

Chronic Heart Failure and Risk Factors in Myocardial Infarction Dataset

Ariane, Alona, Minsu, and Jadey

Introduction

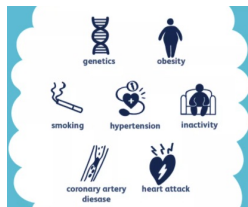
Chronic heart failure (CHF)

According to CDC,

- More than 6 million adults in the USA have heart failure.



- About half of Americans (47%) have at least one of key risk factors.



(Figure(up): <https://www.disability-benefits-help.org/resources/medical-evidence/chronic-heart-failure>)

(Figure(down):<https://www.verywellhealth.com/heart-failure-causes-and-risk-factors-1746181>)

Question: How are the predictors of our interest associated with Chronic heart failure?

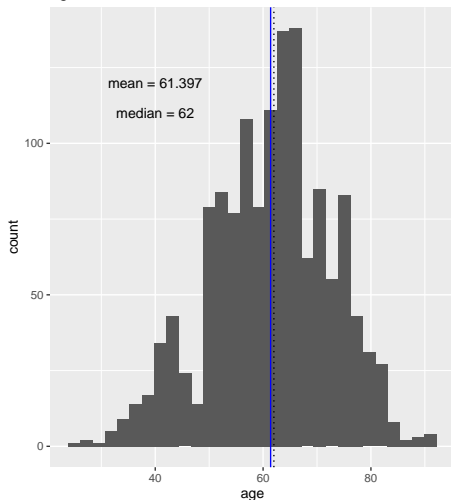
- Test independence of demographics with regards to CHF
- Association of duration of arterial hypertension and CHF
- Build a multiple logistic regression model by adding more predictors and identify the best model
- Modeling the relationship between death outcome and selected variables

Dataset Overview

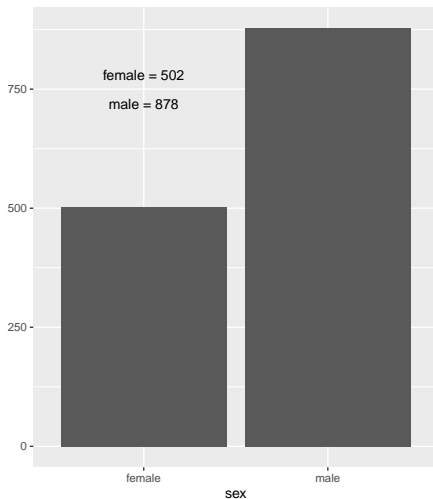
Descriptive statistics

- Demographic information

age distribution



distribution of sex

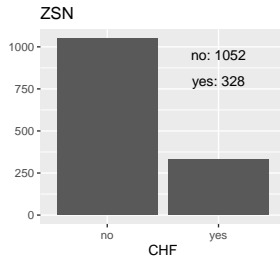
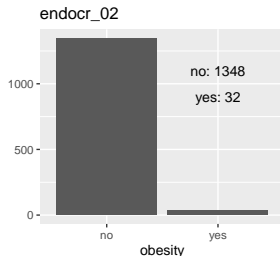
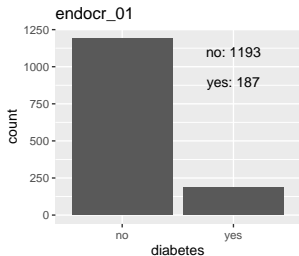
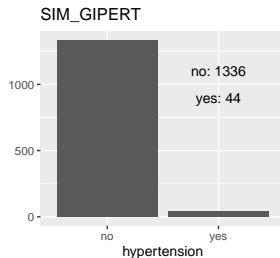
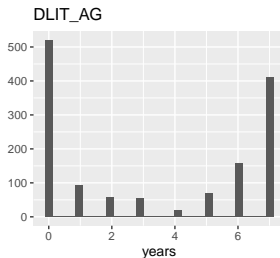
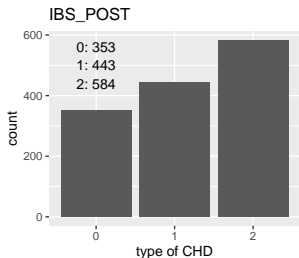


