

Statistical Models Homework 1

Exploratory Data Analysis

After setting up a virtual environment suitable for carrying out the scientific analysis, we load some useful packages to improve the capabilities of the analysis software. In particular, as suggested we load *ISLR*, *ISLR2*, *ROCR* and *tidyverse*. For the purposes of this research, we will also use *e1071*, *caret*, *corrplot*, *patchwork*, *class*

The first step in our analysis is to load the data contained in *chd.csv*. With the function *glimpse()* obtain an overview of the data.

```
Rows: 4,238
Columns: 13
$ sex      <chr> "Male", "Female", "Male", "Female", "Female", "Female", "Fem~
$ age      <dbl> 39, 46, 48, 61, 46, 43, 63, 45, 52, 43, 50, 43, 46, 41, 39, ~
$ education <dbl> 4, 2, 1, 3, 3, 2, 1, 2, 1, 1, 1, 2, 1, 3, 2, 2, 3, 2, 2, 2, ~
$ smoker   <dbl> 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, ~
$ cpd      <dbl> 0, 0, 20, 30, 23, 0, 0, 20, 0, 30, 0, 0, 15, 0, 9, 20, 10, 2~
$ stroke   <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
$ HTN      <dbl> 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, ~
$ diabetes <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
$ chol     <dbl> 195, 250, 245, 225, 285, 228, 205, 313, 260, 225, 254, 247, ~
$ DBP      <dbl> 70.0, 81.0, 80.0, 95.0, 84.0, 110.0, 71.0, 71.0, 89.0, 107.0~
$ BMI      <dbl> 26.97, 28.73, 25.34, 28.58, 23.10, 30.30, 33.11, 21.68, 26.3~
$ HR       <dbl> 80, 95, 75, 65, 85, 77, 60, 79, 76, 93, 75, 72, 98, 65, 85, ~
$ CHD      <chr> "No", "No", "No", "Yes", "No", "No", "Yes", "No", "No", "No"~
```

From the output we observe that the variables *sex* and *CHD* are represented as character variables, while the remaining predictors are stored as double-precision numerical values (real numbers). This suggests that most of the variables are already in a format suitable for analysis, although the character variables may have to be converted to factors depending on the modelling approach.

First of all, let us proceed with some basic checks on the data. The wise thing to do first is always to look for missing values.

sex	age	education	smoker	cpd	stroke	HTN	diabetes
0	0	105	0	29	0	0	0
chol	DBP	BMI	HR	CHD			
50	0	19	1	0			

Since we have a small data set to use, we think it is best to use an imputator to replace the missing values with the actual expected values. We will now make a copy of our data set, *data_simple*, on which we will apply the imputer function. We impute missing values with basic metrics such as mean, median (if the data distribution is skewed) for continuous data and mode for the categorical ordinal *education*. Mode imputation is generally reasonable when the distribution is highly imbalanced and the number of missing values is relatively small, making it acceptable to substitute with the most frequent category. However, mode imputation ignores the ordinal nature of the variable and does not take into account potential correlations with other variables. In our case, the number of missing values is small, so these limitations are unlikely to significantly affect the results. The output below confirms that the missing values were correctly replaced by the imputer function.

sex	age	education	smoker	cpd	stroke	HTN	diabetes
0	0	0	0	0	0	0	0
chol	DBP	BMI	HR	CHD			
0	0	0	0	0			

Now that we have handled the missing values, we will search for possible significant relationships between variables. For data exploration we'll plot only information from the *data_simple()* copy of the data. The first step will be to draw a heat map of the Pearson correlation matrix between the variables. This will help us identify potential multicollinearity issues between predictors and assess the strength of linear relationships between each predictor and the response variable. Since correlation requires numerical inputs, the categorical variables *CHD* and *SEX* must first be converted into binary numerical variables (e.g., 1 = Yes, 0 = No). Once transformed, they can be included in the correlation matrix. From Figure 1, it can be seen that there are no strong linear relationships between the predictors and the response variable. The correlation with CHD is weak for all characteristics, with *AGE* ($r = 0.23$), *HTN* ($r = 0.18$) and *DBP* ($r = 0.15$) showing the highest associations. These weak positive correlations suggest limited linear discriminative power, although they could still contribute useful information in a multivariate model. As expected, positive correlations were also found between *SMOKER* and *CPD* ($r = 0.77$), *HNT* and *DBP* ($r = 0.61$), *DBP* and *BMI* ($r = 0.38$). This indicates potential multicollinearity, which should be monitored during model fitting, particularly in the context of regression analysis, where correlated predictors may inflate the variance in coefficient estimates and reduce interpretability.

