

IBM Introduction to Machine Learning

Supervised Learning: Regression

Brief description and data exploration

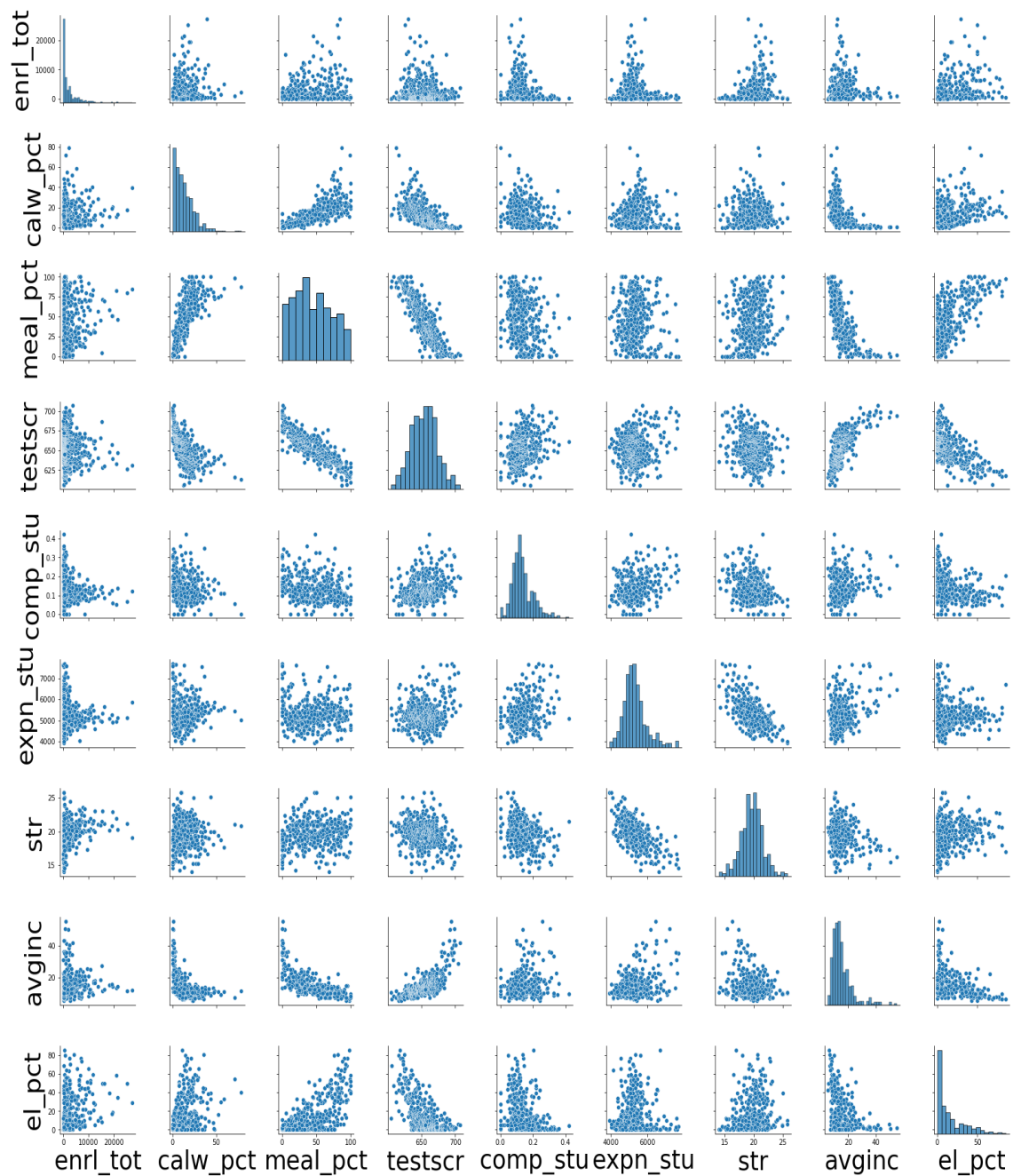
The California Standardized Testing and Reporting dataset contains data on test performance, school characteristics and student demographic backgrounds. The data used here are from 420 districts (45 counties) in California with data available for 1998 and 1999. Our main objective of the analysis will be focusing on interpretation. Python notebook code can be found on GitHub:

<https://github.com/estebanarboni/IBM-Introduction-to-Machine-Learning/blob/master/Supervised%20Learning%20Regression%20Final%20Project.ipynb>

