

Exploring the role of alcohol and other factors in determining student achievement :D

Group 6 member:

Lei Zhu, Ruiyuan Ming, Minghao Mo, Yiming Xu, Xinrui Zhong

DATA: <https://www.kaggle.com/datasets/uciml/student-alcohol-consumption/data>

Variables:

school, sex, age, address,

famsize(family size)

Pstatus(parent's cohabitation status)

Medu(mother's education)

Fedu(father's education)

Mjob(mother's job)

Fjob(father's job)

guardian(student's guardian)

traveltime(home to school travel time)

studytime(weekly study time)

failures(number of past class failures)

schoolsup(extra educational support)

famsup (family educational support)

paid (extra paid classes within the course subject activities)

nursery - attended nursery school

higher - wants to take higher education

internet - Internet access at home

romantic - with a romantic relationship

famrel - quality of family relationships

freetime - free time after school

goout - going out with friends

Dalc - workday alcohol consumption

Walc - weekend alcohol consumption

health - current health status

absences - number of school absences

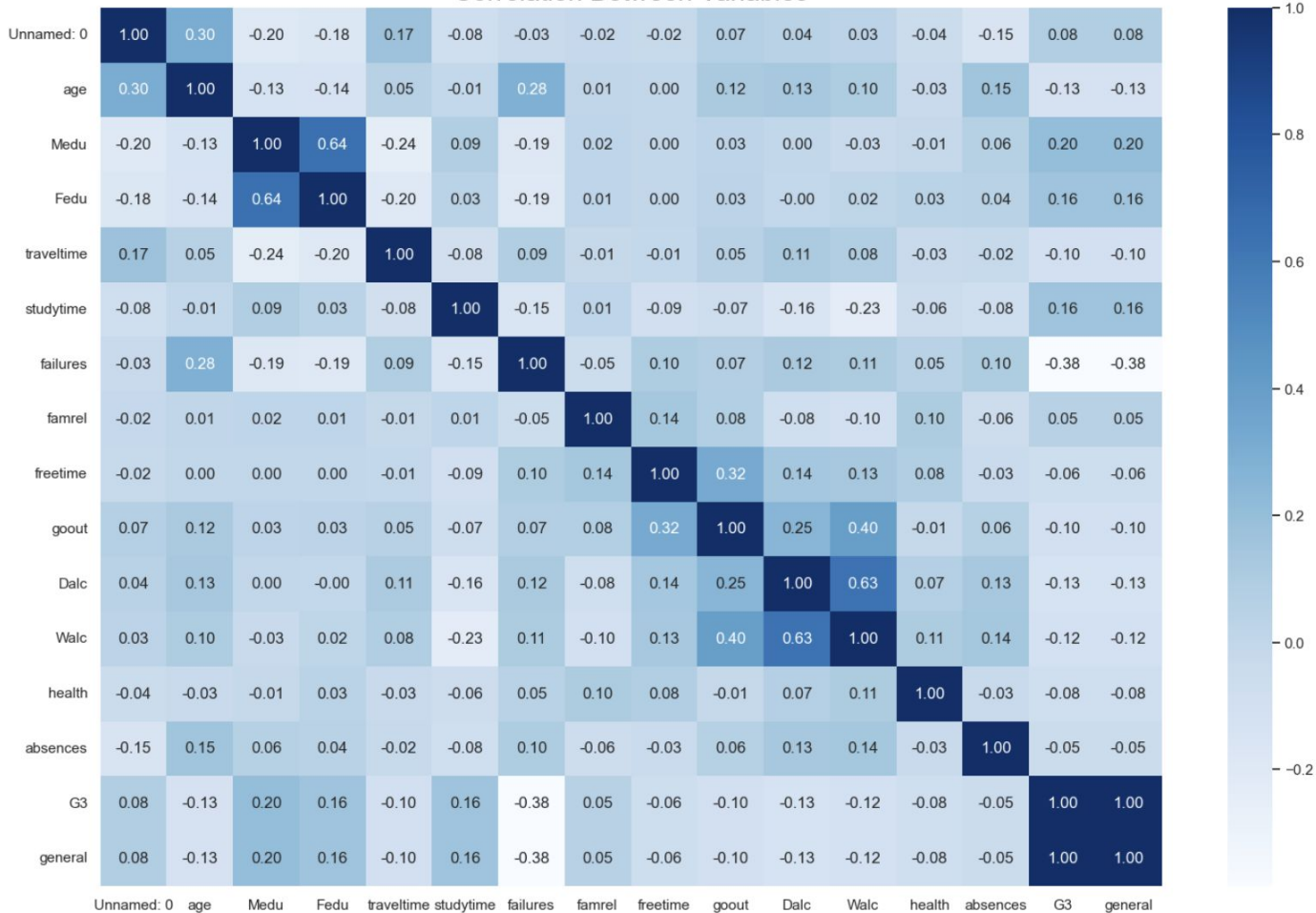
```
Data Shape: (1044, 32)
First Few Records:
  Unnamed: 0  school sex  age address  famsize Pstatus  Medu  Fedu  Mjob  \
0           0    GP  F  18      U    GT3      A    4    4  at_home
1           1    GP  F  17      U    GT3      T    1    1  at_home
2           2    GP  F  15      U    LE3      T    1    1  at_home
3           3    GP  F  15      U    GT3      T    4    2  health
4           4    GP  F  16      U    GT3      T    3    3  other

  ... internet romantic famrel  freetime  goout  Dalc  Walc  health  absences  \
0  ...      no         no      4          3      4      1      1      3          6
1  ...      yes        no      5          3      3      1      1      3          4
2  ...      yes        no      4          3      2      2      3      3         10
3  ...      yes        yes     3          2      2      1      1      5          2
4  ...      no         no      4          3      2      1      2      5          4

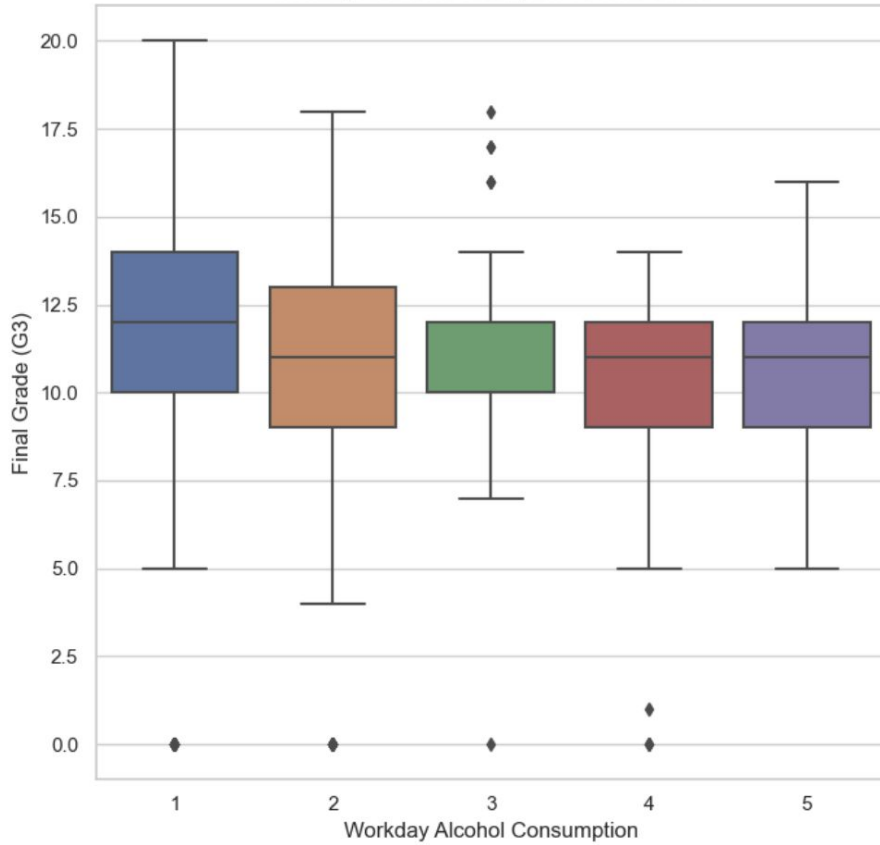
  G3
0    6
1    6
2   10
3   15
4   10

[5 rows x 32 columns]
```