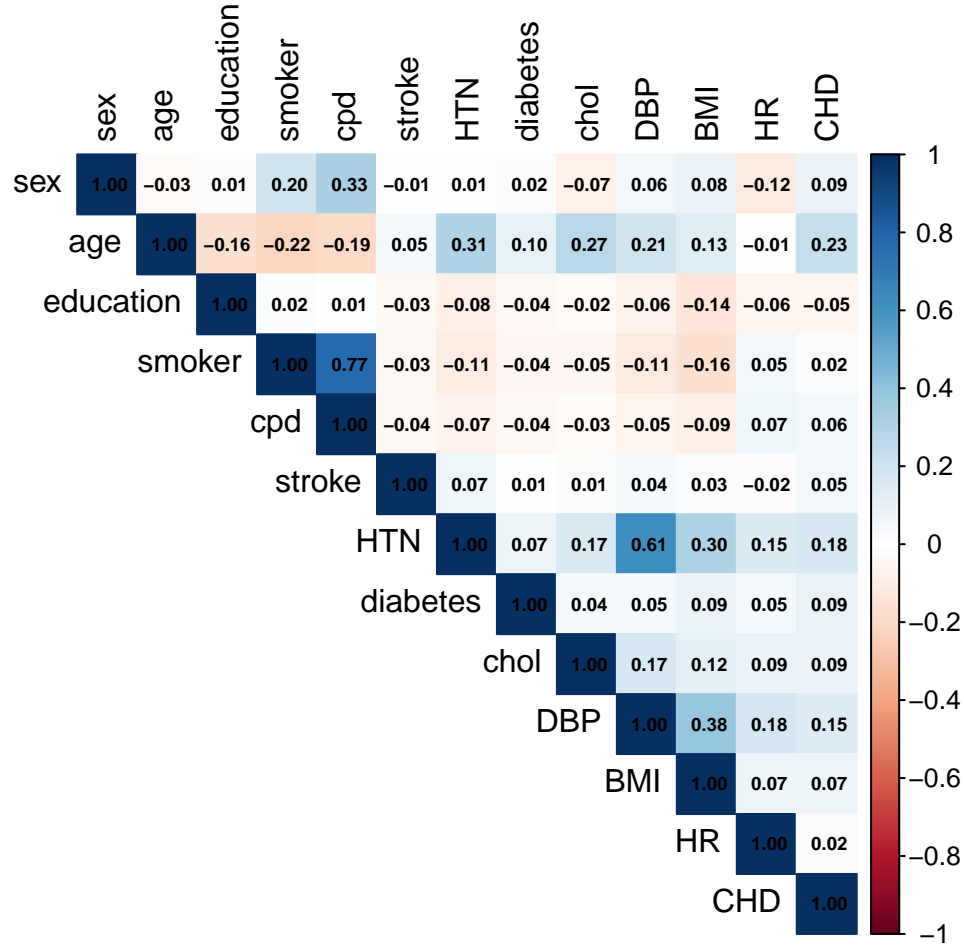


Figure 1: Correlation heatmap (including CHD)

To further investigate the discriminatory power of continuous and categorical predictors, we examined the distributions of those variables that showed some correlation with CHD. Specifically, we used boxplots for continuous predictors and bar plots with proportions for categorical ones. As shown in Figures 2 and 3, the distribution of CHD cases varies across different levels of the predictors *AGE*, *DBP*, *HTN* and *SEX*, providing additional visual evidence of a relationship between these variables and the response. These differences support the hypothesis that certain predictors may contribute useful information to CHD risk modeling.