	0	1	2	3	4
Observation Number	1	2	3	4	5
dist_cod	75119	61499	61549	61457	61523
county	Alameda	Butte	Butte	Butte	Butte
district	Sunol Glen Unified	Manzanita Elementary	Thermalito Union Elementary	Golden Feather Union Elementary	Palermo Union Elementary
gr_span	KK-08	KK-08	KK-08	KK-08	KK-08
enrl_tot	195	240	1550	243	1335
teachers	10.9	11.15	82.9	14	71.5
calw_pct	0.5102	15.4167	55.0323	36.4754	33.1086
meal_pct	2.0408	47.9167	76.3226	77.0492	78.427
computer	67	101	169	85	171
testscr	690.8	661.2	643.6	647.7	640.85
comp_stu	0.34359	0.420833	0.109032	0.349794	0.12809
expn_stu	6384.91	5099.38	5501.95	7101.83	5235.99
str	17.8899	21.5247	18.6972	17.3571	18.6713
avginc	22.69	9.824	8.978	8.978	9.08033
el_pct	0	4.58333	30	0	13.8577
read_scr	691.6	660.5	636.3	651.9	641.8
math_scr	690	661.9	650.9	643.5	639.9

School characteristics include enrollment, number of teachers, number of computers per classroom, and expenditures per student. Demographic variables include the percentage of students in the public assistance program, the percentage of students that qualify for a reduced-price lunch, and the percentage of students that are English Learners. We drop nonnumeric variables for data description and create a pairplot to analyze the relation between features:

	count	mean	std	min	25%	50%	75%	max
enrl_tot	420.0	2628.792857	3913.104985	81.000000	379.000000	950.500000	3008.000000	27176.000000
teachers	420.0	129.067376	187.912679	4.850000	19.662499	48.564999	146.350002	1429.000000
calw_pct	420.0	13.246042	11.454821	0.000000	4.395375	10.520450	18.981350	78.994202
meal_pct	420.0	44.705237	27.123381	0.000000	23.282200	41.750700	66.864725	100.000000
computer	420.0	303.383333	441.341298	0.000000	46.000000	117.500000	375.250000	3324.000000
testscr	420.0	654.156548	19.053348	605.550049	640.049988	654.449982	666.662506	706.750000
comp_stu	420.0	0.135927	0.064956	0.000000	0.093767	0.125464	0.164466	0.420833
expn_stu	420.0	5312.407541	633.937053	3926.069580	4906.180054	5214.516602	5601.401367	7711.506836
str	420.0	19.640425	1.891812	14.000000	18.582360	19.723208	20.871815	25.799999
avginc	420.0	15.316588	7.225890	5.335000	10.639000	13.727800	17.629001	55.327999
el_pct	420.0	15.768155	18.285927	0.000000	1.940807	8.777634	22.970003	85.539719
read_scr	420.0	654.970477	20.107980	604.500000	640.400024	655.750000	668.725006	704.000000
math_scr	420.0	653.342619	18.754202	605.400024	639.375015	652.449982	665.849991	709.500000

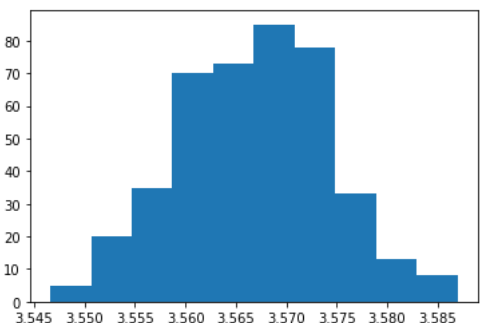
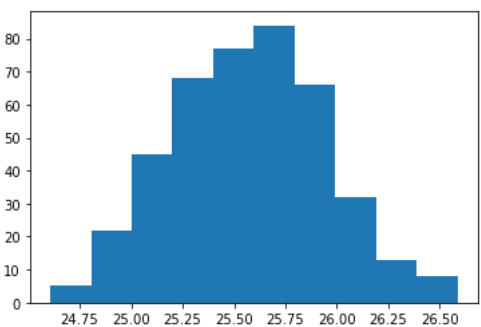
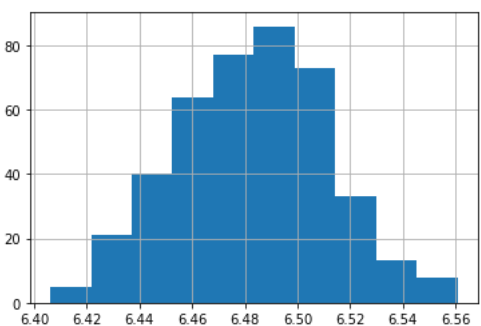
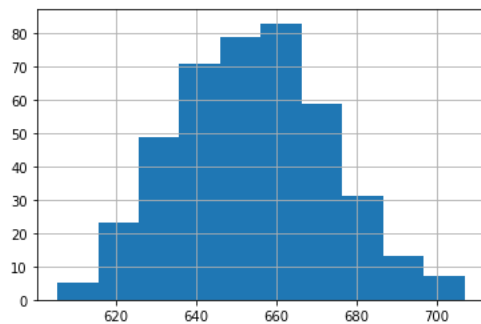


Objective of the analysis

The main objective of the analysis will be focusing on interpretation. We want to examine and determine the coefficients that do better at explaining the target variable (test scores). At a first glance we can guess that when student-teacher ratio (number of students per teacher) increases, test scores would be lower. Then we should focus on socioeconomic variables, where scores might be higher in counties with high average income. Variables like percentage of English learners (not native English speakers), computers per student or percentage of children qualifying for reduced-price lunch could be useful to determine the economic context. With this information, it would be possible to determine where to focus investment in a more efficient way in order to achieve better results

Data cleaning and feature engineering

We check if the target variable seems normally distributed and make log, square root and box cox transformations. We find that both log and box cox transformation give us a significantly more normally distributed (according to p value) than the original one.