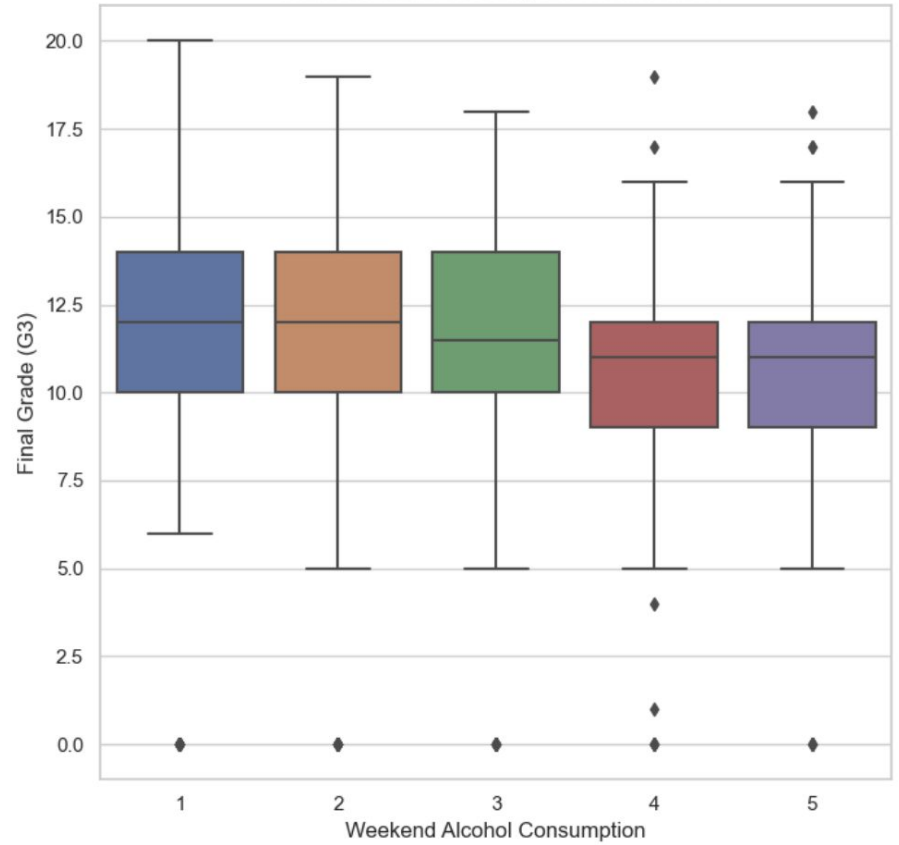
Correlation Between Variables



Workday Alcohol Consumption vs Final Grade



Weekend Alcohol Consumption vs Final Grade



Feature engineering

1. Check for missing values – no missing values
2. Split data: 80% for training and 20% for testing
3. One-hot encoding for qualitative variables
4. Min-max normalization for quantitative variables

