

ASSIGNMENT 3 SOLUTION

Done By: Nayan Man Singh Pradhan

Exercise 1. (Modelling inputs / outputs)

In this exercise you work with two data sets:

- Statlog (Shuttle) Data Set
- Computer Hardware Data Set,

which are both available in the [UCI Machine Learning Repository](#). For each of them, perform the following tasks:

- Briefly describe the data set and all involved variables in your own words. If some information is missing on the UCI Repository site, do your own search for these details.
- Model the data set via input and output random variables / vectors.
- Formulate a question that can be solved using machine learning on this data set and give the type of machine learning (supervised / unsupervised / regression / classification) that will allow to answer the question.

(8 Points)

• Statlog (Shuttle) Data Set

- The Statlog (Shuttle) Data Set is a multivariate data set consisting of 9 attributes (columns) and 58,000 instances (rows) of integer values provided by Jason Catlett from NASA's shuttle datasets. According to the Data Set documentation, the shuttle dataset contains 9 attributes, all of which are numerical values. The first attribute is "time", and the last column is the class / label. There are in total 7 classes / labels. They are:

- | | |
|--------------|--------------|
| 1. Rad Flow | 5. Bypass |
| 2. Fpv Close | 6. Bpv Close |
| 3. Fpv Open | 7. Bpv Open |
| 4. High | |

The documentation specifies that approximately 80% of the data belongs to class 1. I have attached a labelled screenshot of the downloaded data file for clarification.

(time)	(Attributes)									(class/label)
1	2	3	4	5	6	7	8	9	10	
55	0	81	0	-6	11	25	88	64	4	← High
56	0	96	0	52	-4	40	44	4	4	← High
50	-1	89	-7	50	0	39	40	2	1	← Rad Flow
53	9	79	0	42	-2	25	37	12	4	← High
55	2	82	0	54	-6	26	28	2	1	← Rad Flow
41	0	84	3	38	-4	43	45	2	1	← Rad Flow
37	0	100	0	36	-8	63	64	2	1	
46	0	83	0	46	0	37	36	0	1	
44	0	79	0	42	-17	35	37	2	1	
44	-1	78	0	44	0	34	34	0	1	
55	0	81	0	54	-10	25	26	2	1	
56	0	95	0	52	-3	40	44	4	4	
37	0	100	0	34	6	64	67	4	1	
42	2	86	2	44	14	42	42	0	1	

b) Input :

$$X : \mathcal{D} \rightarrow \mathbb{R}^9$$

Output

$$G : \mathcal{D} \rightarrow \{1, 2, 3, 4, 5, 6, 7\}$$

c) Supervised Classification Task:

Given 9 features/attributes of the Statlog shuttle, predict the class (rad flow, fpv close, fpv open,.....) of the shuttle.