Descriptive statistics

- Patient physiological attributes
- IBS_POST: coronary heart disease in recent weeks before admission to hospital
- 0: there was no CHD
- 1: exertional angina pectoris
- 2: unstable angina pectoris
- DLIT_AG: duration of arterial hypertension
- 0: there was no arterial hypertension
- 1: one year
- 2: two years
- 3: three years
- 4: four years
- 5: five years
- 6: 6-10 years
- 7: more than 10 years

- SIM_GIPERT: systematic hypertension; 0 - no, 1 - yes
- endocr_01: diabetes mellitus in the anamnesis; 0 - no, 1 - yes
- endocr_02: obesity in the anamnesis; 0 - no, 1 - yes
- ZSN: chronic heart failure; 0 - no, 1 - yes

Descriptive statistics



Tests for Independence of Demographics

Analysis of Sex and Chronic Heart Failure: Overview

Question: Is there an association between sex and chronic heart failure?

Sex	Chronic Heart Failure	
	No	Yes
Female	353	149
Male	699	179

Analysis of Sex and Chronic Heart Failure: Tests

Pearson χ^2 Test of Independence:

X-squared
14.71773

p-value = 0.00012

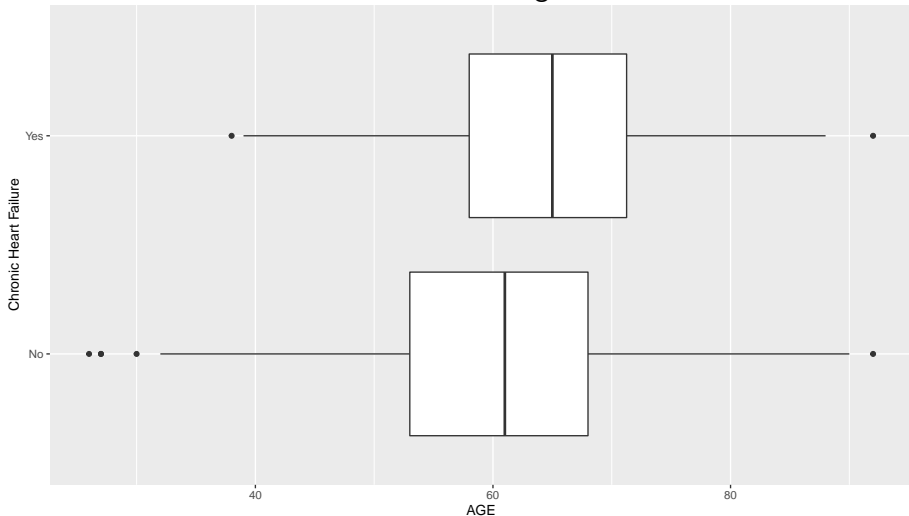
Likelihood Ratio Test of Independence:

G
14.93905

p-value = 0.00011

Analysis of Age(Continuous) and Chronic Heart Failure: Overview

Question: Is there an association between age and chronic heart failure?



Analysis of Age(Continuous) and Chronic Heart Failure: Summary Statistics

	Chronic Heart Failure	
	No	Yes
Min.	26	38
1st Qu.	53	58
Median	61	65
Mean	60.42586	64.51220
3rd Qu.	68.00	71.25
Max	92	92

Analysis of Age(Continuous) and Chronic Heart Failure: Test

Analysis was done using a two sided Wilcoxon Rank Sum Test to test if there is a difference in Chronic Heart Failure outcome across age.

W

136546.5

p-value = 1e-08

Analysis of Age(Categorical) and Chronic Heart Failure: Overview

Question: Is there an association between age(decade) and chronic heart failure?