Feature selection

Method: VarianceThreshold

Parameter: threshold=0.24 vs 0.25

Number of selected variables: 14 vs 7

Threshold=0.24

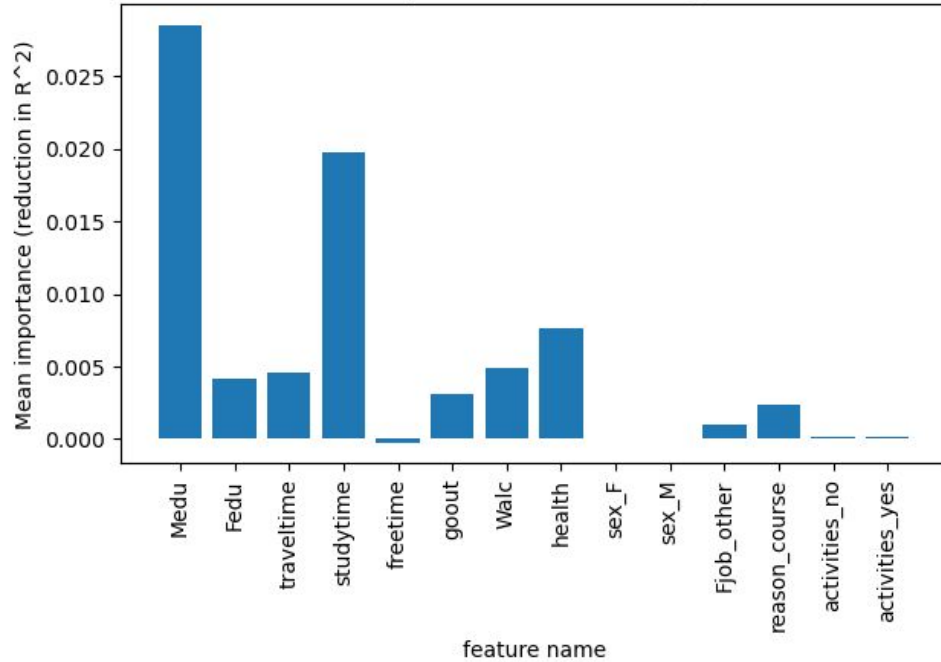
	Medu	Fedu	traveltime	studytime	freetime	goout	Walc	health
0	-0.5	-0.5	-1.000000	-0.333333	0.5	-0.5	-0.5	1.0
1	0.5	0.5	-1.000000	-0.333333	-0.5	0.0	-0.5	0.0
2	1.0	0.0	-1.000000	0.333333	0.0	0.0	0.0	-1.0
3	0.5	0.5	0.333333	-0.333333	0.0	0.0	0.0	-0.5
4	-0.5	-0.5	-1.000000	0.333333	0.0	0.5	-1.0	1.0
sex_F	sex_M	Fjob_other	reason_course	activities_no	activities_yes			
1.0	0.0	0.0	1.0	1.0	0.0			
0.0	1.0	1.0	0.0	0.0	1.0			
1.0	0.0	1.0	0.0	0.0	1.0			
0.0	1.0	1.0	0.0	1.0	0.0			
1.0	0.0	1.0	0.0	0.0	1.0			

Threshold=0.25

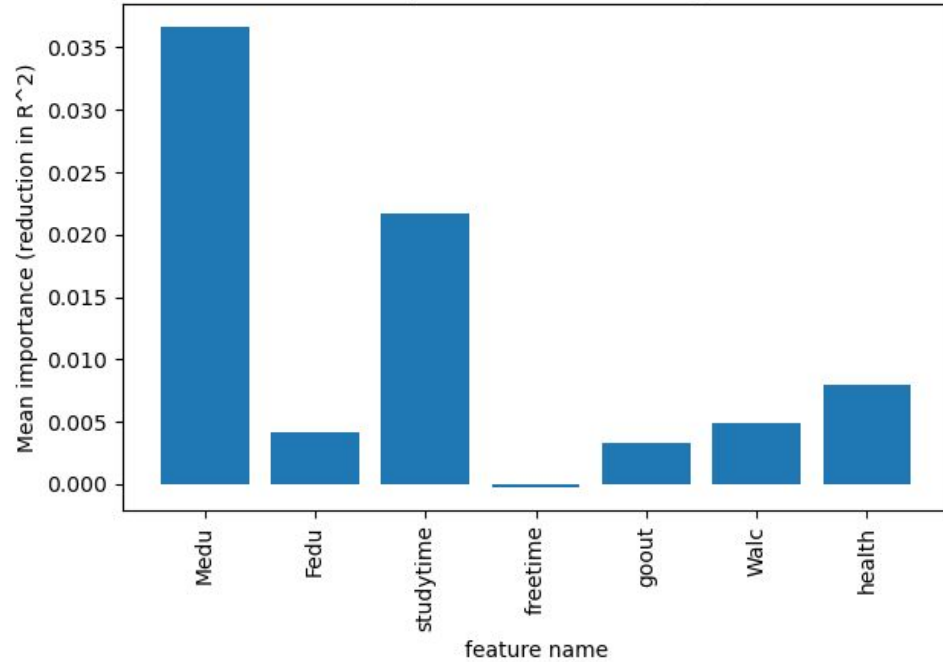
	Medu	Fedu	studytime	freetime	goout	Walc	health
0	-0.5	-0.5	-0.333333	0.5	-0.5	-0.5	1.0
1	0.5	0.5	-0.333333	-0.5	0.0	-0.5	0.0
2	1.0	0.0	0.333333	0.0	0.0	0.0	-1.0
3	0.5	0.5	-0.333333	0.0	0.0	0.0	-0.5
4	-0.5	-0.5	0.333333	0.0	0.5	-1.0	1.0

