

Learning Analytics Tools
Professor Ramkumar Rajendran
Department of Educational Technology
Indian Institute of Technology, Bombay
Lecture 44
Linear Regression - Example

In this video, we will see the example of how this linear regression and logistic regression is used.

(Refer Slide Time: 0:26)



The slide is titled "Linear Regression - Application" in a green header. Below the title, there is a bulleted list with one item: "Saqr, M., Fors, U., & Tedre, M. (2017). How learning analytics can early predict under-achieving students in a blended medical education course. *Medical teacher*, 39(7), 757-767." To the right of the text is a small inset video of Professor Ramkumar Rajendran, who is wearing a checkered shirt and has his hands clasped. The bottom of the slide features the IIT Bombay logo on the left, the text "Learning Analytics" in the center, and the number "2" on the right.

So let us look at this paper how learning analytics can early predict underachieving students in blended medical education course. It is hard to find a paper that explains linear regression in much detail in 2017 or recently because linear regression has been used in the education field for a long time.

But since the authors might be using it for the first time and it is in the medical education course, they explained it in detail. So this is a good paper to look at it, to understand our linear regression. But I do not recommend you writing a paper in such detail because

(Refer Slide Time: 1:09)

3

10

2017 KENNEDY CENTER LECTURE SERIES: "THE ART OF SCIENCE EDUCATION"

7/26/2017, 10:52 AM EDT

an early example of building a warning and feedback system for students and teachers using LA principles. It has been reported to have a positive effect on student retention and teacher assumptions of these students (Phillips & Arnold 2010).

LA has been shown to enable effective, automatic tracking of student engagement across the course (MacLachlan & Dawson 2010, 2012, Woolf et al. 2013, Cruz-Benito et al. 2015, Tatematsu et al. 2015, Galloway et al. 2016). The insights generated by LA can be shared by course teachers, academic supervisors, and administrators (Arnold & Farrell 2012, Howard et al. 2016, Bentley et al. 2016). Those insights could help identify students at risk of underachievement where an early intervention can lead to a meaningful change (MacLachlan & Dawson 2010, LeVine et al. 2013, Tatematsu et al. 2015, Galloway et al. 2016, Howard et al. 2016, Bentley et al. 2016). Although individual assessment methods offer this kind of feedback, their results must often come too late for a possible action or a significant intervention (MacLachlan & Dawson 2010, Cruz-Benito et al. 2015).

Education in the healthcare sector is under a lot of pressure to respond efficiently and timely to the rapidly changing scientific, societal, and cultural environment, as well as to keep programs modern and connected to the community it serves (Bilasy & Wastson 2008, Ellaway et al. 2014, Vainio et al. 2014, Li et al. 2016). Another challenge facing medical schools is understanding and potential actions on issues that may be a symptom of present performance in the medical education in selection of students, curricula, teaching methods, assessment or policies (Chaffin et al. 2011). While the problem of attrition in medical education seems to incur a substantial cost, it is still poorly studied and most of the published studies have focused on students' attributes at the point of admission (Chaffin et al. 2011, Jorgensen et al. 2015), which only explained 10% of variance in performance, recent research indicates that LA can significantly improve the predictability of academic performance, and, hence, also the value of the assessment.

The research questions of this study are:

1. Which tracking variables best correlate with performance?
2. To what extent can the analysis of students' online activity be used to predict student grades, and identify the potential risk of a student failing or dropping a course?

Methodology

The analysis of student data followed a standard procedure used in data mining research and analysis (Ewen 2014, Gordon & Woolf 2014, Woolf et al. 2013):