Age	Chronic Heart Failure	
	No	Yes
20s	3	0
30s	44	2
40s	114	24
50s	294	67
60s	365	126
70s	197	86
80s	32	22
90s	3	1

Analysis of Age(Categorical) and Chronic Heart Failure: Test

Pearson χ^2 Test of Independence:

X-squared
35.41942

p-value = 1e-05

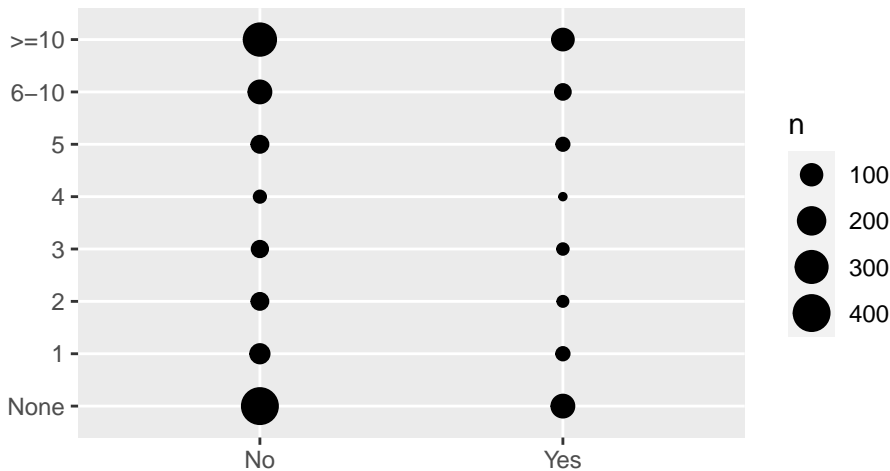
Likelihood Ratio Test of Independence:

G
38.86163

p-value = 2.08e-06

Association of duration of arterial hypertension and CHF

Examining the relationship between Duration of Arterial Hypertension and CHF



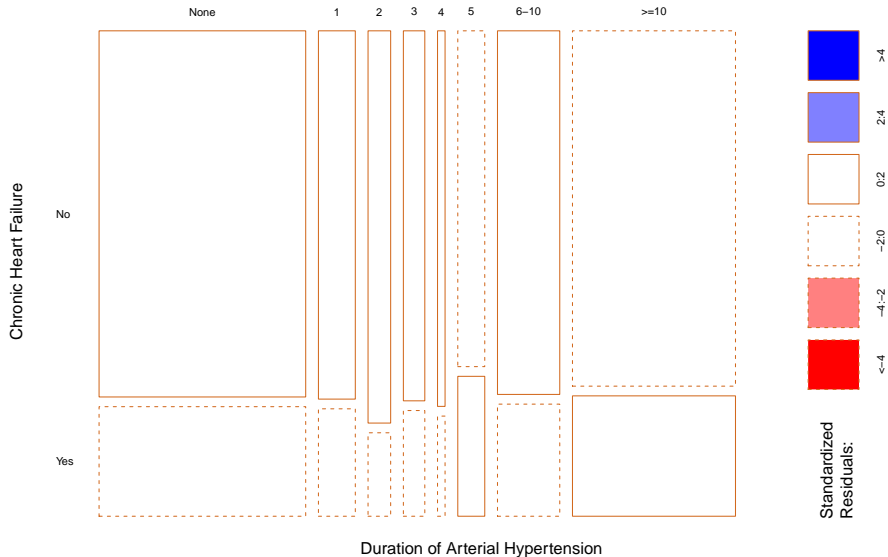
- The two classes of CHF have similar count distributions across the levels of duration of arterial hypertension.
- We will further test the hypothesis that there is an association between the two variables

Inference for contingency table.

Table 1: Duration of Arterial Hypertension by Chronic Heart Failure

	No	Yes
None	401	120
1	72	21
2	47	10
3	42	12
4	15	4
5	48	20
6-10	120	37
≥ 10	307	104

Examining the Standerdized residuals.



For $I \times 2$ tables, testing for a linear trend in either response category, we use the Cochran-Armitage trend test.

Cochran-Armitage test for trend

```
data:  dlitag  
Z = -0.99455, dim = 8, p-value = 0.32  
alternative hypothesis: two.sided
```

Issues to consider: Ordinal variable with unequal intervals so trend test on the original classification provides information about the direction but ignores the unequal spacing in the last two categories.

Logistic Regression model

x - Duration of Arterial Hypertension.