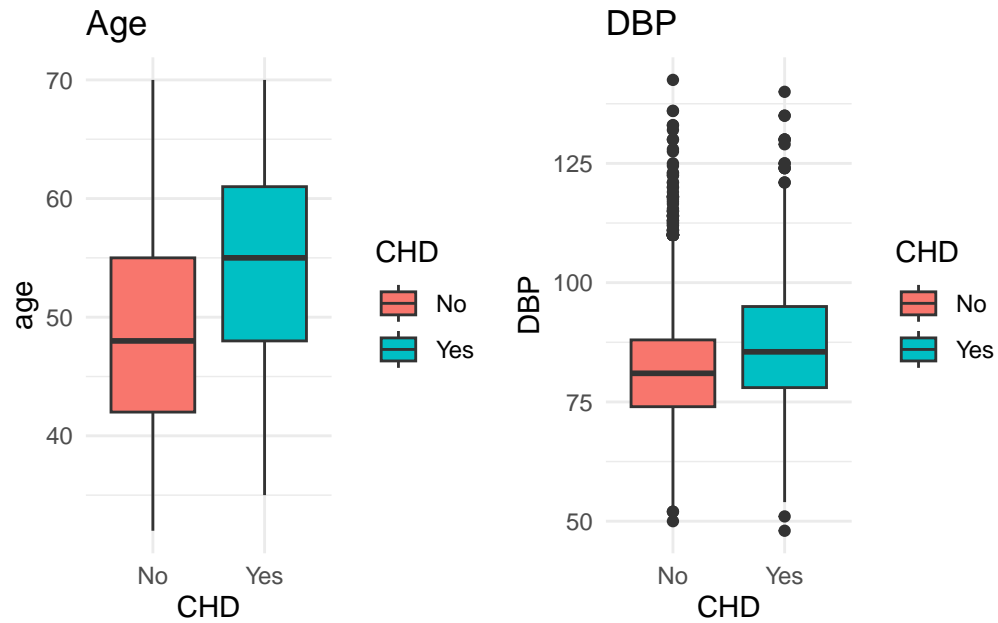


Figure 2: Distribution of continuous predictors against CHD

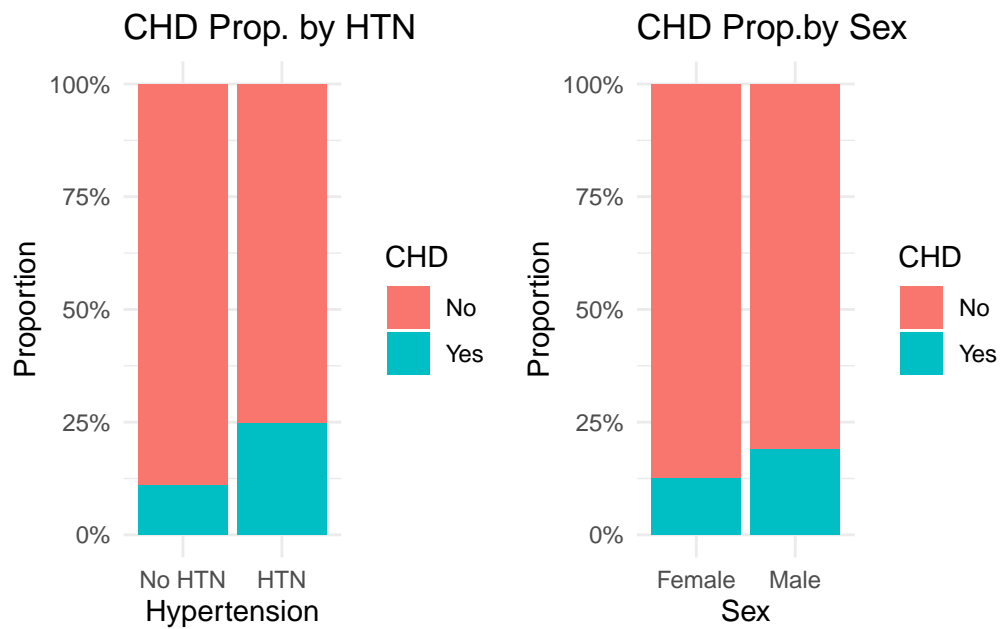


Figure 3: CHD proportions by hypertension status

We proceed with the analysis of the nature of the response variable. Since the response variable is categorical and binary we want to assess the level of balance inside the two categories by computing the proportion of each one of them on the total number of observations.