Linear Regression

Feature importance for Linear Regression



Feature importance for Linear Regression

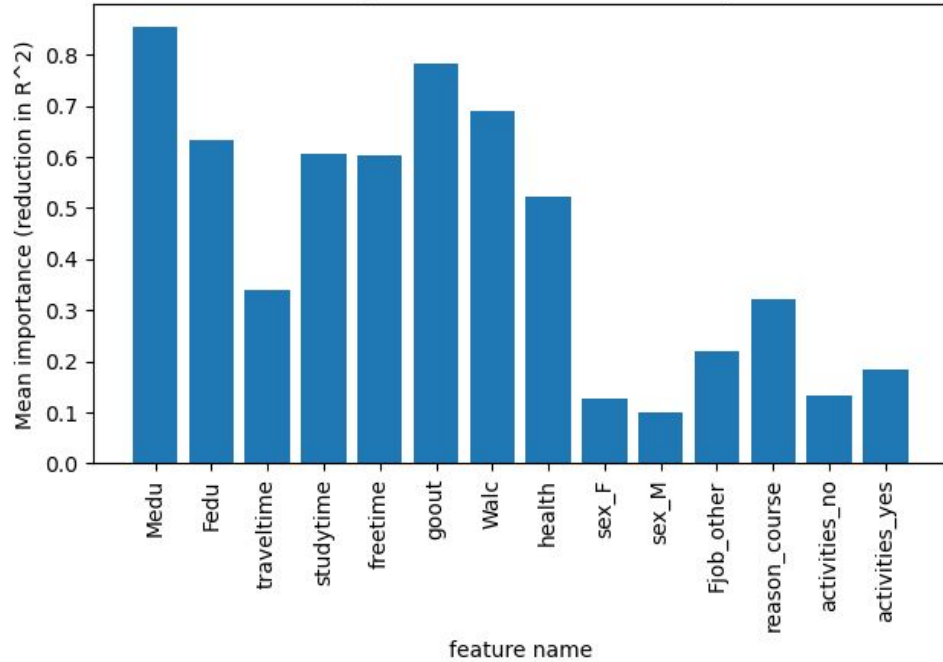


threshold = 0.24
MSE = 7.956

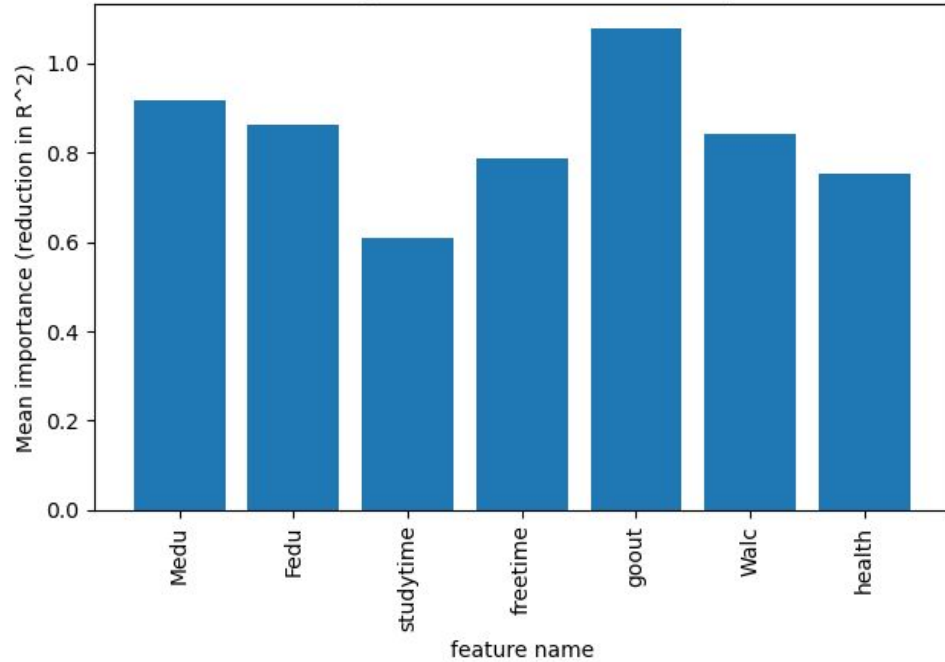
threshold = 0.25
MSE = 7.972

DecisionTree Regressor

Feature importance for DecisionTreeRegressor



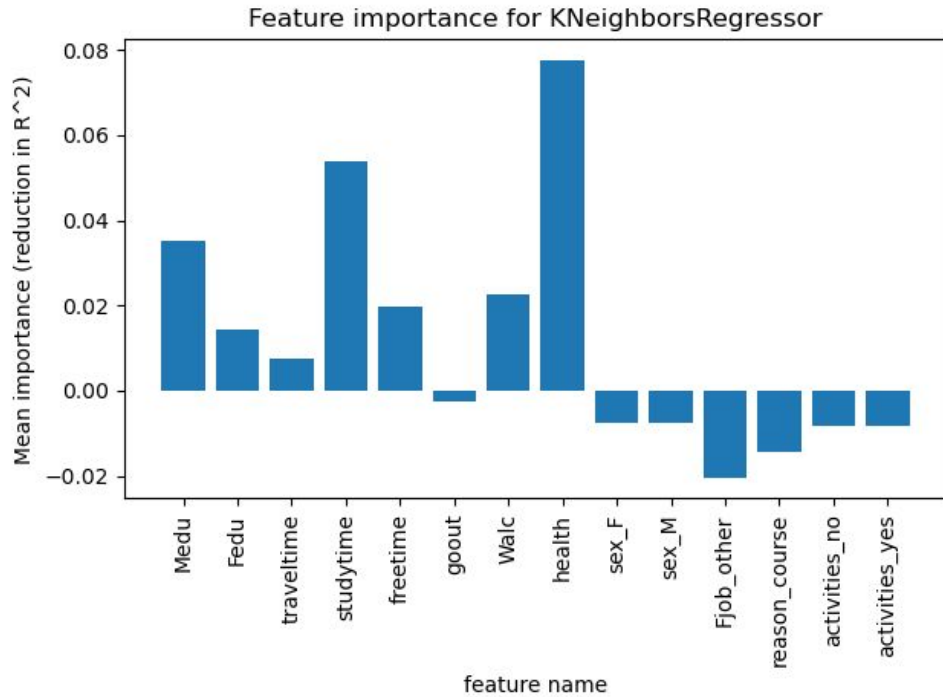
Feature importance for DecisionTreeRegressor



