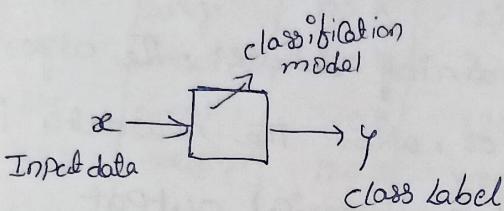


UNIT - IV

Classification



It is the task of assigning objects to one of several predefined categories.

It is the task of learning a target function T that maps each attribute set to one of predefined classification function.

Purpose of classification model:-

Descriptive modeling :-

A classification technique (or classifier) is a systematic approach to build classification models from an input data set. It is to identify category (or) class label.

In classification data is grouped into categories based on a training data set. Using the training data set, the algorithm derives a model or the classifier, the derived model can be a decision tree.

In test data set usually unlabeled data is given to the model. Then it finds the class to which it belongs.

e.g. Computer sales (Y/N)

Prediction:-

Prediction is the process of identifying the missing or unavailable numerical data for a new observation.

According to the training dataset, the algorithm derives the model or a predictor, when the new data is given the model should blend a numerical output.

Unlike in classification, this method does not have the class label. The model predicts a continuous-valued function or ordered values.

Eg:- computer sales (How much he/she spend)

Classification eg:-

computer sales - profession, salary, requirement - Y/N

Customer loan - salary designation - ?

Diagnosis - Fever, Cold, Cough, tests - ?

online shopping - Previous search, item interested.

Illustrating Classification Task

Building a classification model (General approach)

Tid	Attri1	Attri2	Attri3	class
1	Yes	large	125K	NO
2	NO	medium	100K	NO
3	NO	small	70K	NO
4	Yes	medium	120K	YES
5	NO	large	95K	YES
6	NO	Medium	60K	NO
7	YES	Large	220K	NO
8	NO	small	85K	YES
9	NO	medium	75K	NO
10	NO	small	90K	YES

learning
Algorithm

Learn
model

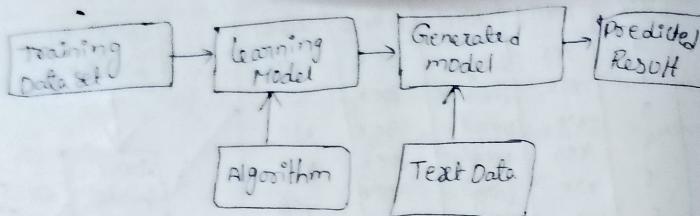
Model

Apply
Model

Deduction *

Tid	Attri4	Attri5	Attri6	class
11	NO	small	55K	?
12	YES	medium	80K	?
13	YES	large	110K	?
14	NO	small	95K	?
15	NO	large	67K	?

General approaches to solving a classification Problem:



Classification is one of the methods used for data analysis. we analyze the data and classify it based on our requirement

e.g.: if we want to know the Performance of the university, we classify the student database based on their Performance as above average, average and below average students.

If the classification shows that the no. of students under "below average" category are more, then the University needs to improve

Data classification process:

Let us consider the data classification where a decision is to be made to increasing the pay scale of employees in organization based on their performance level and current pay scale

Training Data:

Name	Age	Performance	Current-Pay	Like-decision
Ram	30	Good	20,000	Yes
Raj	40	Good	25,000	Yes
Ramya	25	Good	27,000	Yes
Ravi	30	Bad	20,000	No

Classification algorithm

Rules of classification

If Performance = Bad Then Like-decision = No
 If Current-Pay-Scale < 26,000 and Performance = Good Then Like-decision = Yes
 If Age ≥ 30 and Current-Pay-Scale < 30,000 Then Like-decision = Yes.

Test Data:

Name	Age	Performance	Current-Pay-Scale	Like-decision
Vikki	25	Good	20,000	Yes
Yash	35	Good	28,000	Yes
Ajju	35	Good	22,000	Yes
Vaishu	27	Bad	24,000	No

Rules of classification

Rahul, 30, Good, 20,000 Like-decision?
 ↘ increase P2-1

General approach to classification.

