

Capstone Project II

Predicting individual annual cost of medical insurance by ML method

By
Lok Hang Ronald, Wong

Background Information

- ❑ Medical cost is one of the largest burdens of Americans
- ❑ Health care spending per person surpassed \$10,000 in 2016
- ❑ March steadily higher to \$14,944 in 2023.



Statistic source from Centers for Medicare and Medicaid Services (CMS)

Who might Care?

- ❑ Start-up Companies
 - Provide services such as consulting
 - Finding potential customers
- ❑ General Public
 - Evaluation of their current plans
 - Better planning and strategy

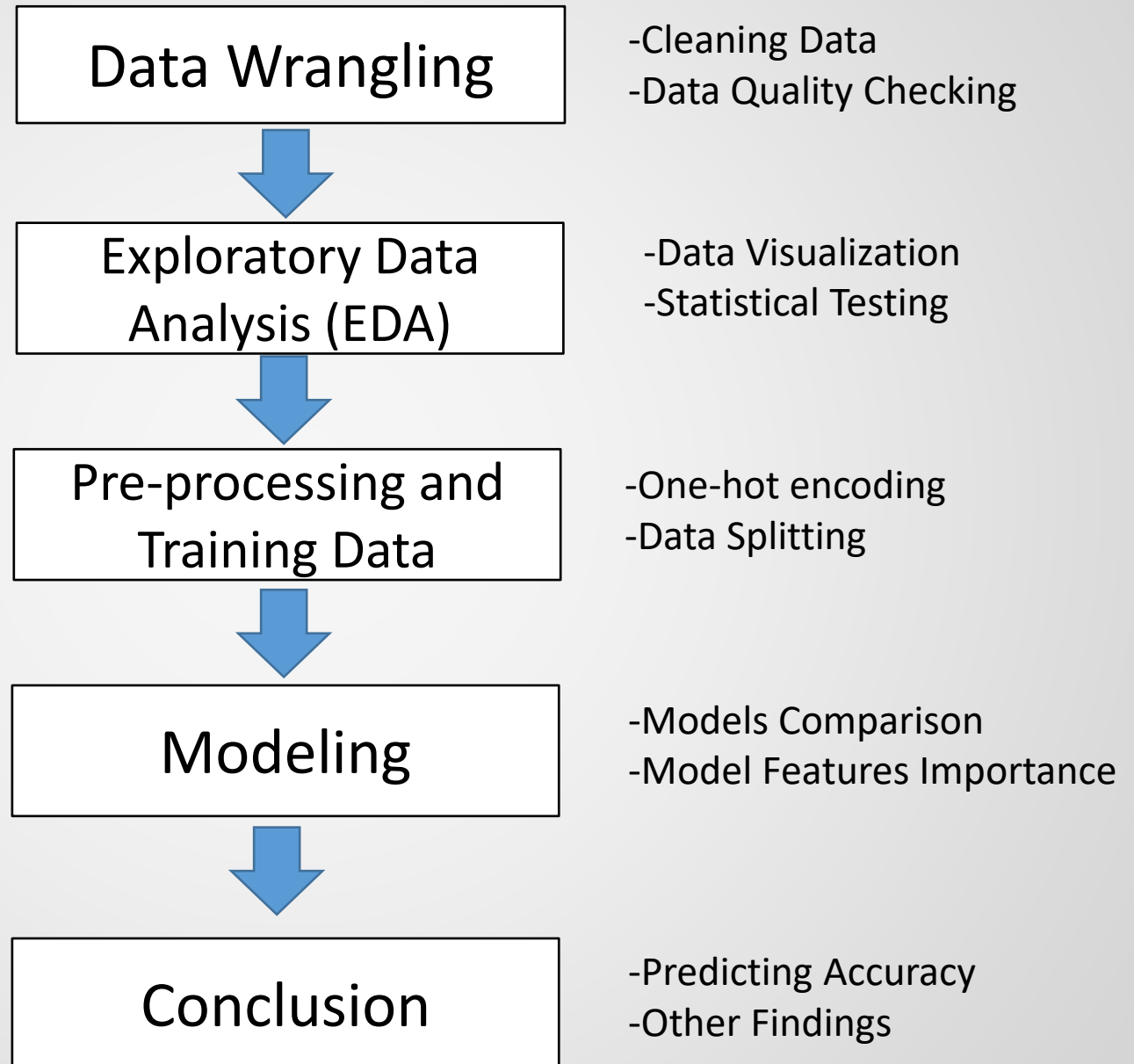


Dataset

1337 customers with their information obtained from Kaggle

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520

Project Outline



Exploratory Data Analysis (EDA) – t-test

- ❑ For comparing categorical variable with numerical variable
- ❑ For 2 unique values in categorical variable
- ❑ Null Hypothesis: Mean of two groups are the same
- ❑ Adopted for gender and smoking behavior

Exploratory Data Analysis (EDA) – One-way ANOVA

- ❑ For comparing categorical variable with numerical variable
- ❑ For 3 or more unique values in categorical variable
- ❑ Null Hypothesis: Mean of all groups are the same
- ❑ Adopted for region and family size (0 to 5 children)

Exploratory Data Analysis (EDA) – chi-squared test

- ❑ For checking correlation between categorical variables
- ❑ Similar with heatmap for numerical variables
- ❑ Null Hypothesis: The categorical variables are not correlated
- ❑ Findings: Smoker and Gender, bmi and Gender are correlated

Modelling

Approach 1	Approach 2
Ordinary least squares (OLS)	Random Forest Regressor (RFR)