	No	Yes
	0.8480415	0.1519585

The output clearly shows that the response variable is strongly skewed toward the “No” class. This means that most observations do not have the event of interest (e.g., no coronary artery disease within 10 years). This has some implications for modeling, such as the fact that standard classifiers, e.g., logistic regression, will likely favor the most representative class, and perform poorly in making predictions about the minority class. In this scenario, the Accuracy is misleading, since the prediction of the majority class alone would still provide high accuracy (e.g., 84%) even if all cases of the minority class were missing.

In addition, we must prevent the random partitioning of data into a training and test set from leading to unbalanced sets where the true proportions of the sample are not actually represented correctly. To do this we use a technique called stratified sampling, which is aimed precisely at maintaining the proportions of the categories of a variable in both subsets of the sample. In practice, we will divide the sample into two groups filtered by the two categories of the response variable (in this case “yes” and “no”). These two groups are called *strata*. We then branch out some instances from the two groups and place them in our desired training or test set in a certain proportion (we chose to put 70% of the instances in the training set and the other 30% in the test). This maintains the original proportion in both sets. In R, we use `caret::createDataPartition()` which is the function provided to do this. After that, we verify the new proportions of both the sets.

	No	Yes
	0.8479946	0.1520054

	No	Yes
	0.8481511	0.1518489

Logistic Regression Model

We fit the following GLM model:

$$\text{logit}(E(\text{CHD})) = \beta_0 + \beta_1 \text{sex} + \beta_2 \text{age} + \beta_3 \text{education} + \dots + \beta_{12} \text{HR}$$