- Acquisitive and recording: Acquiring the data from different sources.
- Preparing the data: Matching and cleaning misaligned data, including incomplete records and appropriate annotation of data types, combining the data into one master table.
- Performing exploratory data analysis (EDA): Exploring data by testing the interrelation between different variables to discover possible correlations, patterns, rules that could help identify the potential predictors. EDA does not require a prior hypothesis in contrast to hypothesis-driven scientific methodology that tests a previously known theory (Friedman & Hogg 2012).
- Building the predictive model: Predicting students' outcomes and identifying at-risk students using appropriate predictive models. The word regression models, as they are among the most common predictive model used in education research at large (Hogg et al. 2003), and in analytics research (Gordon & Woolf 2014, MacLachlan & Dawson 2010, Galloway et al. 2016, Howard et al. 2016), available in most statistical packages, and can be evaluated in several ways (Friedman et al. 2003, Bower et al. 2005). Two types of regression models according to the type of outcome to be predicted:

3

10

Linear Regression.pdf

C:\Users\Nigel\Desktop\UAT\2020-2021\Prof.Sankar Kumar-Lin...-Regression.pdf

Downloaded by [New York University] at 20:40 03 Oct 2020

tify the possible indicators, followed by the prediction of student outcomes, then we try to predict the risk status of the end of course and whether it is possible at mid-course or not.

Correlations

In Table 2, the findings related to correlations are displayed and the most interesting findings were as follows:

Content creation/innovation: There was a positive and significant correlation between the students' final grade and innovation/content creation variables. The most important were total edits or created content ($r(33) = 0.31$, $p < 0.05$), number of edits in the first half of course ($r(33) = 0.3$, $p < 0.05$), total posts initiated by student ($r(33) = 0.26$, $p < 0.05$) and total posts and replies ($r(33) = 0.29$, $p < 0.05$).

Table 1. Detailed description of collected data.

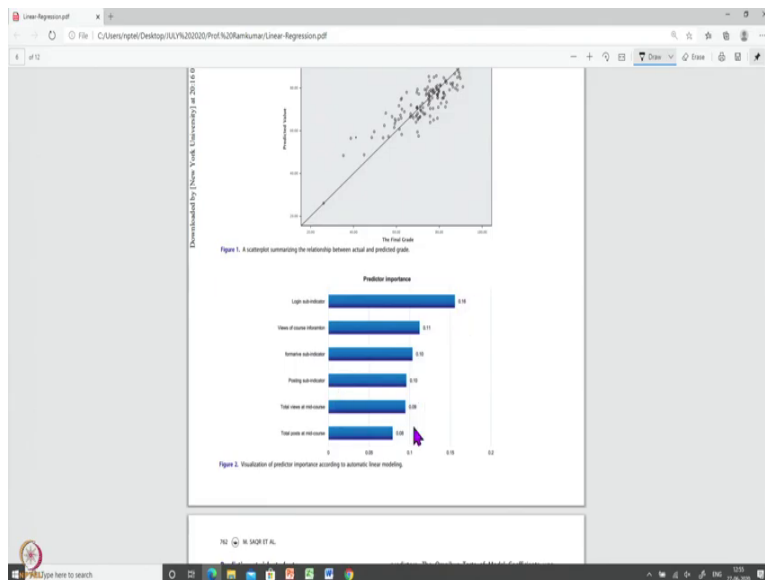
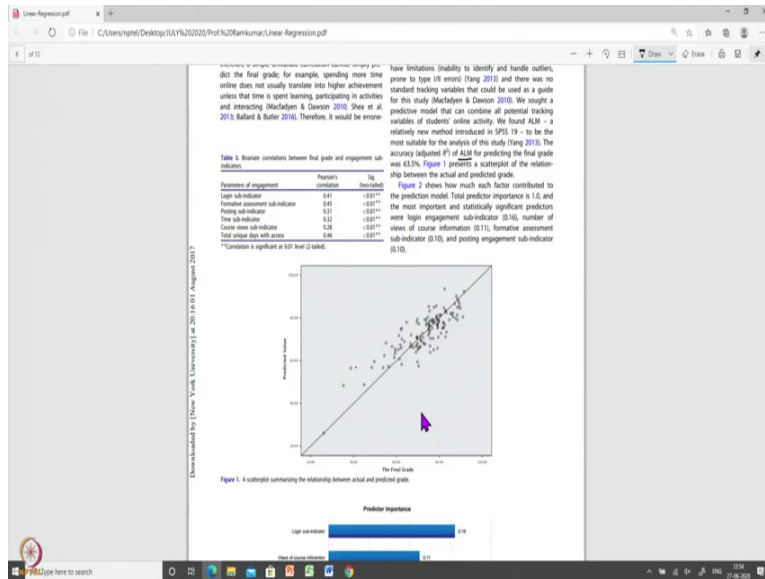
Parameter	Collected data
Inputs	<ul style="list-style-type: none"> Weekly mid-course, and total course logs Logins before and after the end of the course log and password Total number of days with course access Daily number of views, weekly midcourse Total course view, number of unique accesses Accessed and type of the accessed resources
Views (Hz)	Weekly, at mid-course and the overall total of sum-
Outputs	<ul style="list-style-type: none"> Weekly, at mid-course and the overall total of sum- Final grade Results, inputs and total number of edits made in course (number of new content) Interim scores/ranks was calculated by Moodle, i.e. up to the time reflecting the participation of a student was in the form, we obtained Weekly, at mid-course and the overall total view online and educational materials, time taken getting output or viewing non-educational materials, was included Grade of each formative assessment and percentage in assessment regardless of the submission of the answers