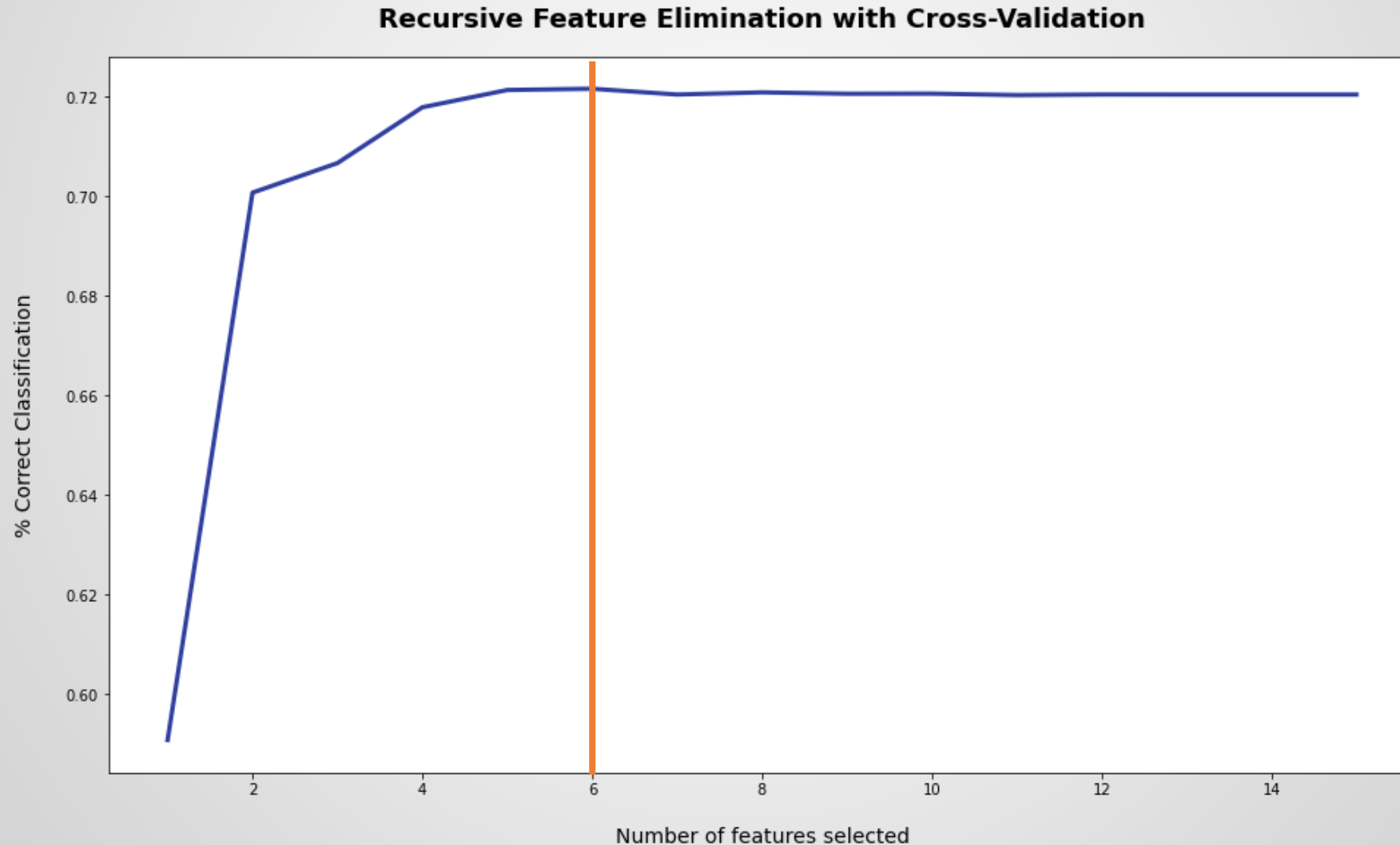
Modelling – Feature Selection

OLS Summary Report

	coef	std err	t	P> t	[0.025	0.975]
const	6462.7392	305.143	21.179	0.000	5863.884	7061.594
age	3686.3005	202.075	18.242	0.000	3289.720	4082.881
bmi	1772.7240	212.844	8.329	0.000	1355.008	2190.440
sex_1	70.4981	402.658	0.175	0.861	-719.735	860.731
smoker_1	2.268e+04	508.044	44.648	0.000	2.17e+04	2.37e+04
children_0_1	42.0528	416.949	0.101	0.920	-776.227	860.333
children_1_1	263.7529	474.098	0.556	0.578	-666.684	1194.190
children_2_1	1563.5188	523.927	2.984	0.003	535.290	2591.747
children_3_1	873.4268	600.356	1.455	0.146	-304.796	2051.650
children_4_1	2677.4248	1172.073	2.284	0.023	377.184	4977.666
children_5_1	1042.5631	1334.109	0.781	0.435	1575.679	3660.805
SW_1	1426.1022	358.018	3.983	0.000	723.478	2128.726
SE_1	1242.4228	364.633	3.407	0.001	526.816	1958.030
NW_1	2113.1498	364.828	5.792	0.000	1397.160	2829.139
NE_1	1681.0644	364.648	4.610	0.000	965.427	2396.702

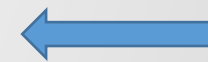
Modelling – Feature Selection

Recursive Feature Elimination and Cross-Validation Selection (RFECV)



Result – OLS Model

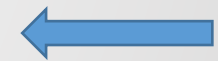
Model	R2 Score	RMSE
All Features	0.786	6033
Best 6 Features	0.643	7793
Drop 5 Features	0.785	6042



Best OLS Model

Result – RFR Model

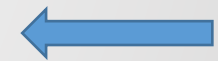
Model	R2 Score	RMSE
All Features	0.862	4840
Best 6 Features	0.862	4844



