

EdX Dataset : Data Wrangling

Sukanya Chandramouli

Dataset

The HarvardX-MITx dataset ('HMXPC13_DI_v2_5-14-14.csv') contains de-identified data from the first year of 13 MITx and HarvardX courses offered on the edX platform. This large scale dataset contains more than six hundred thousand aggregated records.

<https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/26147&version=10.0>

Description of the Dataset

<i>Variable</i>	<i>Description</i>	<i>Data Source</i>
<i>course_id</i>	unique identifier with info about institution, course and semester	<i>administrative</i>
<i>userid_DI</i>	identifies user / course participant	<i>administrative</i>
<i>registered</i>	registered for course, (=1 for all records)	<i>administrative</i>
<i>viewed</i>	indicates the user has viewed the course content	<i>administrative</i>
<i>explored</i>	indicates user has accessed at least half of the chapters in the course	<i>administrative</i>
<i>certified</i>	indicates the user has completed the course with passing grade	<i>administrative</i>
<i>final_cc_cname_DI</i>	indicates the country / region of the user	<i>administrative</i>
<i>LoE</i>	level of education of the user	<i>user provided</i>
<i>YoB</i>	year of birth	<i>user provided</i>
<i>gender</i>	Possible values: m (male), f (female) and o (other)	<i>user provided</i>
<i>grade</i>	final grade of the user	<i>administrative</i>
<i>start_time_DI</i>	date of course registration	<i>administrative</i>
<i>last_event_DI</i>	date of last interaction with course	<i>administrative</i>
<i>nevents</i>	number of interactions with the course	<i>administrative</i>
<i>ndays_act</i>	number of unique days student interacted with course	<i>administrative</i>
<i>nplay_video</i>	number of play video events within the course	<i>administrative</i>
<i>nchapters</i>	number of chapters completed by the student	<i>administrative</i>
<i>inconsistent_flag</i>	identifies records that are internally inconsistent	<i>administrative</i>

Data Wrangling

Missing Values :

As a first step of data clean up, we look at the missing values in the data set. There are quite a few records for which the user provided data (Year of birth, education background, gender) is missing. We replace these values with the median value (for that particular course). Some of the administrative variables such as 'nchapters', 'ndays_act' etc which were missing are replaced with zero.

Data Types

The date of birth which is classified as a text is converted to a "date" format. Other numeric data types are set to appropriate data types (int64 , categorical , float). The grades which are specified between 0 and 1 is converted to percentage grades . The categorical order for educational level is set.

New Columns :

Using the Date of Birth information , we create a new column for the “age” of participants. The course-id is decoded to obtain the institution , course and semester/year. The course-id is mapped to the course title.

Renamed Columns:

Some of the column names are renamed for readability