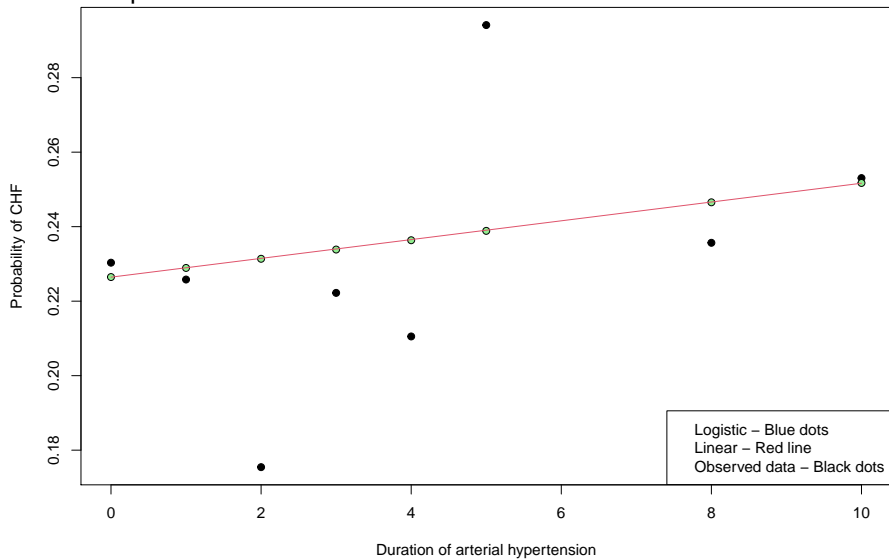
Table 2: Parameter Estimates for Logit link

	Estimate	Std. Error	z value	$\Pr(> z)$
(Intercept)	-1.2283412	0.0915051	-13.4237468	0.0000000
x	0.0138949	0.0143812	0.9661872	0.3339505

Table 3: Parameter Estimates for Identity link

	Estimate	Std. Error	z value	$\Pr(> z)$
(Intercept)	0.2264438	0.0160982	14.0664047	0.0000000
x	0.0025212	0.0026207	0.9620338	0.3360326

Predicted probabilities for the fitted models and the observed data.



We tested the Linear probability model for the subset: Duration of arterial hypertension between 1 and 5.

Table 4: Parameter Estimates for subset analysis

	Estimate	Std. Error	z value	$\Pr(> z)$
(Intercept)	0.1850895	0.0483670	3.826774	0.0001298
DLIT_AG_N	0.0167632	0.0161478	1.038107	0.2992204

Multiple Logistic Regression and Model Selection

Multiple Logistic Regression

- Coefficient estimates of the multiple logistic regressions of all predictors

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.02031	0.46737	-6.46239	0.00000
AGE	0.03221	0.00671	4.79990	0.00000
SEX	-0.17273	0.15006	-1.15112	0.24968
IBS_POST	-0.03283	0.08284	-0.39632	0.69187
DLIT_AG_N	-0.02809	0.01632	-1.72118	0.08522
SIM.fyes	-0.40006	0.40869	-0.97888	0.32764
endocr_01.fyes	0.75213	0.17773	4.23174	0.00002
endocr_02.fyes	0.15009	0.41093	0.36525	0.71493

- Only AGE and endocr_01 are statistically significant.
- The P-value for the overall test is much less than 0.0001, thus there is strong evidence that at least one predictor has an effect.

Multiple Logistic Regression - Goodness of Fit

Fit a multiple logistic regression model by adding AGE and endocr_01 to the logistic regression model with only DLIT_AG:

$$\text{logit}[P(ZSN = 1)] = \alpha + \beta_1 DLIT_AG + \beta_2 AGE + \beta_3 endocr_01.f.$$

- Goodness of Fit

G.square	df	P-value
1458.899	1376	0.0591161

The model has $G^2 = 1459$ with degree of freedom $df = 1376$ (P-value=0.059 > 0.05), which indicates a decent fit.

Multiple Logistic Regression - ANOVA test

Comparing this additive model with the initial model with DLIT_AG only,

- ANOVA Result

Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1378	1512.63	NA	NA	NA
1376	1458.90	2	53.73	0

the likelihood ratios test statistic is 53.73 with degree of freedom 2, producing very tiny p-value ($P < 0.001$). Thus, the model with AGE and endocr_01 in addition to DLIT_AG improves the goodness-of-fit.

Multiple Logistic Regression - Model selection

We perform stepwise model selection to see if there is effect of interaction between predictors.

