within DATA steps and PROC SQL
Date Cleaning	WEEKDAY(), DATEPART(), MONTH()
Date Validation	INTCK()
Exploratory Data Analysis	PROC MEANS(), PROC FREQ(), PROC SGPLOT()
Data Analysis	%MACRO to repeat analysis for each population and model
Training and Validation Data Allocation	PROC SURVEYSELECT
Logistic Regression	PROC LOGISTIC with a SELECTION = STEPWISE, SLENTY = 0.15, PLOTS=ROC
Validation	PROC FREQ, PROC LOGISTIC with PLOTS=ROC

Figure 1. Percentage of Visits Resulting in No-Shows Across Patient Ages



Note: The graph plots the percentage of appointments that are no-shows at each patient age, reported in years. 75 years old was considered a cutoff since it was the average life expectancy of Brazil in 2015 (World Bank, nd.).

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Authors: Joe Lorenz & Kayla Hawkins
University: Grand Valley State University
Advised By: Dr. Robert Downer

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.