

# Machine Learning Engineer Nanodegree

## Capstone Project

Sanjeev Yadav

17 March 2018

## I. Definition

### Project Overview

Supervised learning is one of the most popular areas of machine learning in which much development has already taken place. In this project we are trying to identify the university-level factors which predict the presence of a strong retention and graduation rate. As the leader of the big data revolution, Google gathers information through clicks on the Internet and uses this information to personalize advertising to individual users<sup>[1]</sup>.

The link to the data source is [here](#). The name of the file is data.csv. Data was collected from [data.gov](#), but for ease of access we have downloaded it and pushed it to this GitHub repo.

### Problem Statement

One of the most pressing issues facing American universities is the number of students who fail to graduate. Nearly one out of five four-year institutions graduate fewer than one-third of its first-time, full-time degree-seeking first-year students within six years. Although there are various explanations for attrition, we will try to identify the most important feature which affects the retention and graduation rates in 4-year institutions.

We have two target variables:

1. Graduation rate, and
2. Retention rate

Both are continuous variable so this is a regression task. We will train same regression models for both target variables but the final model will be chosen based on the `r2_score`. It may be the case that one model works good for graduation rate and some other works good for retention rate.

### Metrics

We will use `r2_score` as the metric for performance of our model. `feature_importances_` method from sklearn will show the importance of each feature in predicting the target variables.

## II. Analysis

## Data Exploration

Name of the input data file is data.csv. It has 7593 observations and 123 variables.

Information about all the variables can be seen in metadata.xlsx file. Let us discuss the variable in our input data.

Features:

**UNITID:** Unit ID for institution

**OPEID:** 8-digit OPEID ID for institution

**OPEID6:** 6-digit OPEID for institution

**INSTNM:** Institution name

**CITY:** city

**STABBR:** State postcode

**INSTURL:** URL for instution's homepage

**NPCURL:** URL for institution's net price calculator

**HCM2:** Schools that are on Heightened Cash Monitoring 2 by the Department of Education

**PREDDEG:** Predominant undergraduate degree awarded. Can take 5 values:

1. Not classified
2. Predominantly certificate-degree granting
3. Predominantly associate's-degree granting
4. Predominantly bachelor's-degree granting
5. Entirely graduate-degree granting

**HIGHDEG:** Highest degree awarded. Can take 5 values:

1. Non-degree-granting
2. Certificate degree
3. Associate degree
4. Bachelor's degree
5. Graduate degree

**CONTROL:** Control of institution. Can take 3 values:

1. Public
2. Private non-profit
3. Private for-profit

**LOCALE:** Locale of institution. Can take 12 values:

1. City: Large (population of 250,000 or more)
2. City: Midsize (population of at least 100,000 but less than 250,000)
3. City: Small (population less than 100,000)
4. Suburb: Large (outside principal city, in urbanized area with population of 250,000 or more)
5. Suburb: Midsize (outside principal city, in urbanized area with population of at least 100,000 but less than 250,000)
6. Suburb: Small (outside principal city, in urbanized area with population less than 100,000)
7. Town: Fringe (in urban cluster up to 10 miles from an urbanized area)
8. Town: Distant (in urban cluster more than 10 miles and up to 35 miles from an urbanized area)
9. Town: Remote (in urban cluster more than 35 miles from an urbanized area)
10. Rural: Fringe (rural territory up to 5 miles from an urbanized area or up to 2.5 miles from an urban cluster)
11. Rural: Distant (rural territory more than 5 miles but up to 25 miles from an urbanized area or more than 2.5 and up to 10 miles from an urban cluster)
12. Rural: Remote (rural territory more than 25 miles from an urbanized area and more than 10 miles from an urban cluster)

**HBCU:** Flag for historically Black College and University.

**PBI:** Flag for predominantly black institution.

**ANNHI:** Flag for Alaska Native Native Hawaiian serving institution.

**TRIBAL:** Flag for tribal college and university

**AANAPII:** Flag for Asian American Native American Pacific Islander-serving institution

**HSI:** Flag for Hispanic-serving institution

**NANTI:** Flag for Native American non-tribal institution

**MENONLY:** Flag for men-only college

**WOMENONLY:** Flag for women-only college

**RELAFFIL:** Religious affiliation of the institution. It can take 65 values:

1. Not reported
2. Not applicable
3. American Evangelical Lutheran Church
4. African Methodist Episcopal Zion Church
5. Assemblies of God Church
6. Brethren Church
7. Roman Catholic
8. Wisconsin Evangelical Lutheran Synod
9. Christ and Missionary Alliance Church
10. Christian Reformed Church
11. Evangelical Congregational Church
12. Evangelical Covenant Church of America
13. Evangelical Free Church of America
14. Evangelical Lutheran Church