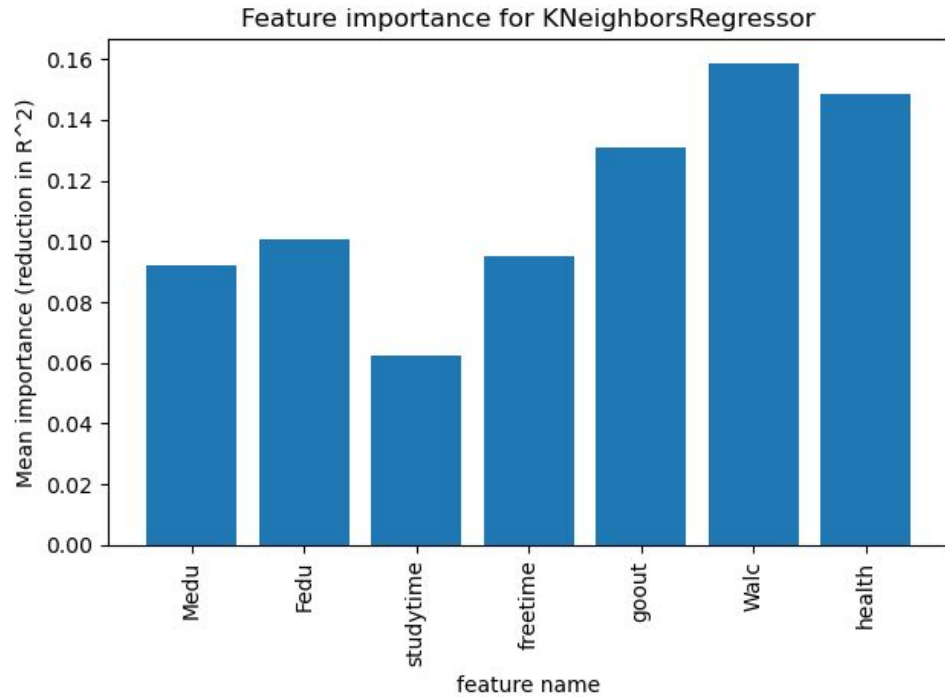
threshold = 0.24
MSE = 19.228

threshold = 0.25
MSE = 17.650

KNeighbors Regressor



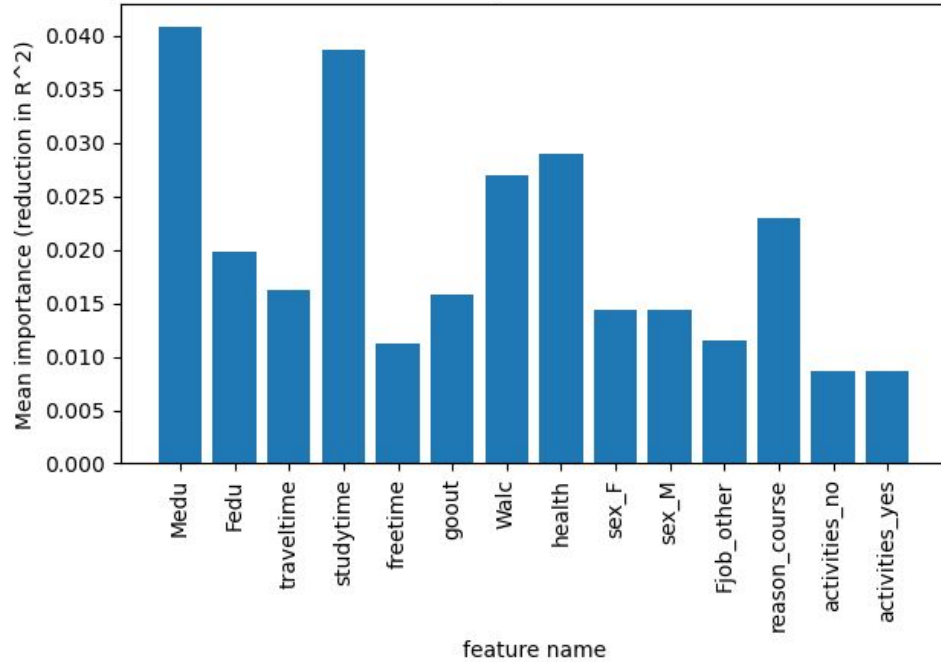
threshold = 0.24
MSE = 9.775



threshold = 0.25
MSE = 10.523

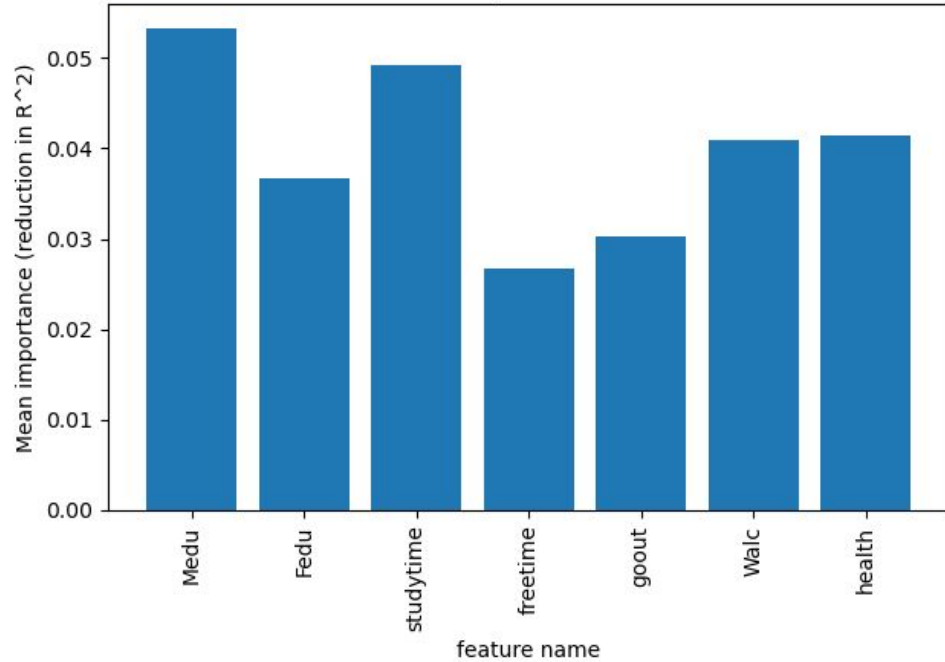