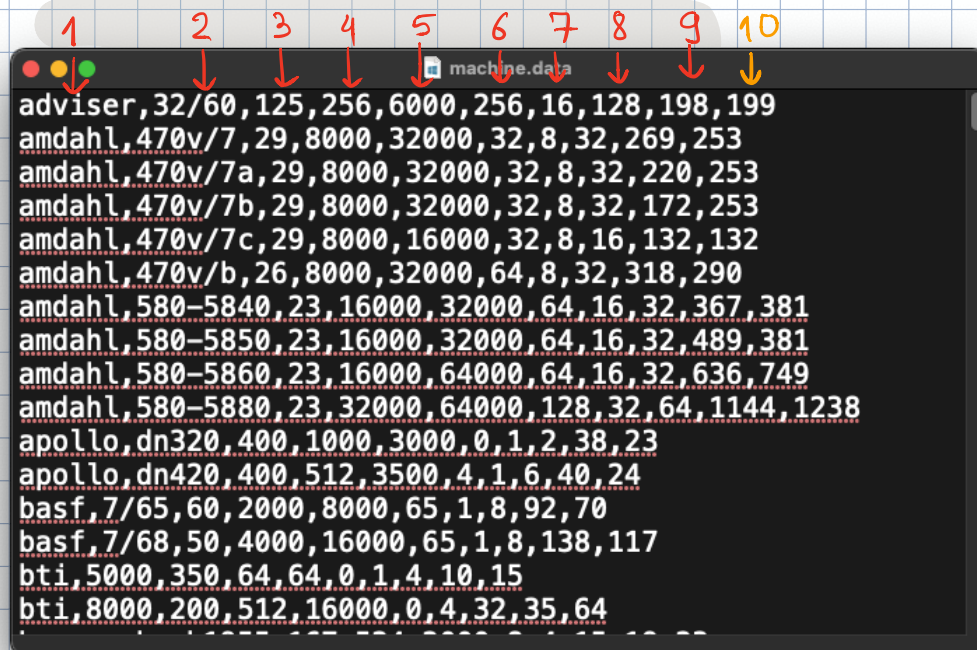
• Computer Hardware Data Set

a) The Computer Hardware Data Set is a multivariate data set consisting of 9 attributes (columns) and 209 instances (rows) of integer values provided by Phillip Ein-Dor and Jacob Feldmesser. According to the Data Set documentation, out of the 9 attributes, 2 are non-predictive, 6 are predictive and 1 is goal field. The last column is the estimated relative performance estimated by the authors using linear regression method. There are no missing

attributes. A description of the attributes :

1. Vendor name: Total of 30 vendor names : adviser, amdahl, apollo,...
2. Model name: Many unique symbols
3. MYCT : Machine Cycle time in nanoseconds (integer)
4. MMIN : Minimum main memory in kilobytes (integer)
5. MMAX : Maximum main memory in kilobytes (integer)
6. CACH : Cache memory in kilobytes (integer)
7. CHMIN : Minimum channels in units (integer)
8. CHMAX : Maximum channels in units (integer)
9. PRP : Published relative performance (integer)
10. ERP : Estimated relative performance by the authors using linear regression (integer)

(attributes)



1	2	3	4	5	6	7	8	9	10
adviser	32	60	125	256	6000	256	16	128	198,199
amdahl	470v/7	29,8000	32000	32	8	32	269	253	
amdahl	470v/7a	29,8000	32000	32	8	32	220	253	
amdahl	470v/7b	29,8000	32000	32	8	32	172	253	
amdahl	470v/7c	29,8000	16000	32	8	16	132	132	
amdahl	470v/b	26,8000	32000	64	8	32	318	290	
amdahl	580-5840	23,16000	32000	64	16	32	367	381	
amdahl	580-5850	23,16000	32000	64	16	32	489	381	
amdahl	580-5860	23,16000	64000	64	16	32	636	749	
amdahl	580-5880	23,32000	64000	128	32	64	1144	1238	
apollo	dn320	400	1000	3000	0	1	2	38	23
apollo	dn420	400	512	3500	4	1	6	40	24
basf	7/65	60	2000	8000	65	1	8	92	70
basf	7/68	50	4000	16000	65	1	8	138	117
bti	5000	350	64	64	0	1	4	10	15
bti	8000	200	512	16000	0	4	32	35	64

b) Input:

$$X : \Omega \rightarrow \mathbb{R}^9$$

Output:

$$Y : \Omega \rightarrow \mathbb{R}_{\geq 0}^K$$

c) Supervised Regression Task:

Given 9 features/attributes of the hardware of a computer, predict the estimated relative performance of the computer.

Exercise 2. (SPAM e-mail representation)

The **Spambase Data Set** is a SPAM classification data set that has exactly the 57 input variables that are roughly described in Example 2.9 of the lecture.

- a) In Example 2.9, we did not mention the specific words and characters that are used in the features. Use the information of the UCI Repository to make a complete description of these variables, i.e. give all the key words, etc.
- b) Search the web for alternative features that can be used to describe (SPAM) emails. Pick one example feature set, cite the source, and describe these features.

(4 Points)