15. International United Pentecostal Church
16. Free Will Baptist Church
17. Interdenominational
18. Mennonite Brethren Church
19. Moravian Church
20. North American Baptist
21. Pentecostal Holiness Church
22. Christian Churches and Churches of Christ
23. Reformed Church in America
24. Episcopal Church, Reformed
25. African Methodist Episcopal
26. American Baptist
27. American Lutheran
28. Baptist
29. Christian Methodist Episcopal
30. Church of God
31. Church of Brethren
32. Church of the Nazarene
33. Cumberland Presbyterian
34. Christian Church (Disciples of Christ)
35. Free Methodist
36. Friends
37. Presbyterian Church (USA)
38. Lutheran Church in America
39. Lutheran Church - Missouri Synod
40. Mennonite Church
41. United Methodist
42. Protestant Episcopal
43. Churches of Christ
44. Southern Baptist
45. United Church of Christ
46. Protestant, not specified
47. Multiple Protestant Denomination
48. Other Protestant
49. Jewish
50. Reformed Presbyterian Church
51. United Brethren Church
52. Missionary Church Inc
53. Undenominational
54. Wesleyan
55. Greek Orthodox
56. Russian Orthodox
57. Unitarian Universalist
58. Latter Day Saints (Mormon Church)
59. Seventh Day Adventists
60. The Presbyterian Church in America
61. Other (none of the above)
62. Original Free Will Baptist
63. Ecumenical Christian
64. Evangelical Christian
65. Presbyterian