SVR

Feature importance for SVR



threshold = 0.24
MSE = 7.696

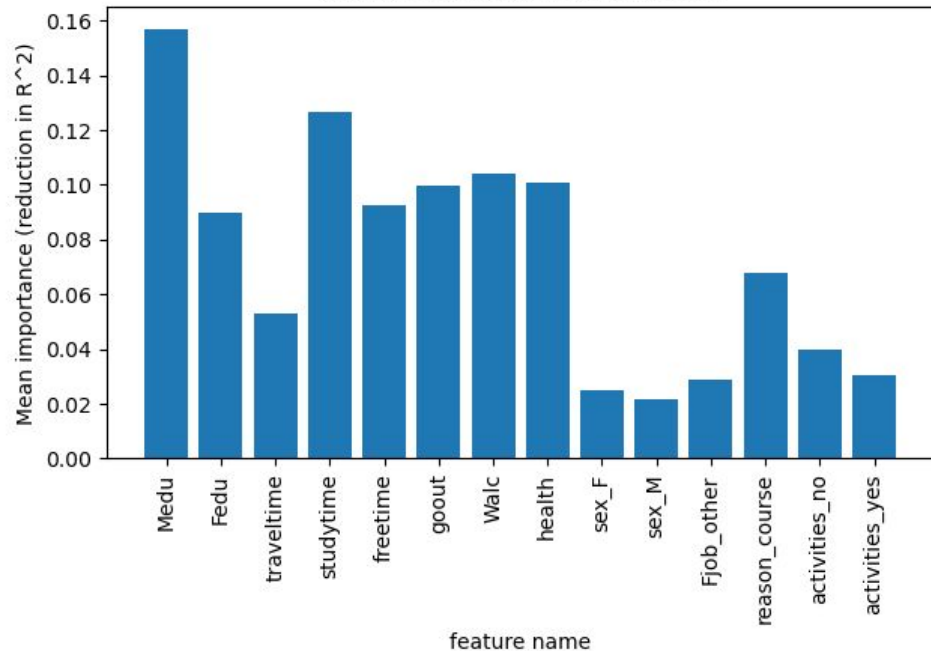
Feature importance for SVR



threshold = 0.25
MSE = 8.044

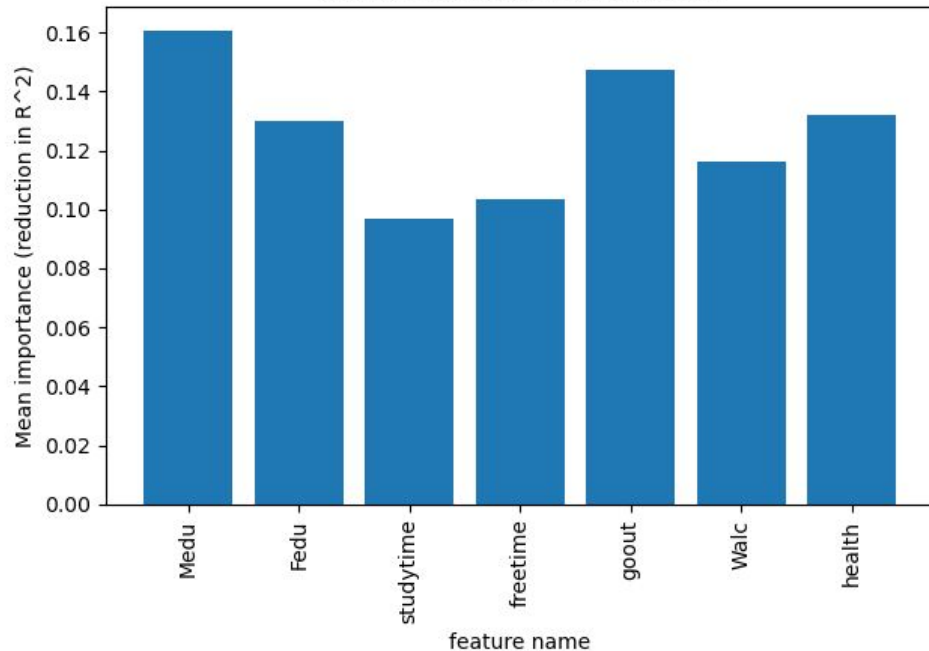
Stacking

Feature importance for stacking