- Backward selection

Step	Df	Deviance	Resid. Df	Resid. Dev	AIC
	NA	NA	1372	1450.891	1466.891
- DLIT_AG_N:AGE:endocr_01.f	1	0.225	1373	1451.116	1465.116
- AGE:endocr_01.f	1	0.594	1374	1451.710	1463.710
- DLIT_AG_N:AGE	1	0.421	1375	1452.131	1462.131

Multiple Logistic Regression - Model selection

- Forward selection

Step	Df	Deviance	Resid. Df	Resid. Dev	AIC
	NA	NA	1379	1513.563	1515.563
+ AGE	-1	34.509	1378	1479.054	1483.054
+ endocr_01.f	-1	17.399	1377	1461.655	1467.655
+ DLIT_AG_N	-1	2.755	1376	1458.899	1466.899
+ DLIT_AG_N:endocr_01.f	-1	6.769	1375	1452.131	1462.131

Based on the AIC, both backward elimination and forward selection choose the model of

$$\text{logit}[P(ZSN = 1)] = \alpha + \beta_1 DLIT_AG + \beta_2 AGE + \beta_3 endocr_01.f + \beta_4 DLIT_AG * endocr_01.f.$$

Predictive Power - ROC curves

- ROC curves of the selected model with interaction and the additive model

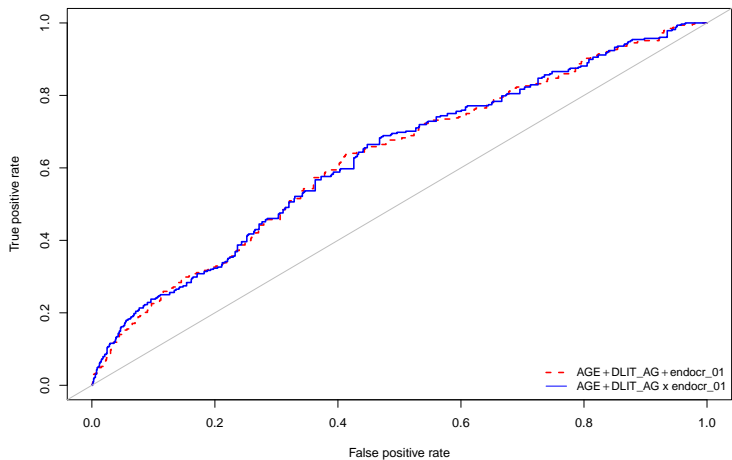


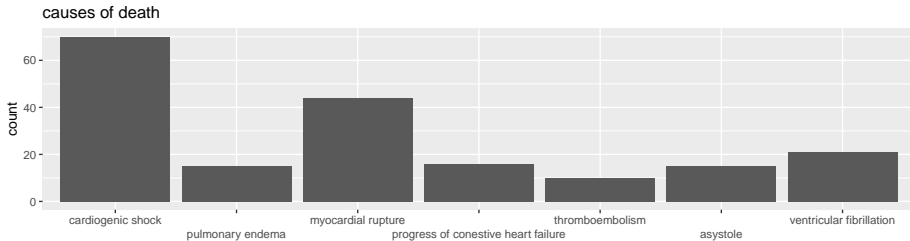
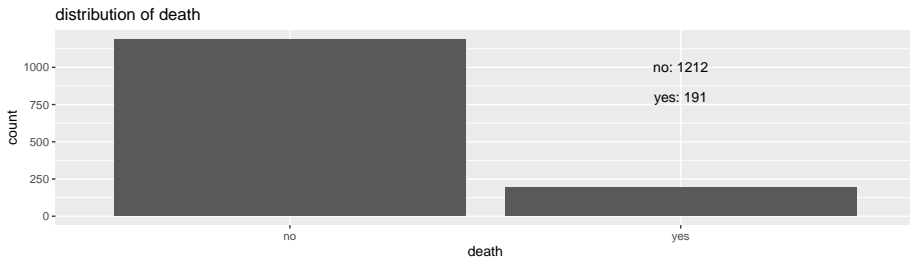
Figure 1: ROC curves