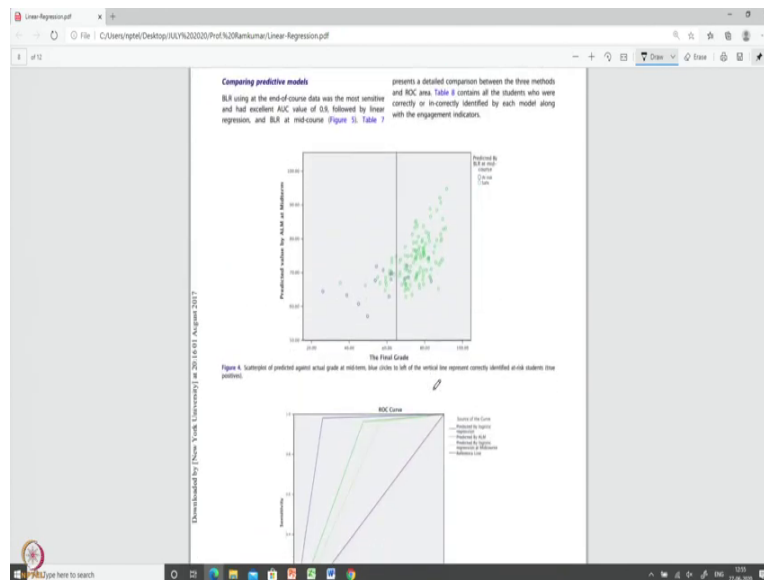
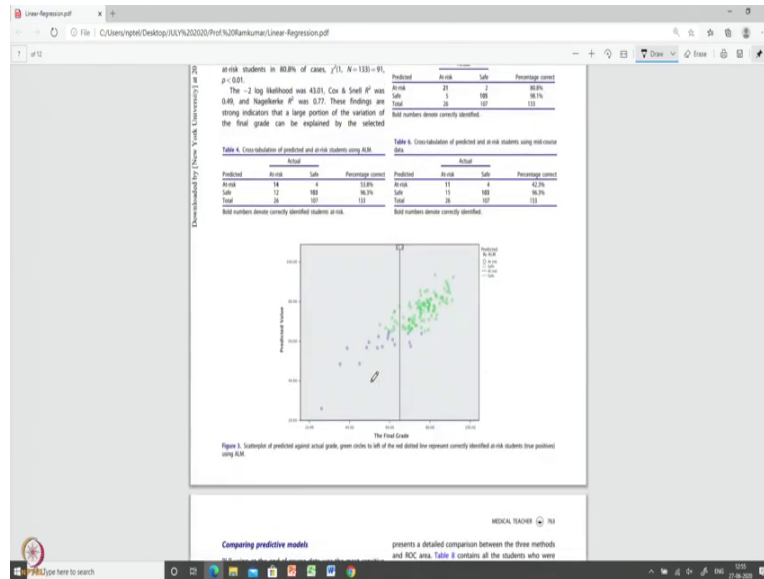
simple scenario or otherwise; especially, nor the parameters that measured student consistency of using the LMS resources (see Table 3). The highest was total unique days of access ($r(33) = 0.46$, $p < 0.05$), followed by the formative assessment sub-indicator ($r(33) = 0.45$) logo sub-indicator ($r(33) = 0.41$, $p < 0.05$), $r = 0.41$ and 0.46 (Table 3).