threshold = 0.24
MSE = 7.572

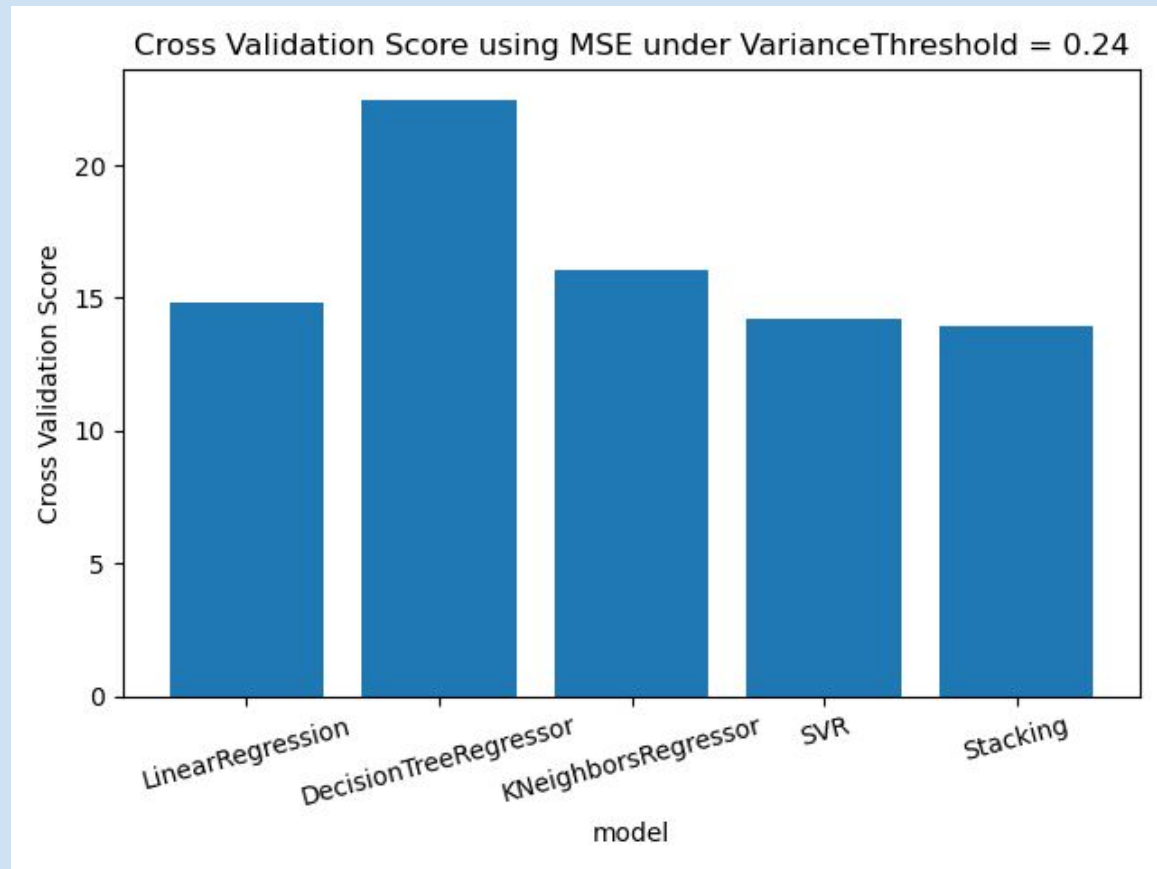
Feature importance for stacking



threshold = 0.25
MSE = 7.850

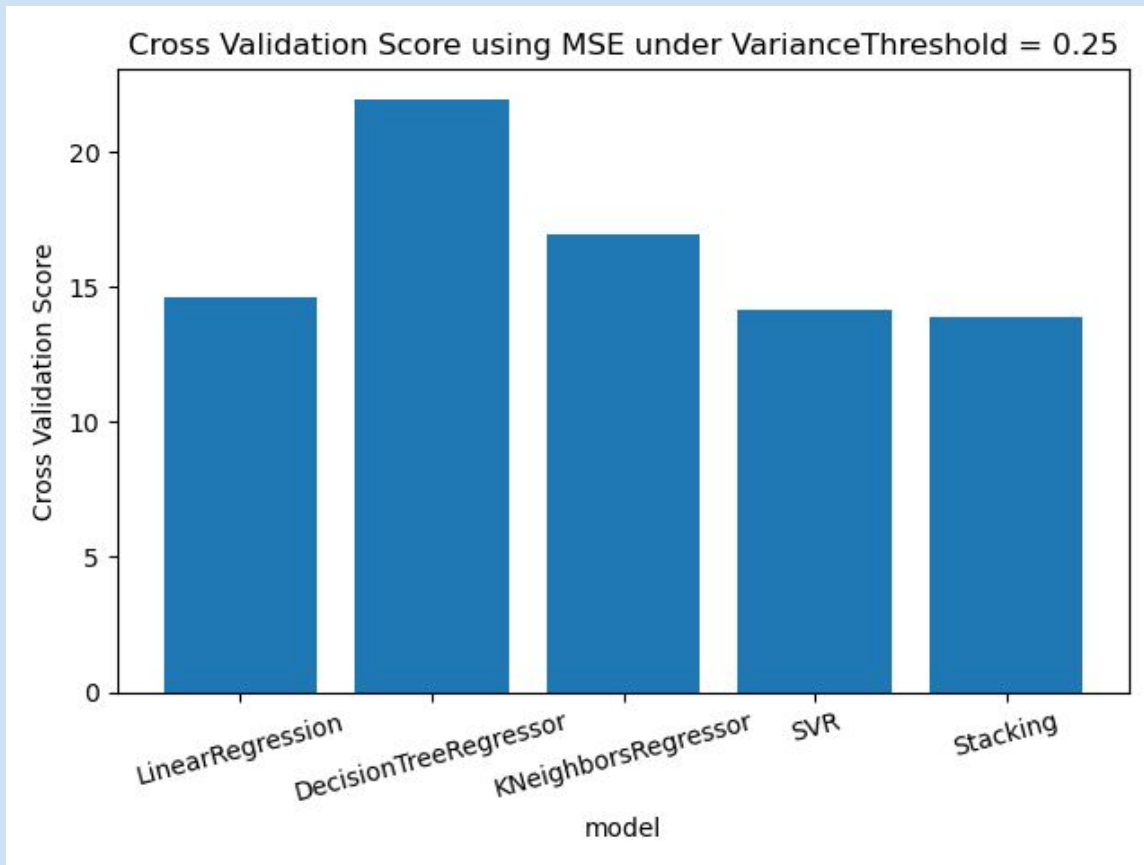
Cross Validation Score using MSE under VarianceThreshold = 0.24

