

ML Assignment 02

Name: Kiran K. Patil

Sem: 06

ID: 211070904

course: ML Lab

Aim: Implement and demonstrate the Find-s algorithm for finding the most specific hypothesis based on a given set of training data samples. Read the training data from a .csv file.

Theory :

What is Find-s Algorithm in Machine learning:

In order to understand Find-s algorithm, you need to have a basic idea of the following as concepts as well.

1. Concept learning :
2. General hypothesis
3. Specific hypothesis.

1. concept learning :-

lets try to understand concept learning with an example. Most of the humans learning is based on past instances or experiences. For example, we are able to identify any type of vehical based on a certain features like make, model etc that are

defined over a large set of features.

These special features differentiate the set of cars, trucks etc from the larger set of vehicles. These features that define the set of cars, trucks etc are known as concepts. So the concept learning can be formulated as a problem of searching through a predefined space of potential hypothesis for the hypothesis that best fits the training example.

In simple words, concept learning is inferring a boolean-valued function from training examples of its input and output.

2. General hypothesis:

Hypothesis in general is an explanation of something. The hypothesis basically states that general relationship between the major variables. For example, a general hypothesis for ordering food would be "I want a burger".
$$G = \{ '?', '?', '?', '?' \}$$

3. Specific hypothesis:

The specific hypothesis fill in all the important details about the variables given in the general hypothesis. The more specific details into the example given above would be "I want a cheeseburger with lot of lettuce."

$$S = \{ '\emptyset', '\emptyset', '\emptyset', '\emptyset' \}$$

FIND-S Algorithm:

1. Initialize h to the most specific hypothesis in H

2. For each positive training instance x

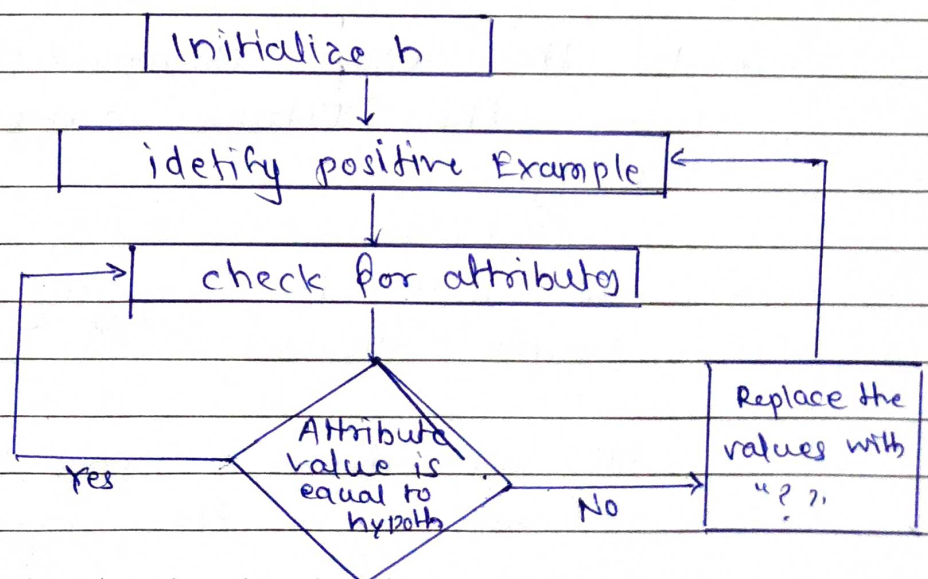
For each attribute constraint a_i in h

if the constraint is satisfied by x
Then do nothing

Else replace a_i in h by the next
more general constraint that is satisfied
by x

3. Output hypothesis h .

How does it work:



• Limitations of Find-s Algorithm:

- There are few limitations of the find-s algorithm listed down below.

1. There is no way to determine if the hypothesis is consistent throughout the data.
2. Inconsistent training set can actually mislead the find-s algorithm, since it ignores the negative examples.
3. Find-s algorithm doesn't provide a backtracking technique to determine the best possible changes that could be done to improve the resulting hypothesis.

• Operations with the dataset.

- In this example / Experiment we are using the "Titanic" dataset, available at Kaggle.
- We will be applying find-s on the Titanic dataset.
- If we directly apply the algorithm then it will give the hypothesis $\langle '?', '?', '?' \rangle$.

- So in order to get some meaningful hypothesis we have to filter the data accordingly.

1. Removing the unnecessary attributes / columns.

- for removing unnecessary columns will use.
- `data.drop('passengerId', inplace=True, axis=1)`
- likewise we will remove, Passenger ID, Name, ticket, fare,
- Now we are having the dataset as follows

Survived	Sex	Age	Adoles	sibsp	Parch	Cabin	Embarked
0	male	14	3	0	0	NaN	S
1	Female	30	3	0	0	NaN	S
0	male	28	2	1	1	NaN	S
1	female	18	3	0	0	NaN	C

2. The for above data-set when find-s algorithm applied by reducing the num of rows.

- `data = data[:25]`
- we get hypothesis as :

`< 1, '?', 'Female', '?', '?', '0', '?', '?' >`

Further applied some more functions.

3. Sort the data according to age.

```
data-frame.sort_values('Age')
```

4. Drop the Null value records.

```
valid-df = data-frame.dropna()
```

5. Removing Rows / unnecessary records.

```
sorted-df = valid-df[:25]
```

6. Group by age (created the class intervals for the age)

Survived	places	sex	Age	Parch	Cabin	Emb.
1	1	female	30.0	1	C54	C
1	1	female	30.0	0	B46	S
0	1	male	30.0	0	B24	G
0	1	male	30.0	0	C31	S
0	1	male	40.0	0	B45	S

7. After applying Rind-s on this filtered data got the hypothesis as follows.

$\langle 1, 1, 'female', 30, '?', '?', '?', '?' \rangle$

Conclusion: Thus From This Experiment we have applied Rind-s Algorithm on the Titanic Dataset.

- Without Filtering when we had applied the algo. we got hypothesis as $\langle ?, ?, ?, ?, ?, ? \rangle$
- After applying some filters on the dataset and removing unnecessary columns we got the hypothesis as $\langle 1, '?', 'female', '?', '?', '0', '?', ? \rangle$
- Then applied grouping and classified the data according to age group. and applied Rind-s algorithm got the hypothesis as $\langle 1, 1, 'female', 30, '?', '?', '?', '?' \rangle$
- So from above observation in all three hypothesis we have got female as common attribute thus mostly females are survived.