Table 2. Bivariate correlations between LMS tracking variables and final grade.

Parameter	Correlation	Tig
(X)	(Y)	(Desired level)
Innovation and content creation		
Total posts initiated by a student	0.31**	<0.05
Total posts and replies by a student	0.29**	<0.05
Number of posts in the first half of course	0.3**	<0.05
Number of posts in the second half of course	0.27**	<0.05
Total edits of content created	0.31**	<0.05
Number of edits in the first half of course	0.3**	<0.05
Number of edits in the second half of course	0.22*	<0.05
Days and sites		
Total time from views read	0.27**	<0.05
Total hit on resources	0.28**	<0.05
Total hits on course information	0.22*	<0.05
Total course hits	0.24**	<0.05
Total views before course started	0.037	0.872
Number of hits in the first half of course	0.24**	<0.05
Number of hits in the second half of course	0.26**	<0.05
Number of courses accessed	0.107	0.321
Logins and course access		
Total logins	0.46**	<0.001
Number of logins in the first half of course	0.34**	<0.05
Number of logins in the second half of course	0.34**	<0.05
Time		
Total spent online	0.32**	<0.05
Total time spent online first half of course	0.32**	<0.05
Total time in the second half of course	0.18*	0.054
Formative assessment		
Formative assessment grade	0.45**	<0.001
Mid-course formative assessment grade	0.42**	<0.05
Number attempted the formative quiz	0.41**	<0.05
End of course formative assessment grade	0.32**	<0.05
Communications		
Response to discussions	0.05	0.517
Total times of new communications	0.088	0.544

*Correlation is significant at 0.05 level (2-tailed). **Correlation is significant at 0.01 level (2-tailed).





But this paper is explaining in detail. Let us look at this paper. The paper shows how learning analytics can early predict underachieving students in a blended medical education course. So in this paper, they want to predict that students who are at risk i.e. were going to score below 65 per cent in the final exam. So potentially safe students who have finally scored more than 65 per cent will be considered as pass students.

How to predict the students were going to get less than 65 marks in the final score, that are at risk. So there are 145 students, however, they have excluded many, so they have

only one category of students over the period of six weeks. So what is the data they are collecting, the data they collected is in the blended learning approach, that is, they use MOOC and classroom. The data they collected is log in like weekly, mid-course, total course logins.

How many times he logs in a week, in mid-course how many times he logs in before the mid sem, how many time total login times and login before and after the end course of the exam. Also views, the number of views daily, weekly and mid sem till mid-course or also total course views, number of unique resource accessed.

A number of unique resource means how many resources he looked at it? How many papers, how many videos he watched everything and the data type of this resource he accessed. All the forums, like number of posts created, read or replace, number of edits made in the post, the number of likes and also how influential a particular post he created is, based on the number of people who reply to that post.

And the time he spent on each of these sessions, like weekly meets and over all the time on this particular learning environment also the grades at each formative assessment and participation assessment like assignment, regardless of submission of the answers whether the students submit the answer or not. If he participated in the quiz or answers the pop-up questions in the video, if he answers those questions, is also considered as a formative assessment, not the answer score.

So they took all this and multiple features from this like if you see the login can be classified into total login, number of login in the first half and second half. They computed the correlation with the final score. So this you can compute using the correlation matrix.

So you know how to compute it now or they compute each one individually and also they identify whether this has significant value or not because there are 132 students we know whether this correlation is significant or not, significant is not telling the strength of the correlation, significant tells that whether this correlation is reliable or not.