Call:

```
glm(formula = CHD ~ sex + age + education + smoker + cpd + stroke +  
      HTN + diabetes + chol + DBP + BMI + HR, family = binomial,  
      data = data_glm)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-7.6576275	0.7633802	-10.031	< 2e-16 ***
sexMale	0.3795087	0.1170652	3.242	0.001188 **
age	0.0770830	0.0070773	10.892	< 2e-16 ***
education	0.0277723	0.0539584	0.515	0.606764
smoker	0.0313558	0.1698117	0.185	0.853503
cpd	0.0201528	0.0067220	2.998	0.002717 **
stroke	0.4801136	0.5358126	0.896	0.370227
HTN	0.4321926	0.1419158	3.045	0.002324 **
diabetes	0.8800290	0.2484495	3.542	0.000397 ***
chol	0.0001026	0.0012284	0.084	0.933422
DBP	0.0054746	0.0054715	1.001	0.317033
BMI	0.0184361	0.0137308	1.343	0.179375
HR	0.0046910	0.0046094	1.018	0.308818

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 2528.9 on 2966 degrees of freedom
Residual deviance: 2277.6 on 2954 degrees of freedom
AIC: 2303.6

Number of Fisher Scoring iterations: 5

Interpretation of the Logistic Regression Model

The output summarizes the results of a fitted *logistic regression model* estimating the log-odds of developing *CHD* as a function of an intercept and 12 predictor variables. The estimated *Intercept* (-7.597) represents the log-odds of developing *CHD* for individuals in all reference categories (e.g., female, no high school degree, non-smoker, no history of stroke, hypertension, or diabetes). Its strongly negative value reflects a very low baseline probability of developing *CHD* in this subgroup. The coefficient for *sexMale* (0.440) is positive and highly significant ($p < 0.001$), suggesting that, all else equal, males have a higher risk of developing *CHD* than females. Similarly, *age* (0.0703) is highly significant, with each additional year increasing the

log-odds of *CHD*; the corresponding odds ratio is $\exp(0.0703) = 1.072$, indicating a 7.2% increase in odds per year of age. The variable *cpd* (cigarettes per day) is significant ($p < 0.001$), while *smoker* status is not. This implies that smoking intensity may be more predictive of *CHD* risk than smoking status alone. A history of *stroke* (1.023) is statistically significant ($p = 0.02$), with an odds ratio of approximately $\exp(1.023) = 2.78$. This indicates that individuals with prior strokes have nearly three times the odds of developing *CHD*. Similarly, *HTN* (hypertension) and *diabetes* are both highly significant, with odds ratios of about 1.55 and 2.36, respectively. The variable *chol* (cholesterol) has a small but borderline significant effect ($p = 0.046$). Although the per-unit increase has a minor effect on risk, extremely high cholesterol values could meaningfully affect *CHD* probability. *DBP* (diastolic blood pressure) is significant ($p = 0.006$), with an odds ratio of $\exp(0.0129) = 1.013$, meaning that each mmHg increase in diastolic pressure is associated with a 1.3% increase in *CHD* odds. By contrast, *education*, *smoker*, *BMI*, and *HR* are not statistically significant in this model, indicating that they may not provide additional predictive value for *CHD* in the presence of the other variables.

K-NN Classifier

Since we know that models based on clustering perform poorly with features with different scale, before the fitting process we standardize all the continuous variable. Without standardization, features like cholesterol or age might dominate the distance metric simply due to larger numeric ranges. We extract the mean and s.d. from the train set, and we use them to standardize both the train and test set. After this, we follow this procedure to fit a *K-NN* model: we set up a training control with a 5-fold Cross Validation, we apply grid search to fine tune the *k* parameter (the tuning grid for *k* will be $[5,30]$) and finally we fit the model after setting the seed to 42. We produce an output showing the best value of *k* and the accuracy result.

Highest accuracy of 0.8487 with *k* = 16 .

Performance evaluation

The next step is to evaluate the models. Given the nature of the response variable, it is clear that *accuracy* alone is not an appropriate evaluation metric. In particular, since this is a medical study, we are especially concerned with not missing high-risk patients, while we are more tolerant of issuing a false alarm. In this context, a more meaningful evaluation metric is the *FNR* (False Negative Rate), which measures the proportion of patients who developed *CHD* but were not identified by the system. The *FNR* is defined as:

$$FNR = \frac{FN}{FN + TP} = 1 - \text{Sensitivity}$$

We produce the confusion matrices of both models to extract this metric:

	Actual	
LR	No	Yes
No	1070	185
Yes	8	8

Logistic Regression Accuracy: 0.8482

Logistic Regression FNR: 0.9585

	Actual	
KNN	No	Yes
No	1075	191
Yes	3	2

k-NN Accuracy: 0.8474

k-NN FNR: 0.9896

Conclusion

In conclusion, the two models exhibit comparable performance in terms of overall *accuracy*. However, considering the metric of primary interest—the *False Negative Rate (FNR)*—the *logistic regression* model outperforms *k-NN*. Specifically, the *k-NN* model fails to correctly identify any true positive cases, misclassifying all actual *CHD* cases as negative. Although the *logistic regression* model also demonstrates a high false negative rate, incorrectly classifying approximately 96% of the true positives, it nonetheless provides a modest improvement over *k-NN* in identifying high-risk patients. Consequently, in an application where correctly identifying the minority class is crucial, *logistic regression* appears to be the more suitable model. The first limitation we wish to address concerns some missing values that were present in the dataset and were handled by simple imputation techniques. However, if the mechanism of missingness is not well captured by the distribution of the observed data, this may introduce bias into the analysis. Secondly, a major limitation is the imbalance in the outcome variable: approximately 85% of the observations correspond to the absence of *CHD*. This imbalance can lead the models to favor the majority class, thereby increasing the risk of overlooking high-risk individuals. Another limitation relates to potential multicollinearity among predictors, as correlations between key variables may distort coefficient estimates and reduce model interpretability. Additionally, some predictors, such as *diabetes*, occur infrequently in the dataset. This low prevalence can result in unstable coefficient estimates and limit the variable’s discriminative power. Finally, the dataset may omit relevant risk factors such as dietary habits, alcohol consumption, physical activity, or genetic predisposition. The absence of these variables could reduce the predictive accuracy of the models and overlook important contributors to *CHD* risk.