- training set consists of records with known class labels
- training set is used to build a classification model
- A labeled test set of previously unseen data records is used to evaluate the quality of the model.
- the classification model is applied to new records with unknown class labels.

Evaluation of classifiers & classification techniques.

A classification technique (or) classifier is a systematic approach to building classification models from an input dataset.

Examples

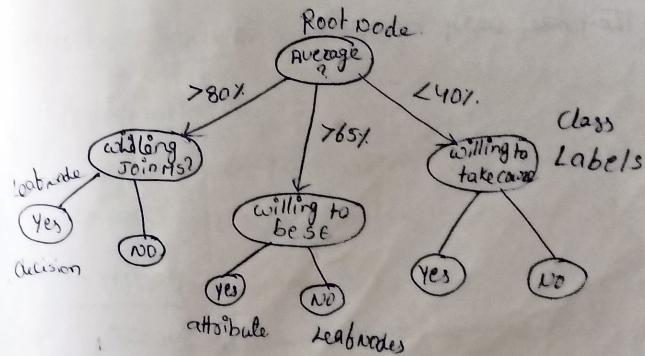
- Decision tree classifier
- Rule based classifier
- Neural Networks
- Support Vector machines
- Naïve Bayes classifier

Each technique employs a learning algorithm to identify a model that best fits the relationship between the attribute set and class label of the input data.

The model generated by a learning algorithm should both fit the input data well and correctly predict the class labels of records it has never seen before.

Decision Trees - Decision tree construction:-

- A decision tree is a tree-structure
- where each non-leaf node represents the test on an attribute
- branches represent the outcome of the test and the test and the leaf nodes represents the class labels.



The decision tree enables the organization to identify the no. of students who are going to join a software company. Some decision trees are binary and some trees are non-binary.

Decision trees are mostly used for classification rules for tuples which don't have class labels identifier for them.

The class predictions can be made by traversing from node to the leaf node.

Advantages of Decision tree :-

- (i) They don't require domain knowledge
- (ii) They are easy to understand.
- (iii) Handles high dimensional data
- (iv) Classification and learning becomes simpler when decision trees are used.
- (v) They are very accurate.

Decision tree uses Information Gain

$$I(P, n) = -\frac{P}{S} \log_2 \frac{P}{S} - \frac{n}{S} \log_2 \frac{n}{S} \quad S = (P, n)$$

$$\text{Entropy} - E(A) = \sum_{i=1}^V \frac{P_i + n_i}{P+n} I(P_i, n_i)$$

$$\text{Gain}(A) = I(P, n) - E(A)$$

$$\left\{ \begin{array}{l} \log_2 = \frac{\log x}{\log 10} \\ \log_2 \end{array} \right\}$$

Decision tree induction:-

→ Flow chart like tree structure

→ supports in taking decisions

→ it defines the rules visually in form of tree

Types of node:-

1. Root node : Main question
2. Branch node : Intermediate processes
3. Leaf node : Answer

Attribute selection measures :-

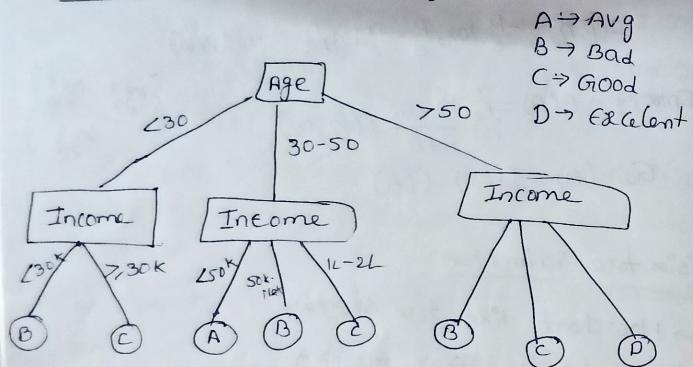
1. Information Gain :-

How much information does the answer to the specific question provide.