Significant does not mean this correlation is high or not that is a very important thing we have to understand. So what they do, they see this 0.29, the two stars(asterisk) significant by 0.01, 0.01 level, yes here. So all of them are like average correlation. Not really good maybe this is a good correlation number of times you log in and the assessment is a very good correlation with the final scores.

Yes, as expected assessment is really highly correlated with the final score based on the performance in the quizzes or mid-term assessment grade and formative grades. So what they did, they combined these feature values again this split given to us, six variables. Let us look at it. They made it into six parameters of engagement. See the login indicator so formative assessment, posting time, course views and total uniqueness.

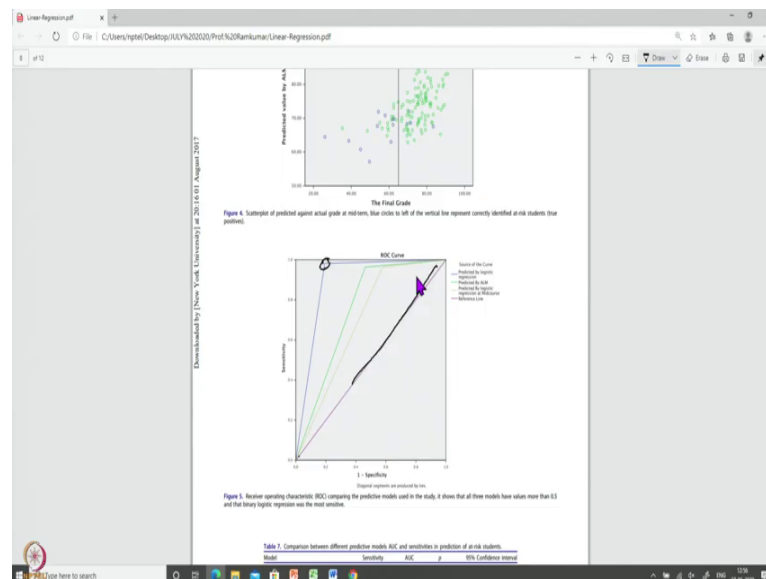
They made into six variables $x_1, x_2, x_3, x_4, x_5, x_6$. They used SPSS software for automated linear modelling that they call this ALM here. But SPSS is proprietary software I do not recommend using that. But if you have access to SPSS 19, please go ahead and use it. But they use SPSS and they are reporting the values here. Let us look at the values. So the predicted final grade is here and the actual final grade is given in this given in these dots but the regression scale is like this.

So this is the difference between the final grade and the predicted value. So also, they are giving the weight of each variable. But as I mentioned, this is the paper that explains every detail about linear regression, which is not needed for the current setting because

everyone knows what is LR. So but yeah see the weight 0.16, 0.1 and formative assessment is 0.10.

So there is a like a weight of each of these variables and the intercept is not given. And that is a problem you know, the intercept is not given. We do not know what is the value but we do not want to interpret the intercept. So we want to interpret only this value. So login is the most strong indicator of the performance and the answer.

(Refer Slide Time: 7:07)



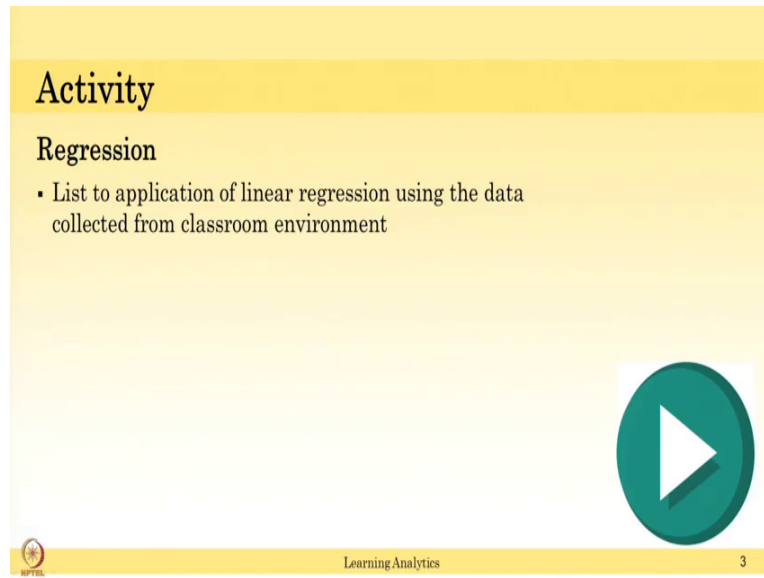
So they computed the students' performance at risk, and they also used the logistic regression to do that, and they also plotted ROC curve, predicted using-