**SATVR25:** 25th percentile of SAT scores at the institution (critical reading)

**SATVR75:** 75th percentile of SAT scores at the institution (critical reading)

**SATMT25:** 25th percentile of SAT scores at the institution (math)

**SATMT75:** 75th percentile of SAT scores at the institution (math)

**SATWR25:** 25th percentile of SAT scores at the institution (writing)

**SATWR75:** 75th percentile of SAT scores at the institution (writing)

**SATVRMID:** Midpoint of SAT scores at the institution (critical reading)

**SATMTMID:** Midpoint of SAT scores at the institution (math)

**SATWRMID:** Midpoint of SAT scores at the institution (writing)

**ACTCM25:** 25th percentile of the ACT cumulative score

**ACTCM75:** 75th percentile of the ACT cumulative score

**ACTEN25:** 25th percentile of the ACT English score

**ACTEN75:** 75th percentile of the ACT English score

**ACTMT25:** 25th percentile of the ACT math score

**ACTMT75:** 75th percentile of the ACT math score

**ACTWR25:** 25th percentile of the ACT writing score

**ACTWR75:** 75th percentile of the ACT writing score

**ACTCMMID:** Midpoint of the ACT cumulative score

**ACTENMID:** Midpoint of the ACT English score

**ACTMTMID:** Midpoint of the ACT math score

**ACTWRMID:** Midpoint of the ACT writing score

**SAT\_AVG:** Average SAT equivalent score of students admitted

**SAT\_AVG\_ALL:** Average SAT equivalent score of students admitted for all campuses rolled up to the 6-digit OPE ID

**PCIP01:** Percentage of degrees awarded in Agriculture, Agriculture Operations, And Related Sciences.

**PCIP03:** Percentage of degrees awarded in Natural Resources And Conservation.

**PCIP04:** Percentage of degrees awarded in Architecture And Related Services.

**PCIP05:** Percentage of degrees awarded in Area, Ethnic, Cultural, Gender, And Group Studies.

**PCIP09:** Percentage of degrees awarded in Communication, Journalism, And Related Programs.

**PCIP10:** Percentage of degrees awarded in Communications Technologies/Technicians And Support Services.

**PCIP11:** Percentage of degrees awarded in Computer And Information Sciences And Support Services.

**PCIP12:** Percentage of degrees awarded in Personal And Culinary Services.

**PCIP13:** Percentage of degrees awarded in Education.

**PCIP14:** Percentage of degrees awarded in Engineering.

**PCIP15:** Percentage of degrees awarded in Engineering Technologies And Engineering-Related Fields.

**PCIP16:** Percentage of degrees awarded in Foreign Languages, Literatures, And Linguistics.

**PCIP19:** Percentage of degrees awarded in Family And Consumer Sciences/Human Sciences.

**PCIP22:** Percentage of degrees awarded in Legal Professions And Studies.

**PCIP23:** Percentage of degrees awarded in English Language And Literature/Letters.

**PCIP24:** Percentage of degrees awarded in Liberal Arts And Sciences, General Studies And Humanities.

**PCIP25:** Percentage of degrees awarded in Library Science.

**PCIP26:** Percentage of degrees awarded in Biological And Biomedical Sciences.

**PCIP27:** Percentage of degrees awarded in Mathematics And Statistics.

**PCIP29:** Percentage of degrees awarded in Military Technologies And Applied Sciences.

**PCIP30:** Percentage of degrees awarded in Multi/Interdisciplinary Studies.

**PCIP31:** Percentage of degrees awarded in Parks, Recreation, Leisure, And Fitness Studies.

**PCIP38:** Percentage of degrees awarded in Philosophy And Religious Studies.

**PCIP39:** Percentage of degrees awarded in Theology And Religious Vocations.

**PCIP40:** Percentage of degrees awarded in Physical Sciences.

**PCIP41:** Percentage of degrees awarded in Science Technologies/Technicians.

**PCIP42:** Percentage of degrees awarded in Psychology.

**PCIP43:** Percentage of degrees awarded in Homeland Security, Law Enforcement, Firefighting And Related Protective Services.

**PCIP44:** Percentage of degrees awarded in Public Administration And Social Service Professions.

**PCIP45:** Percentage of degrees awarded in Social Sciences.

**PCIP46:** Percentage of degrees awarded in Construction Trades.

**PCIP47:** Percentage of degrees awarded in Mechanic And Repair Technologies/Technicians.

**PCIP48:** Percentage of degrees awarded in Precision Production.

**PCIP49:** Percentage of degrees awarded in Transportation And Materials Moving.

**PCIP50:** Percentage of degrees awarded in Visual And Performing Arts.

**PCIP51:** Percentage of degrees awarded in Health Professions And Related Programs.

**PCIP52:** Percentage of degrees awarded in Business, Management, Marketing, And Related Support Services.

**PCIP54:** Percentage of degrees awarded in History.

**DISTANCEONLY:** Flag for distance-education-only education

**UGDS:** Enrollment of undergraduate certificate/degree-seeking students

**UGDS\_WHITE:** Total share of enrollment of undergraduate degree-seeking students who are white

**UGDS\_BLACK:** Total share of enrollment of undergraduate degree-seeking students who are black

**UGDS\_HISP:** Total share of enrollment of undergraduate degree-seeking students who are Hispanic

**UGDS\_ASIAN:** Total share of enrollment of undergraduate degree-seeking students who are Asian

**UGDS\_AIAN:** Total share of enrollment of undergraduate degree-seeking students who are American Indian/Alaska Native

**UGDS\_NHPI:** Total share of enrollment of undergraduate degree-seeking students who are Native Hawaiian/Pacific Islander

**UGDS\_2MOR:** Total share of enrollment of undergraduate degree-seeking students who are two or

more races

**UGDS\_NRA:** Total share of enrollment of undergraduate degree-seeking students who are non-resident aliens

**UGDS\_UNKN:** Total share of enrollment of undergraduate degree-seeking students whose race is unknown

**PPTUG\_EF:** Share of undergraduate, degree-/certificate-seeking students who are part-time

**CURROPER:** Flag for currently operating institution, 0=closed, 1=operating

**NPT4\_PUB:** Average net price for Title IV institutions (public institutions)

**NPT4\_PRIV:** Average net price for Title IV institutions (private for-profit and nonprofit institutions)

**NPT41\_PUB:** Average net price for \ \$0-\$30,000 family income (public institutions)

**NPT42\_PUB:** Average net price for \ \$30,001-\$48,000 family income (public institutions)

**NPT43\_PUB:** Average net price for \ \$48,001-\$75,000 family income (public institutions)

**NPT44\_PUB:** Average net price for \ \$75,001-\$110,000 family income (public institutions)

**NPT45\_PUB:** Average net price for \ \$110,000+ family income (public institutions)

**NPT41\_PRIV:** Average net price for \ \$0-\$30,000 family income (private for-profit and nonprofit institutions)

**NPT42\_PRIV:** Average net price for \ \$30,001-\$48,000 family income (private for-profit and nonprofit institutions)

**NPT43\_PRIV:** Average net price for \ \$48,001-\$75,000 family income (private for-profit and nonprofit institutions)

**NPT44\_PRIV:** Average net price for \ \$75,001-\$110,000 family income (private for-profit and nonprofit institutions)

