

Topic(s): Decision Tree & Random Forest

- 1.) A cloth manufacturing company is interested to know about the segment or attributes contributing to high sale. Approach - A decision tree & random forest model can be built with target variable 'Sale' (we will first convert it into categorical variable) & all other variables will be independent in the analysis.

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban	US
1	9.50	138	73	11	276	120	Bad	42	17	Yes	Yes
2	11.22	111	48	16	260	83	Good	65	10	Yes	Yes
3	10.06	113	35	10	269	80	Medium	59	12	Yes	Yes
4	7.40	117	100	4	466	97	Medium	55	14	Yes	Yes
5	4.15	141	64	3	340	128	Bad	38	13	Yes	No
6	10.81	124	113	13	501	72	Bad	78	16	No	Yes
7	6.63	115	105	0	45	108	Medium	71	15	Yes	No
8	11.85	136	81	15	425	120	Good	67	10	Yes	Yes
9	6.54	132	110	0	108	124	Medium	76	10	No	No
10	4.69	132	113	0	131	124	Medium	76	17	No	Yes
11	9.01	121	78	9	150	100	Bad	26	10	No	Yes
12	11.96	117	94	4	503	94	Good	50	13	Yes	Yes
13	3.98	122	35	2	393	136	Medium	62	18	Yes	No
14	10.96	115	28	11	29	86	Good	53	18	Yes	Yes
15	11.17	107	117	11	148	118	Good	52	18	Yes	Yes
16	8.71	149	95	5	400	144	Medium	76	18	No	No
17	7.58	118	32	0	284	110	Good	63	13	Yes	No

2.) Divide the data (Diabetes) into training and test datasets and create a Random Forest Model to classify 'Class Variable'.

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
1	6	148	72	35	0	33.6	0.627	50	1
2	1	85	66	29	0	26.6	0.351	31	0
3	8	183	64	0	0	23.3	0.672	32	1
4	1	89	66	23	94	28.1	0.167	21	0
5	0	137	40	35	168	43.1	2.288	33	1
6	5	116	74	0	0	25.6	0.201	30	0
7	3	78	50	32	88	31.0	0.248	26	1
8	10	115	0	0	0	35.3	0.134	29	0
9	2	197	70	45	543	30.5	0.158	53	1
10	8	125	96	0	0	0.0	0.232	54	1
11	4	110	92	0	0	37.6	0.191	30	0
12	10	168	74	0	0	38.0	0.537	34	1
13	10	139	80	0	0	27.1	1.441	57	0
14	1	189	60	23	846	30.1	0.398	59	1
15	5	166	72	19	175	25.8	0.587	51	1
16	7	100	0	0	0	30.0	0.484	32	1
17	0	118	84	47	230	45.8	0.551	31	1

- 3.) Use decision trees & random forest algorithm to prepare a model on fraud data, treating those who have taxable_income ≤ 30000 as "Risky" and others are "Good".

	Undergrad	Marital.Status	Taxable.Income	City.Population	Work.Experience	Urban
1	NO	Single	68833	50047	10	YES
2	YES	Divorced	33700	134075	18	YES
3	NO	Married	36925	160205	30	YES
4	YES	Single	50190	193264	15	YES
5	NO	Married	81002	27533	28	NO
6	NO	Divorced	33329	116382	0	NO
7	NO	Divorced	83357	80890	8	YES
8	YES	Single	62774	131253	3	YES
9	NO	Single	83519	102481	12	YES
10	YES	Divorced	98152	155482	4	YES
11	NO	Single	29732	102602	19	YES
12	NO	Single	61063	94875	6	YES
13	NO	Divorced	11794	148033	14	YES
14	NO	Married	61830	86649	16	YES
15	NO	Married	64070	57529	13	YES
16	NO	Divorced	69869	107764	29	NO
17	YES	Divorced	24987	34551	29	NO

Hints:

1. Business Problem
 - 1.1. Objective
 - 1.2. Constraints (if any)
2. Data Pre-processing
 - 2.1 Data cleaning, Feature Engineering, EDA etc.
3. Model Building
 - 3.1 Partition the dataset
 - 3.2 Model(s) - Reasons to choose any algorithm
 - 3.3 Model(s) Improvement steps
 - 3.4 Model Evaluation
 - 3.5 Python and R codes
4. Deployment
 - 4.1 Deploy solutions using R shiny and Python Flask.
5. Result Share the benefits/impact of the solution - how or in what way the business (client) gets benefit from the solution provided.

Note:

1. For each assignment the solution should be submitted in the format
2. Research and Perform all possible steps for improving the model(s) accuracy
Ex: Feature Engineering, Hyper Parameter tuning etc.
3. All the codes (executable programs) are running without errors
4. Documentation of the module should be submitted along with R & Python codes, elaborating on every step mentioned here