

# Pronóstico de Graduación Estudiantes Subgraduados Universidad de Puerto Rico Recinto Mayagüez

Ramineh Lopez Yazdani<sup>1</sup> & Jesús D. Hernández Londoño<sup>2</sup>

<sup>1</sup>Universidad de Puerto Rico. Departamento de Matemáticas.  
Recinto Mayagüez.

<sup>2</sup>Universidad de Puerto Rico. Departamento de Matemáticas.  
Recinto Mayagüez.

9 de septiembre de 2020



# Contenido

- 1 Objetivo
- 2 Antecedentes
- 3 Metodología
- 4 Referencias



# Contenido

- 1 Objetivo
- 2 Antecedentes
- 3 Metodología
- 4 Referencias



## Objetivo

Elegir un Modelo de Machine Learning, para la predicción de graduación de estudiantes subgraduados a tiempo o en general (150%).



# Contenido

- 1 Objetivo
- 2 Antecedentes
- 3 Metodología
- 4 Referencias



## Tasa de Graduación

Es la proporción de estudiantes que completan sus estudios dentro del 150 % (en general) del tiempo del programa o a tiempo, i.e., 8 años para programas de 5 años (Facultad de Ingeniería) y 6 años para los demás programas, excepto grados asociados.



## ¿Por qué graduarse a tiempo?

Según "Complete College America Alliance of States"

- Grados Asociados pierden aproximadamente \$35,000 en salarios.
- Programas de 4 años \$45,327 en salarios.
- Prestigio y Acreditación

## Causas

- Consejería académica insuficiente o inadecuada.
- Grupos grandes, ausentismo frecuente.
- Deficiencias en las materias básicas.
- Pocos cursos.



## Estrategias

- Establecer un sistema de alerta temprana.
- Consejería académica compuestos por estudiantes.

## Tasa de Graduación UPRM

Según OPIMI en su informe del 2019, el 51,4 % de los estudiantes con cohorte en 2010, habría completado su grado al cabo del 150 % del tiempo.





Average Graduation rates of full-time, undergraduates within 150% of normal time to program completion according to IPEDS: Cohorts entry year 2011 for PR and 2012 for USA

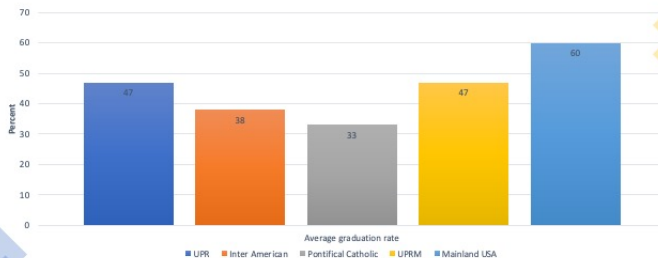


Figura: Autoría propia



# Contenido

- 1 Objetivo
- 2 Antecedentes
- 3 Metodología**
- 4 Referencias



- Se inicia con un preprocesamiento, extracción y creación de variables, eliminación de datos duplicados, eliminación Na's
- Exploración de los datos.
- Comparar Modelos de Machine learning para validar predicciones. (Logistic Regression, Random Forests, Boosting, Support Vector Machines, Deep Learning).



## Software

- Rstudio
- tidyverse, dplyr, rpart, partykit, ROCR, gbm, RandomForest, caret, e1071, nnet.



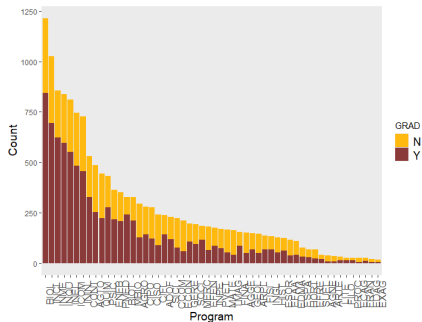
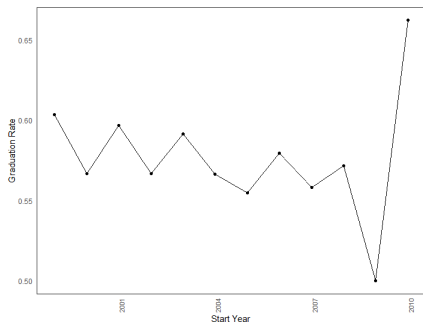
## Data

- Oficina de planeación, investigación y mejoramiento institucional (OPIMI), (43,515 ; 22), (43,515 ; 11), (18,955 ; 8), años de ingreso 1999-2018.
- Observaciones 18413, 21 variables.
- Cohorte 2010.