Best RFR Model

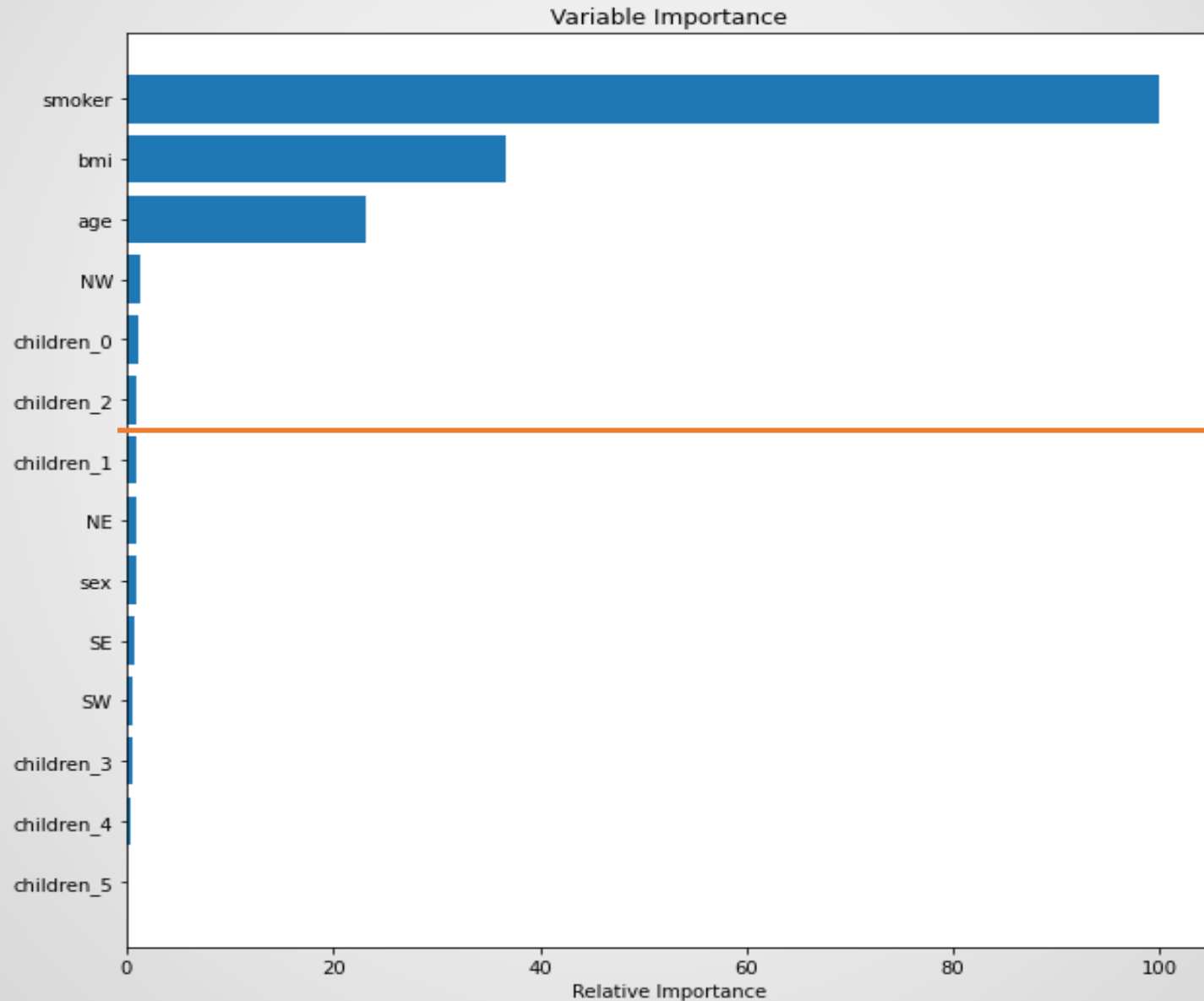
Result –Model Selection

Model	R2 Score	RMSE
Best OLS	0.785	6042
Best RFR	0.862	4844



Selected Model

Result –Feature Importance



**Adopted in
selected
model**

**Not adopted
in selected
model**

Future Studies

- ❑ Increasing data size to prevent over-fitting
- ❑ Comprise more customer information (More features)
- ❑ Comparing more ML methods such as K-NN and Neural Network

Conclusion

- ❑ Established a model for predicting medical insurance cost
- ❑ Explored important and non-essential factors
- ❑ Importance to start-up companies and general public

Q & A