2. Entropy :- (IG ↑ ⇔ ↓ Entropy)

Measures the amount of uncertainty in the information

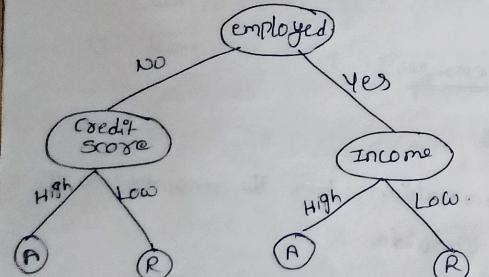
Example: credit score Rating



Dataset \rightarrow Algorithm \rightarrow classifies the data
(decision tree algorithm)

example: Loan system

(decides the loan should be approved or not)



\rightarrow Algorithm (ID3): $\stackrel{\text{(repeatedly)}}{\text{iteratively}}$ $\stackrel{\text{(divide)}}{\text{Divide}}$ $\stackrel{\text{(choose)}}{\text{choose}}$ $\stackrel{\text{(stop)}}{\text{Stop}}$

1. In the given dataset, choose a target attribute
2. Calculate Information gain of target attributes

$$IG = \frac{-P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

3. For remaining attributes, find entropy

$$\text{Entropy} = IG \times \text{Probability}$$

$$E(A) = E \frac{P_i + N_i}{P+N} I(P_i, N_i) \quad (\log_2 = \frac{\log x}{\log_2 10})$$

4. Calculate Gain = $IG - E(A)$

Age	Competition	Type	Profit
old	Yes	S/W	Down
old	No	S/W	Down
old	No	H/W	Down
mid	Yes	S/W	Down
mid	Yes	H/W	Down
mid	No	H/W	Up
mid	No	S/W	Up
new	Yes	S/W	Up
new	No	H/W	Up
new	No	S/W	Up

1453
 1889
 1290
 21
 4682

Step 1: Target Attribute = Profit.

Step 2: Information Gain

$$IG = \frac{-P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$P = \text{Count(Down)} = 5$$

$$N = \text{Count(Up)} = 5$$

$$= -\frac{5}{10} \log_2 \left(\frac{5}{10} \right) - \frac{5}{10} \log_2 \left(\frac{5}{10} \right)$$

$$= -\left(\frac{1}{2} \log_2(2^1) + \frac{1}{2} \log_2(2^1)\right)$$

$$= -\left(\frac{1}{2} \times -1 \log_2 2 + \frac{1}{2} \times -1 \log_2 2\right)$$

$$= -\left(\frac{1}{2} \times -1 + \frac{1}{2} \times -1\right)$$

$$= -(-1 - 1) = 1$$

$$IG = 1$$

Step 3: Calculate entropy for remaining attributes

$$E(A) = \sum_{i=1}^{n-1} P_i \times N_i I(A|N_i) [IG \times \text{probability}]$$

Age: Prepare a table for each attribute

rows: values of undertaken attributes (old, mid, new)

columns: values of target attributes
(down, up)

	down	up
old	3	0
mid	2	2
new	0	3

Entropy = $Ig \times \text{Probability}$

P \rightarrow down count

N \rightarrow up count

$$\begin{aligned} IG(\text{old}) &= -\left(\frac{3}{5} \log_2 \left(\frac{3}{5}\right) + \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right) \\ &= -\left(1 \log_2 1 + 0 \log_2 0\right) = 0 \end{aligned}$$