Original: NormaltestResult(statistic=1.7477124875701222, pvalue=0.4173390809731915)

Log: NormaltestResult(statistic=1.404769712873947, pvalue=0.4954024301856593)

sqrt: NormaltestResult(statistic=1.4963705863719965, pvalue=0.47322453780319285)

Box cox: NormaltestResult(statistic=1.411882390136503, pvalue=0.4936437404868109)

Linear regression models

After normalizing our target variable, we train three linear regression models (simple linear regression, with polynomial effects and using Lasso regression), using the same train and test split. In terms of accuracy, Lasso regression outputs higher R^2 (0.80 against 0.77). Lasso did a better job explaining the additional variation. When we see the magnitude of the Lasso coefficients compared to linear regression, we are also able to reduce the complexity of the model, reduce the number and magnitude of the coefficients, and come up with a solution that's doing a better job of explaining variation on our holdout set, and will generalize better on new data.

Key findings and insights

Based on these results, we find that qualifying for reduced price lunch (MEAL_PCT) and being a native English speaker (EL_PCT) are the most negative affecting features. On the other hand, as income increases (AVGINC), higher scores are predicted.

	0	1
2 meal_pct	-10.175501	
7 el_pct	-3.619577	
1 calw_pct	-0.890676	
5 str	-0.359379	
0 enr_tot	0.001878	
3 comp_stu	0.771574	
4 expn_stu	0.966368	
6 avginc	4.486458	

Suggestions for next steps

By running a correlation matrix `".corr()"` we can analyze deeply the relationship between variables. Scores from reading or math could be analyzed separately and not as the mean of both (testscr) as the model deploys. Also, this dataset could be updated to date and measure the impact of better and more powerful computers in learning. Anyone can suggest or revisit the model to achieve a better explanation or a better prediction:

<https://github.com/estebanarboni/IBM-Introduction-to-Machine-Learning/blob/master/Supervised%20Learning%20Regression%20Final%20Project.ipynb>

	enr_tot	teachers	calw_pct	meal_pct	computer	testscr	comp_stu	expn_stu	str	avginc	el_pct	read_scr	math_scr
enr_tot	1.000000	0.802513	0.112117	0.129234	0.928882	-0.153988	-0.209229	-0.096537	0.298481	0.066636	0.381451	-0.188399	-0.110889
teachers	0.802513	1.000000	0.054434	0.081466	0.775572	-0.108695	-0.312705	-0.091189	0.269540	0.164226	0.418580	-0.150194	-0.059822
calw_pct	0.112117	0.054434	1.000000	0.815308	0.063496	-0.740436	-0.203053	-0.021580	0.099343	-0.725884	0.277650	-0.718566	-0.734058
meal_pct	0.129234	0.081466	0.815308	1.000000	0.061386	-0.868772	-0.207044	-0.049214	0.135203	-0.761504	0.576613	-0.878808	-0.823015
computer	0.928882	0.775572	0.063496	0.061386	1.000000	-0.073736	-0.034653	-0.056007	0.233826	0.125295	0.330575	-0.109005	-0.032950
testscr	-0.153988	-0.108695	-0.740436	-0.868772	-0.073736	1.000000	0.273762	0.177130	-0.226363	0.749586	-0.582960	0.981882	0.979143
comp_stu	-0.209229	-0.312705	-0.203053	-0.207044	-0.034653	0.273762	1.000000	0.268757	-0.310348	0.164493	-0.285790	0.283772	0.252000
expn_stu	-0.096537	-0.091189	-0.021580	-0.049214	-0.056007	0.177130	0.268757	1.000000	-0.625174	0.247891	-0.154721	0.203830	0.141368
str	0.298481	0.269540	0.099343	0.135203	0.233826	-0.226363	-0.310348	-0.625174	1.000000	-0.192305	0.257567	-0.246593	-0.195553
avginc	0.066636	0.164226	-0.725884	-0.761504	0.125295	0.749586	0.164493	0.247891	-0.192305	1.000000	-0.292921	0.743872	0.725517
el_pct	0.381451	0.418580	0.277650	0.576613	0.330575	-0.582960	-0.285790	-0.154721	0.257567	-0.292921	1.000000	-0.633464	-0.505327
read_scr	-0.188399	-0.150194	-0.718566	-0.878808	-0.109005	0.981882	0.283772	0.203830	-0.246593	0.743872	-0.633464	1.000000	0.922901
math_scr	-0.110889	-0.059822	-0.734058	-0.823015	-0.032950	0.979143	0.252000	0.141368	-0.195553	0.725517	-0.505327	0.922901	1.000000