CV score for LinearRegression: 14.8
CV score for DecisionTreeRegressor: 22.5
CV score for KNeighborsRegressor: 16.0
CV score for SVR: 14.2
CV score for Stacking: 13.9

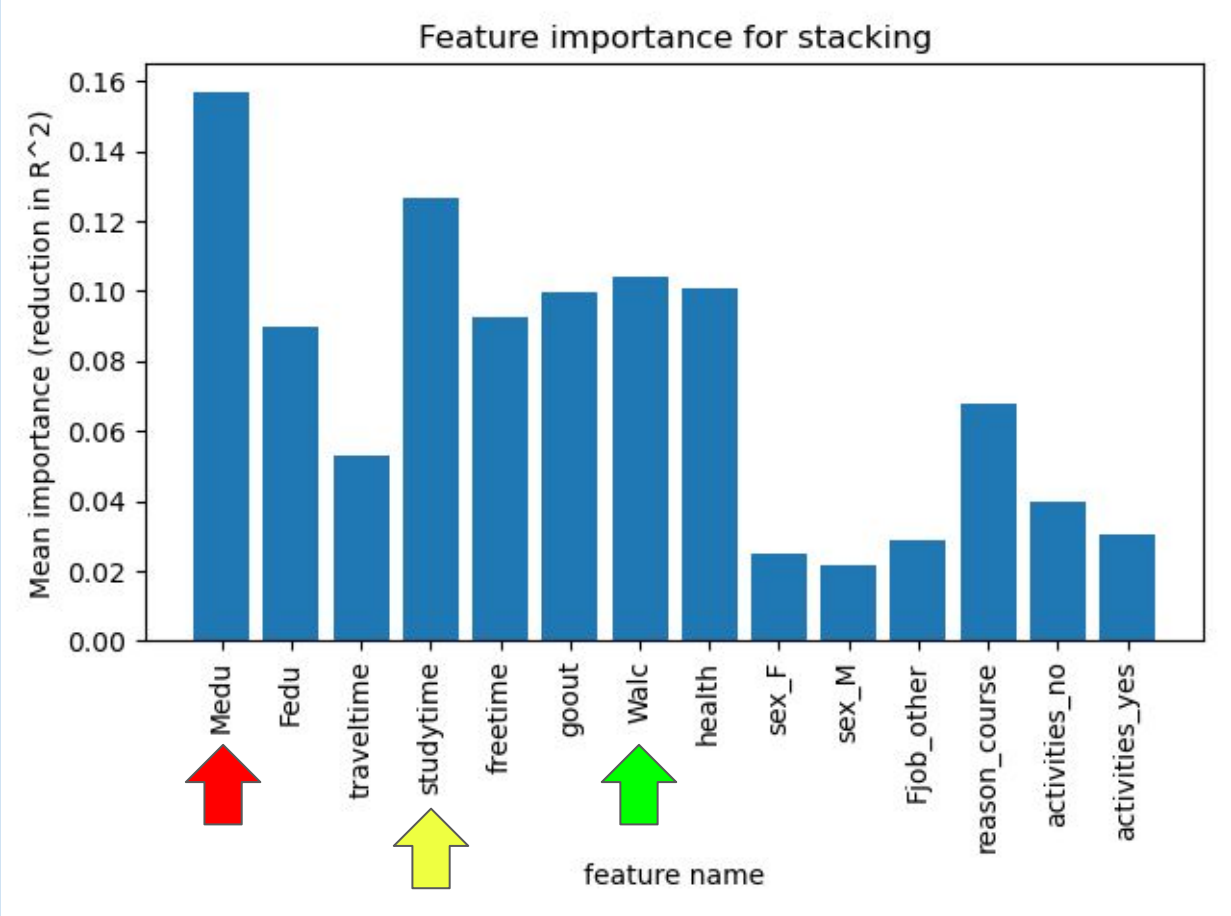


Cross Validation Score using MSE under VarianceThreshold = 0.25

CV score for LinearRegression: 14.6
CV score for DecisionTreeRegressor: 21.9
CV score for KNeighborsRegressor: 16.9
CV score for SVR: 14.2
CV score for Stacking: 13.9



Best model: Stacking under VarianceThreshold = 0.24



Thank you for listening!