$$\text{Probability} = \frac{3}{10}$$

$$\text{Entropy} = 0 \times \frac{3}{10} = 0$$

$$\begin{aligned} IG(\text{mid}) &= -\left(\frac{2}{4} \log_2 \left(\frac{2}{4}\right) + \frac{2}{4} \log_2 \left(\frac{2}{4}\right)\right) \\ &= -\left(\frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \\ &= -\left(\frac{1}{2} \times -1 \log_2 2 + \frac{1}{2} \times -1 \log_2 2\right) \\ &= -\left(\frac{1}{2}(-1) + \frac{1}{2}(-1)\right) \\ &= -(-1 - 1) = 1 \end{aligned}$$

$$\text{Probability} = \frac{4}{10} = \text{Entropy(mid)} = 1 \times \frac{4}{10} = 0.4$$

$$\begin{aligned} IG(\text{new}) &= -\left(\frac{0}{3} \log_2 \left(\frac{0}{3}\right) + \frac{3}{3} \log_2 \left(\frac{3}{3}\right)\right) \\ &= 0 \end{aligned}$$

$$\text{Probability} = \frac{3}{10} \quad \text{Entropy} = 0 \times \frac{3}{10} = 0$$

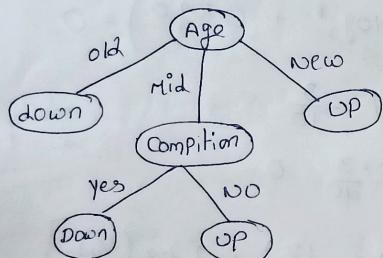
$$\text{Entropy(Age)} = E(O) + E(M) + E(N) = 0 + 0.4 + 0 = 0.4$$

$$④ \text{Gain} = IG - \epsilon(A) = 1 - 0.4 = 0.6$$

$$\text{Gain (Type)} = 0$$

$$\text{Gain (competition)} = 0.124$$

Highest gain \rightarrow root node
(Age)



old \rightarrow all down

mid \rightarrow some down some up

new \rightarrow all up

Competition :-

Prepare table

	UP	down
yes	1	3
mid	4	2
new		

$$\begin{aligned} IG(\text{yes}) &= -\left(\frac{1}{4}\log_2\left(\frac{1}{4}\right) + \left(\frac{3}{4}\right)\log_2\left(\frac{3}{4}\right)\right) \\ &= -\left(\frac{1}{4}\log_2\left(\frac{1}{4}\right) + \left(\frac{3}{4}\right)\log_2\left(\frac{3}{4}\right)\right) \end{aligned}$$

Naïve Bayes classifier :-

→ Classification technique based on Bayes Theorem with an assumption of independence among features.

$$P(A/B) = \frac{P(B/A) P(A)}{P(B)}$$

Fruit = {yellow, sweet, long}

Fruit	yellow	sweet	Long	Total
orange	350	450	0	650
Banana	400	300	350	400
Others	50	100	50	150
Total	800	850	400	1200

$$P(\text{yellow/orange}) = \frac{P(0/Y) P(Y)}{P(O)}$$

$$= \frac{\left(\frac{350}{800}\right) \left(\frac{800}{1200}\right)}{\frac{650}{1200}}$$

$$P(\text{sweet/orange}) = \frac{P(0/S) P(S)}{P(O)} = 0.53$$

$$= 0.69$$

$$P(\text{long/orange}) = 0$$

$$P(\text{Fruit/orange}) = P(Y/O) * P(S/O) * P(L/O)$$

$$= 0.53 * 0.69 * 0$$

$$= 0$$

$$P(\text{Fruit/Banana}) = \frac{P(B/Y) P(Y)}{P(B)}$$

$$= \frac{P\left(\frac{400}{800}\right) \left(\frac{800}{1200}\right)}{P\left(\frac{400}{1200}\right)}$$

$$= \frac{\frac{400}{800} * \frac{8}{12}}{\frac{1}{3}} = \frac{1}{3} * \frac{3}{1} = 1$$

$$P(\text{sweet/Banana}) = \frac{P(B/S) P(S)}{P(B)}$$

$$= \frac{\frac{300}{800} * \frac{850}{1200}}{\frac{400}{1200}} = \frac{1}{4} * \frac{3}{4}$$

