



Unveiling the Predictive Power of Big Data: Predicting Smoking and Drinking Behavior

IST 718 - Group 7

Vedant Patil

Shrish Vaidya

Sahil Wani

Chintan Patel

Introduction

- Investigate health metrics: Deep insights into lifestyle patterns.
- Comprehensive dataset: Essential for understanding health dynamics.
- Explore impact: Correlate lifestyle choices with overall well-being.
- Predictive resource: Informing outcomes in alcohol and tobacco research.
- Support evidence-based decisions: Vital for health professionals and policymakers.



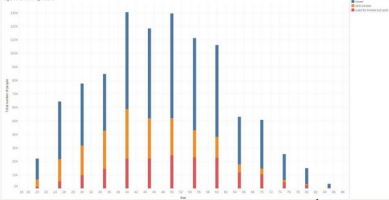


Predictive Modeling Goals

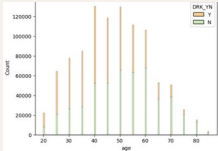
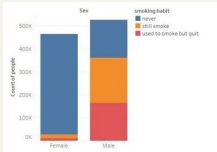
- Predict alcohol consumption.
- Discern the nuances of an individual's smoking habits.
- Predicts one of the three categories: never engaged in smoking, those who once smoked but have since ceased, and active smokers.
- Uncover underlying patterns or trends influencing drinking and smoking behaviors.
- Benefits: Insurance Companies, awareness campaigns, etc.

Data Exploration

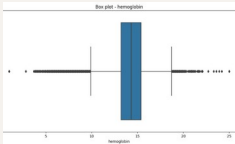
age vs. smoking habit



Data Exploration



Data Exploration



Before Preprocessing

	count	mean	std	min	25%	50%	75%	max
age	602665.0	47.464638	14.262872	20.0	25.0	45.0	65.0	85.0
height	602665.0	162.866017	6.268359	130.0	155.0	165.0	175.0	190.0
weight	602665.0	62.891652	12.105140	25.0	55.0	65.0	75.0	125.0
restingtime	602665.0	60.649154	9.373173	40.0	74.0	81.0	87.0	115.0
right_left	602665.0	0.854767	0.341642	0.1	0.7	1.0	1.2	2.5
right_right	602665.0	0.853218	0.308438	0.1	0.7	1.0	1.2	2.5
hear_left	602665.0	1.830008	0.173664	1.0	1.0	1.0	1.0	2.0
hear_right	602665.0	1.829638	0.170419	1.0	1.0	1.0	1.0	2.0
SBP	602665.0	121.478792	13.654997	70.0	111.0	120.0	130.0	165.0
DBP	602665.0	75.443728	9.307915	47.0	70.0	75.0	81.0	104.0
BLDR	602665.0	87.630817	15.226262	30.0	68.0	85.0	104.0	172.0
tot_chole	602665.0	185.796205	35.743309	81.0	169.0	182.0	217.0	310.0
HDL_chole	602665.0	57.858762	14.204609	12.0	47.0	55.0	65.0	102.0
LDL_chole	602665.0	112.692138	32.663760	12.0	90.0	111.0	134.0	214.0
triglyceride	602665.0	120.180348	87.265137	1.0	72.0	103.0	150.0	409.0
hemoglobin	602665.0	14.201064	1.488754	11.1	13.1	14.2	15.3	18.0
urine_protein	602665.0	1.380163	0.362660	1.0	1.0	1.0	1.0	6.0
crum_creatinine	602665.0	0.848162	0.169907	0.1	0.7	0.8	1.0	2.2
SGOT_AST	602665.0	24.841091	6.328965	1.0	15.0	22.0	27.0	93.0
SGOT_ALT	602665.0	22.773251	12.288759	1.0	14.0	19.0	28.0	79.0

After Preprocessing

Initial Interesting Observations

- Age bracket 40-60 has most drinkers.
- Fewer Females engage in Smoking compared to Men.
- Similar to Alcohol, age 40-60 has most smokers.
- BMI, Age, and Gender, Hemoglobin somewhat influence the smoking and drinking habits the most.



Model Training Pipeline:

Basic Model Training Pipeline:

Creating a Vector using Vector Assembler



Scaling the Vector Created



Parameters Tuning



Model Training



Machine Learning

Logistic Regression

Predicts the probability of a multinomial outcome using a linear combination of features, ideal for classification tasks.

Smoking Prediction Accuracy:
63%

Drinking Prediction Accuracy:
70%



Machine Learning

Decision Tree

Hierarchical model for predictive analysis, visually representing decisions and outcomes through a tree-like structure.

Smoking Prediction Accuracy:

68%

Drinking Prediction Accuracy:

69%



Machine Learning

Random Forest

Hierarchical model for predictive analysis, visually representing decisions and outcomes through a tree-like structure.

Smoking Prediction Accuracy:
66%

Drinking Prediction Accuracy:
69%





Problems Encountered & Future Scope

- Gradient Boosting Algorithm just predicts binary outcomes; however, smoking column has 3 prediction outcomes.
- Biased data to some extent regarding gender.
- Removed the 'sex' variable to avoid model bias due to disproportionate gender representation.
- Focused on identifying patterns based on a broader range of equitable and informative features for improved model fairness and generalizability.

Conclusion

- **Best performing model:**
Decision Tree Accuracy: 68%
- **Drinking Prediction Model:**
Best performing model: Logistic Regression,
Accuracy: 70%
- **Key Feature in Smoking**
Prediction: Hemoglobin identified as crucial feature
in classifying smoking behavior.
- Assists the insurance companies and
health management organization to take
smart decisions to increase business.



Thanks!

Questions?