1. logistic regression
2. a linear model
3. logistic regression at the mid-course

Mid-course means that we don't have to wait till the end rather we use data till mid sem to predict their final exam performance.

If you use the mid sem, so you know you will see this. I hope you know what is this line means. This line indicates that the area under the curve is 0.5 anything below this is not good. This particular curve indicates case 3 (i.e. using logistic regression at the mid-course). The graph shows that predicting using logistic regression is the better.

(Refer Slide Time: 8:21)



Activity

Regression

- List to application of linear regression using the data collected from classroom environment

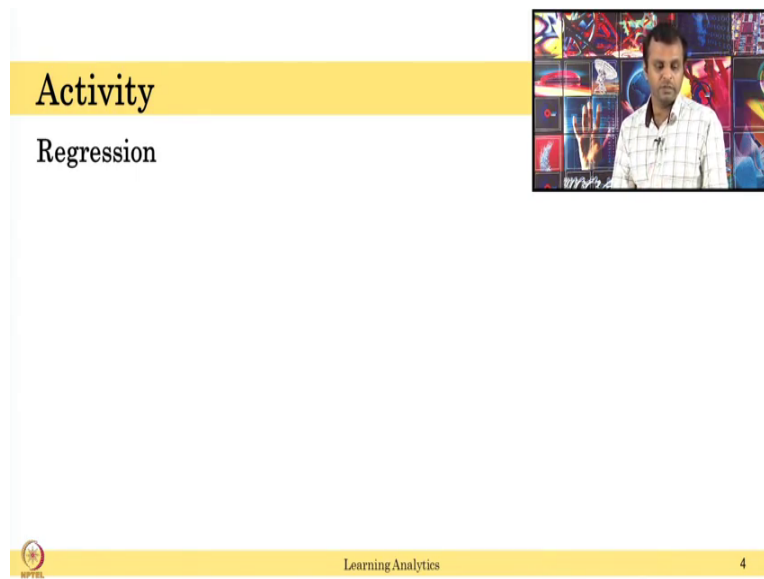
Learning Analytics 3

This paper discusses what we learned in general about our ROC curve, the area under the curve and also logistic regression introduction. Please go and read this paper understand how to collect data, how to collect features and how to use those features to predict the final score. So you saw, what is linear regression, also the logistic regression, can you list down one or two applications of linear regression using data collected from classroom environment?

The paper we just discussed now is applying linear regression on blended learning that is, the student is interacting with the MOOC kind of environment where they have to log in, watch videos, read things, post in discussion forums in a Moodle something like that.

Can you think of the application of using linear regression in a classroom environment and which data you can collect?

(Refer Slide Time: 9:17)



Activity

Regression

Learning Analytics 4

So list down your answers and resume to continue. So I am not writing any answers here because there are a lot of things that can be predicted I also discussed that debate at the starting of this week's course. You can predict students' performance, student's engagement. If you listed down something else, that is also good.

If you can access this data, you can go and collect the data, please go ahead and collect data and see which one works.

(Refer Slide Time: 10:02)

Summary

- Linear Regression - Application



In this week we discussed only linear and logistic regression and not in detail. This week is kind of less on new learning but I request you to go and explore the tools demo to you and use linear regression, logistic regression, collect data and use the existing data. Go and check Internet for the data, use the data and try to apply and understand. Thank you.