$$P(\text{long/Banana}) = \frac{P(B/L) P(L)}{P(B)} = \frac{\frac{350}{800} * \frac{400}{1200}}{\frac{400}{1200}} = \frac{\frac{35}{8} * \frac{3}{1}}{\frac{40}{12}} = \frac{0.875}{0.833} = 1$$

$$P(\text{Fruit/Banana}) = 1 * 0.75 * 0.89 = 0.65$$

$$\begin{aligned} P(\text{Fruit/Others}) &= P(4\%) * P(5\%) * P(1\%) \\ &= 0.33 * 0.66 * 0.33 = 0.072 \end{aligned}$$

Banana having most probability

Confusion Matrix:-

The Confusion matrix is a matrix used to determine the performance of the classification models for a given set of test data. It can only be determined if the true value for test data are known.

The matrix itself can be easily understood but the related terminologies may be confusion.

It shows the errors in the model performance in the form of a matrix, hence also known as an error matrix.

Features of confusion Matrix:-

- For the 2 prediction classes of classifiers, the matrix is of 2×2 table. For three classes, it is 3×3 table.
- The matrix is divided into two dimensions, that are Predictive values and actual values along with the total no. of predictions.
- Predicted values are those values, which are predicted by the model, and actual values are the true values for the given observations.

It Looks Like.

$n = \text{total predictions}$	Predicted: No	Predicted: Yes
Actual: No	True Negative	False Positive
Actual: Yes	False Negative	True Positive

True Negative: Model has given prediction No, and the real or actual value was also No.

True Positive: The model has predicted Yes and the actual value was also true.

False Negative :- The model has predicted no, but the actual value was yes, it is also called as Type-II error.

False Positive :- The model has predicted yes, but the actual value was no, it is also called a type-I error.

Need for Confusion matrix :-

- It evaluates the performance of the classification models when they make predictions on test data and tell how good our classification model is.
- It not only tells the errors made by the classifier but also the type of errors such as it is either type-1 or type-2 errors.
- With the help of the confusion matrix we can calculate the different parameters for the model, such as accuracy, precision etc.

We can perform various calculations for the model, such as the model's accuracy, using the matrix.

Classification accuracy :-

It defines how often the model predicts the correct output. It can be calculated as the ratio of the no. of correct predictions made by the classifier to all no. of predictions made by the classifiers.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

Misclassification rate :- It is also termed as Error rate, and it defines how often the model gives the wrong predictions. The value of error rate can be calculated as the no. of the predictions made by the classifier.

$$\text{Error Rate} = \frac{FP + FN}{\text{Total}(TP + TN + FP + FN)}$$

Precision :- It can be defined as the no. of correct outputs provided by the model out of all positive classes that have predicted correctly by the model. How many of them were actually true.

$$\text{Precision} = \frac{TP}{TP + FP}$$

Recall :-

If it is defined as the out of total positive classes, how our model predicted correctly. The recall must be high as possible.

$$\text{Recall} = \frac{TP}{TP+FN}$$

F-Measure :- If two models have low precision and high recall. It is difficult to compare these models. For this we use F-score. This score helps us to calculate the recall and precision at the same time. The F-score is maximum if the recall is equal to the precision.

$$F\text{-Measure} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$$

Example :-

		Predicted		
		No	Yes	165
Actual	No	50	100	150
	Yes	5	100	105
		55	110	

Accuracy :-

$$\frac{TP+TN}{\text{Total}} = \frac{100+50}{165} = \frac{150}{165} = 0.91$$

Error Rate :- $1 - \text{accuracy} = 1 - 0.91 = 0.09$
(%)

$$\frac{FP+FN}{\text{Total}} = \frac{5+10}{165} = \frac{15}{165} = 0.09$$

Precision :-

$$\frac{TP}{TP+FP} = \frac{100}{100+10} = \frac{100}{110} = 0.91$$

(Predicted Yes)