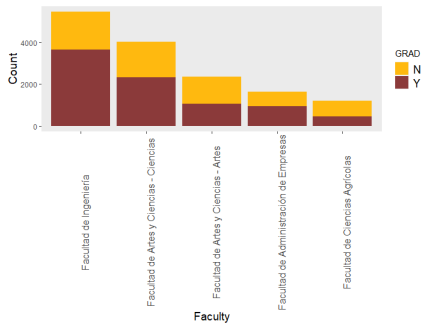
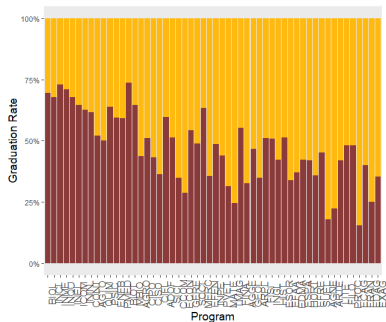
## Variables

STUDENT RECORD KEY, Year, Faculty, Major, Apt Verbal, Aprox Matem, Apt Matem, Aprox Espanol, Aprox Ingles, Highschool GPA, INGRESO FAMILIAR, EDUC PADRE, EDUC MADRE, Gender, School Type, GPA 1ER ANO, GRAD, Year Grad, Rel Stud GPA, Rel School GPA, grad intime

# Metodología



# Metodología



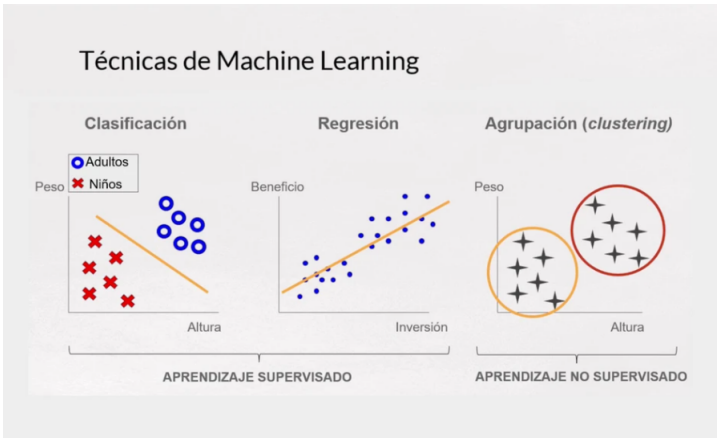


Figura: Extraída de OpenWebinars





- ➊ Clasificadores de Bayes
  - Linear Discriminant Analysis.
  - Naive Bayes.
  - k-nearest neighbors.
  - Logistic Regression.
- ➋ Decision Trees.
- ➌ Ensembles.
- ➍ Neural Networks and Deep Learning.
- ➎ Support vector machines.



## Matriz de confusión

Es una herramienta que permite conocer el desempeño de un algoritmo.

	Actual class=Yes	Actual class=No
Predicted class =Yes	True Positives=TP	False Positives=FP
Predicted class=No	False Negative=FN	True Negatives=TN

Figura: Acuña E.

- $\text{Accuracy} = (\text{TP} + \text{TN}) / \text{Total}$ .
- $\text{Missclassification rate} = (\text{FP} + \text{FN}) / \text{Total} = 1 - \text{Accuracy}$ .
- $\text{True Positives Rate} = \text{Sensitivity} = \text{Recall} = \text{TP} / (\text{TP} + \text{FN})$ .
- $\text{False Positives Rate} = \text{FP} / (\text{FP} + \text{TN})$ .
- $\text{Specifity} = \text{TN} / (\text{FP} + \text{TN}) = 1 - \text{False Positive Rate}$



## ROC

Es el punto de corte en el que se alcanza la sensibilidad y especificidad más alta. Evalúa la capacidad del modelo para clasificar correctamente. Un valor de 1 significa que el modelo es perfecto, un valor de 0.5 indica que el modelo no es útil.

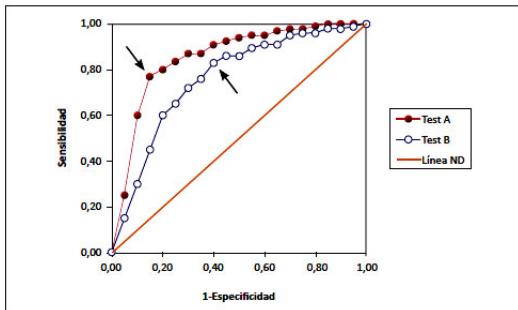


Figura: Jaime Cerda y Lorena Cifuentes.



## AUC

Es el área bajo la curva, Esta refleja que tan bueno es la prueba o modelo para discriminar pacientes con y sin la enfermedad.



# Contenido

- 1 Objetivo
- 2 Antecedentes
- 3 Metodología
- 4 Referencias





Tibshirani R., Hastie T., Witten D., James G.  
*An Introduction to Statistical Learning with applications in R.*



Acuña E.  
*Data Mining and Machine Learning.*  
Decision Trees



Retención de Estudiantes.  
<http://ponce.inter.edu/>.  
Rodríguez A.





University of Puerto Rico-Mayaguez.

<https://nces.ed.gov/ipeds>.

IPEDS (Integrated Postsecondary Education Data System)



Cerda J., Cifuentes L.

*Uso de curvas ROC en investigación clínica. Aspectos teórico-prácticos.*

Rev. chil. infectol. vol.29 no.2 Santiago abr. 2012



Complete College America Alliance of States.

<https://completecollege.org/>.