**NPT45\_PRIV:** Average net price for \ \$110,000+ family income (private for-profit and nonprofit institutions)

**PCTPELL:** Percentage of undergraduates who receive a Pell Grant

**PCTFLOAN:** Percent of all undergraduate students receiving a federal student loan

**UG25ABV:** Percentage of undergraduates aged 25 and above

**MD\_EARN\_WNE\_P10:** Median earnings of students working and not enrolled 10 years after entry



**GT\_25K\_P6:** Share of students earning over \$25,000/year (threshold earnings) 6 years after entry

**GRAD\_DEBT\_MDN\_SUPP:** Median debt of completers, suppressed for n=30

**GRAD\_DEBT\_MDN10YR\_SUPP:** Median debt of completers expressed in 10-year monthly payments, suppressed for n=30

**RPY\_3YR\_RT\_SUPP:** 3-year repayment rate, suppressed for n=30

Target variables:

1. For graduation rates we have two variables in our data. Let us see the difference between those two:

#### 1.1. **rate\_suppressed.four\_year**

Completion rate for first-time, full-time students at four-year institutions (150% of expected time to completion) , pooled in two-year rolling averages and suppressed for small n size.

#### 1.2. **rate\_suppressed.lt\_four\_year\_150percent**

Completion rate for first-time, full-time students at less-than-four-year institutions (150% of expected time to completion) , pooled in two-year rolling averages and suppressed for small n size

We will be making predictions for 4-year institutions.

1. For retention rates we have four variables in our data. Let us see the difference between them:

#### 2.1. **retention\_rate.four\_year.full\_time**

First-time, full-time student retention rate at four-year institutions.

#### 2.2. **retention\_rate.lt\_four\_year.full\_time**

First-time, full-time student retention rate at less-than-four-year institutions.

#### 2.3. **retention\_rate.four\_year.part\_time**

First-time, part-time student retention rate at four-year institutions

#### 2.4. **retention\_rate.lt\_four\_year.part\_time**

First-time, part-time student retention rate at four-year institutions

Retention rate is for full-time students and we are making predictions for 4-year institutions. So, our target variable is **retention\_rate.four\_year.full\_time**.

So there are 2 response variables:

1. **rate\_suppressed.four\_year:** This follows a normal distribution.

2. **retention\_rate.four\_year.full\_time**: This follow a left skewed distribution.

## Algorithms and Techniques

We have one benchmark model and 5 other supervised regression models. Below is the explanation of each model:

### 1. Linear Regression(Benchmark Model)

After data preprocessing, we train the input data and the evaluation metric from this model was considered as the benchmark.

A simple linear regression uses only one independent variable, and it describes the relationship between the independent variable and dependent variable(s) as a straight line.

Here we have 2 dependent variable so we train our model 2 times.

### 1. AdaBoost Regressor

AdaBoost, short for Adaptive Boosting, is a machine learning algorithm formulated by Yoav Freund and Robert Schapire, who won the 2003 Gödel Prize for their work. AdaBoost is sensitive to noisy data and outliers. The individual learners can be weak, but as long as the performance of each one is slightly better than random guessing, the final model can be proven to converge to a strong learner.

### 1. Decision Tree Regressor

Decision tree builds regression or classification models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes.

### 1. Extra Trees Regressor

With respect to random forests, the method drops the idea of using bootstrap copies of the learning sample, and instead of trying to find an optimal cut-point for each one of the K randomly chosen features at each node, it selects a cut-point at random.

### 1. Gradient Boosting Regressor

Gradient boosting is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. It builds the model in a stage-wise fashion like other boosting methods do, and it generalizes them by allowing optimization of an arbitrary differentiable loss function.

### 1. Light GBM

Light GBM grows tree vertically while other algorithm grows trees horizontally meaning that Light GBM grows tree leaf-wise while other algorithm grows level-wise. It will choose the leaf with max delta loss to grow. When growing the same leaf, Leaf-wise algorithm can reduce more loss than a

level-wise algorithm.

## 1. Random Forest Regressor

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests correct for decision trees' habit of overfitting to their training set.

## Benchmark

We will consider benchmark model as the initial linear regressor. We will get the `r2_score` from this model. Then we will use other regression models to improve our score.

# III. Methodology

## Data Preprocessing

We have many column which were unique for every observation. We will choose one of them as the identifier variable for each observation.