Modeling the relationship between death outcome and selected variables

Secondary analysis

- The dataset includes one variable indicating the causes of lethal outcome for the patients
- LET_IS: causes of lethal outcome
 - 0: survive
 - 1: cardiogenic edema
 - 2: pulmonary edema
 - 3: myocardial rupture
 - 4: progress of congestive heart failure
 - 5: thromboembolism
 - 6: asystole
 - 7: ventricular fibrillation
- Build a logistic regression model to predict death of the patients by turning LET_IS to a binary variable “death”
- Build model with multinomial response to investigate the cause of death

Secondary analysis



Secondary analysis

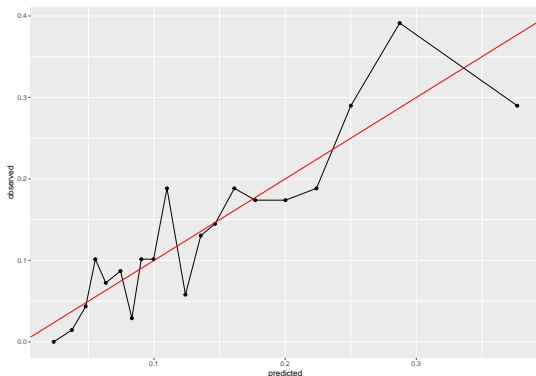
- Full model contains continuous variables age, duration of arterial hypertension, categorical variables SEX, chronic heart disease duration before admission to hospital, systematic hypertension, diabetes, obesity, and the interaction terms between AGE and all the other variables.
- Used stepwise step() to select the best model.
- The best model selected:

$$\log[P(\text{death} = 1)] = \beta_0 + \beta_1 \times \text{age} + \beta_2 \times I(\text{IBS} = 1) + \beta_3 \times I(\text{IBS} = 2) + \beta_4 \times I(\text{SIM} = 1) + \beta_5 \times I(\text{endocr01} = 1) + \beta_6 \times I(\text{endocr02} = 1)$$

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-6.018	0.572	-10.527	0.000
AGE	0.058	0.008	6.943	0.000
factor(IBS_POST)1	0.073	0.249	0.294	0.769
factor(IBS_POST)2	0.696	0.227	3.064	0.002
factor(SIM_GIPERT)1	0.726	0.393	1.846	0.065
factor(endocr_01)1	0.476	0.203	2.345	0.019
factor(endocr_02)1	1.081	0.403	2.679	0.007

Secondary analysis

- Goodness of fit check with Hosmer-Lemeshow test by grouping the observations into 20 groups. The test statistic is 0.4291, indicating an adequate fit of the model to the dataset.
- Plotted the predicted value against the observed value of the 20 groups. Overall the dots follow the diagonal.



Secondary analysis

- Fit baseline category logit model on cause of death. Used predictors selected in the previous analysis.

$$\log \frac{\pi_j(x)}{\pi_J(x)} = \beta_{0j} + \beta_{1j} \times \text{age} + \beta_{2j} \times I(IBM = 1) + \beta_{3j} \times I(IBM = 2) + \beta_{4j} \times I(SIM = 1) + \beta_{5j} \times I(\text{endocr01} = 1) + \beta_{6j} \times I(\text{endocr02} = 1), j = 1, \dots, 6$$

where J = cardiogenic shock, $j = 1$ pulmonary edema, 2 myocardial rupture, 3 progress of congestive heart failure, 4 thromboembolism, 5 asystole, 6 ventricular fibrillation

Secondary analysis

	intercept	AGE	IBS_POST = 1	IBS_POST = 2	SIM_GIPERT = 1	endocr_01 = 1	endocr_02 = 1
2	-5.208	0.050	0.417	-0.261	-14.880	1.348	-14.029
3	-2.662	0.045	-0.873	-1.325	-0.013	0.237	0.668
4	-3.190	0.030	-0.397	-0.722	0.004	1.173	-15.228
5	1.047	-0.038	-0.226	-1.507	-16.165	1.573	-16.311
6	-2.552	0.031	-2.039	-1.308	-16.884	0.824	0.765
7	2.873	-0.057	-0.533	-0.288	0.257	-0.726	-15.708

Conclusion

- Age and Sex are associated with CHF
- Duration of Arterial Hypertension is predictive when included in a multivariate model
- The final multivariable model for CHF is not rejected
- Age, coronary heart disease in recent weeks, symptomatic hypertension, obesity and diabetes are associated with patient death.