For SAT and ACT scores we have 25 percentile, 75 percentile and mid-point values. We will only use the midpoint values for calculation to avoid curse of dimensionality.

We have two columns for SAT Average: `SAT_AVG` and `SAT_AVG_ALL`. Since we have removed `OPEID` column, we will be using `SAT_AVG` column because it provides overall stats, rather than averages based on `OPEID`.

We have average net price for public and private institutions. We will remove the average price based on different family income levels and use the overall average net price for public and private institutions. We are removing columns related to family income levels because we don't think that family income is a university level factor. This is a factor based on student level.

Now I am removing the **Categorical Columns ( dtype object )** which have many levels.

I am keeping categorical columns which have less than 10 levels.

`INSTNM`, `INSTURL` & `NPCURL` are identifiers for colleges in offline or online media and are not university level factors which affect education's status.

`STABBR`, `CITY` & `RELAFFIL` columns have been removed because they had too many levels to be considered.

Calling `.info()` method shows data type `"object"` for last few columns. This is because some values in these columns are `"PrivacySuppressed"`.

There are two types of invalid entries here. First is the `Nan` entry and another is `PrivacySuppressed`.

We will first convert `PrivacySuppressed` to null value and then replace all the null values accordingly.

After converting `PrivacySuppressed` to `NaN`, pandas still treats them as *object* data type. Below is a table showing correct data type of these columns(from metadata.xlsx file):

Column	Data type
MD_EARN_WNE_P10	integer
GT_25K_P6	float
GRAD_DEBT_MDN_SUPP	float
GRAD_DEBT_MDN10YR_SUPP	float
RPY_3YR_RT_SUPP	float
C150_L4_POOLED_SUPP	float
C150_4_POOLED_SUPP	float

Converting `MD_EARN_WNE_P10` to integer data type will throw an error because `NaN` cannot be converted to integer.

So we will convert `MD_EARN_WNE_P10` to *float* type.

Now we will handle the null values.

But before applying median imputation, we need to check for some other categorical columns which might be important for analysis. There are some categorical levels which cannot be ignored while building the model. Those columns are:

1. PREDDEG
2. HIGHDEG
3. CONTROL
4. LOCALE(contains 12 levels but I will reduce them to 4 levels)

We have created dummy variable for these columns.

Finally we renamed the variables to be user-friendly.

## Implementation

The evaluation metric for our models was `r2` score. We used the following algorithms in this project:

1. Linear Regression(Benchmark Model)
2. AdaBoosr Regressor
3. Decision Tree Regressof
4. Extra Trees Regressor
5. Gradient Boosting Regressor
6. Light GBM
7. Random Forest Regressor

## Refinement

Our benchmark model was linear regressor. For this model, `r2_score` for graduation rate was 0.44 and for retention rate was 0.25.

The final model which we chose was Light GBM. Initial `r2` scores for this model were 0.51 for graduation rates and 0.12 for retention rates. We applied hyperparameter tuning by using `num_iteration=gbm.best_iteration_` for the model. Here `gbm` is our regressor model.

Final metrics were 0.83 for graduation rates and 0.57 for retention rates.

## IV. Results

### Model Evaluation and Validation

The final model was Light GBM. It was chosen for graduation rates only. But, when we applied hyperparameter tuning then it was the best model for retention rates also.

### Justification

Below is a comparing of our final model of our final model and the benchmark model showing `r2` score for prediction of graduation and retention rates.

Model	graduation rates	retention rates
Linear Regression(Benchmark model)	0.44	0.25
Light GBM(Final model)	0.83	0.58

We do not expect a very high `r2` score because our aim is not to get a very good prediction. Our aim is to see the important features. So, it will be same when we have best model, whatever high be the `r2` score.

## V. Conclusion

### Reflection

We had 7593 observations and 123 variables. The toughest part was feature selection. We have used domain knowledge to remove features which we think are important or are not university level factor.

One interesting aspect was that we had multiple target labels and we have to narrow down by applying certain conditions, like doing the problem for first-time full-time 4-year institutions.

My initial and final model can be used for this problem in general because they provide a good result of the important university level factors affecting education.

### Improvement

Further improvements can be made if we use hyperparameter tuning on the remaining models. Maybe they will give results than our final model.

We have not used XGBoost algorithm here. Maybe it can be used in further iteration to improve the results. I did not use XGBoost because I did not know in detail how it works.

If we consider our final solution as the new benchmark then I think better solutions exist because we have obtained  $r^2$  scores of 0.83 and 0.58 only. There are good chances of improvement.