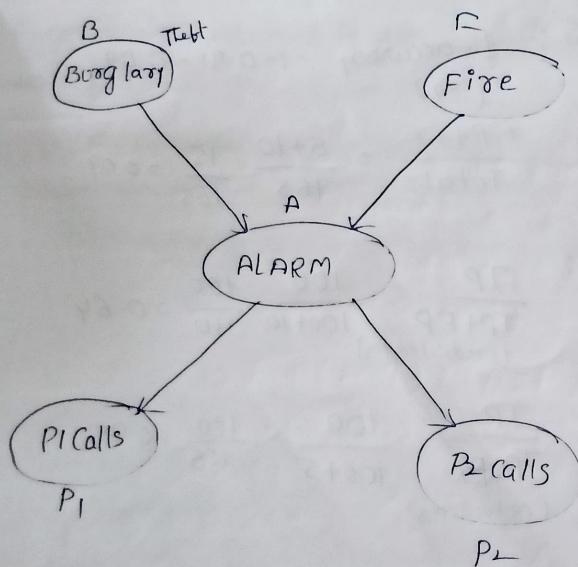
Recall :-

$$\frac{TP}{TP+FN} = \frac{100}{100+5} = \frac{100}{105} = 0.95$$

(Actual Yes)

Bayesian-Belief Networks:-

Bayesian-Belief network is a graphical representation of different probabilistic relationships among random variables in a particular set. It is a classifier with no dependency on attributes i.e. it is condition independent. Due to its feature of joint probability, the probability in Bayesian Belief network is derived based on a condition - $P(\text{attribute}/\text{parent})$ i.e. Probability of an attribute, true over Parent-attribute.



We have an alarm 'A' - a node, say installed in a house of a person 'gfg' which rings upon two probabilities i.e. burglary 'B' and fire 'F', which are - parent nodes of the alarm node. The alarm is the parent node of two probabilities P_1 calls 'p1' & P_2 calls p_2 Person nodes.

→ Upon the instance of burglary and fire. ' P_1 ' and ' P_2 ' call person 'gfg' respectively. But there are few drawbacks in this case, as sometimes ' P_1 ' may forget to call the person 'gfg' even after hearing the alarm, as he has a tendency to forget things quick. Similarly ' P_2 ' sometimes fails to call the person 'gfg', as he is only able to hear the alarm from a certain distance.

Q) Find the Probability that 'P₁' is true (P₁ has called 'gfg') 'P₂' is true (P₂ has called 'gfg') when the alarm 'A' is rung, but no burglary 'B' and fire 'F' has occurred.

$\Rightarrow P(P_1, P_2, A, \sim B, \sim F)$ [where -P₁, P₂ & A are true events and $\sim B$ & $\sim F$ are false events]

Burglary 'B':-

$P(B=T) = 0.001$ ("B" is true i.e. burglary has occurred)
 $P(B=F) = 0.999$ ("B" is false i.e. burglary has not occurred)

Fire 'F':-

$P(F=T) = 0.002$ ("F" is true i.e. fire has occurred)

$P(F=F) = 0.998$ ("F" is False i.e. fire has not occurred)

Alarm A :- $B \quad F \quad P(A=T) \quad P(A=F)$

$\frac{B}{T}$	$\frac{F}{T}$	$\frac{P(A=T)}{0.95}$	$\frac{P(A=F)}{0.05}$
---------------	---------------	-----------------------	-----------------------

T	F	0.94	0.06
---	---	------	------

F	T	0.29	0.71
---	---	------	------

F	F	0.001	0.999
---	---	-------	-------

The alarm 'A' node can be 'true' or 'false' (i.e. may have rung or may not have rung). It has two Parent nodes burglary 'B' and Fire 'F', which can be 'true' or 'false' (i.e. may have occurred or may not have occurred) depending upon different conditions.

Person 'P₁':-

A	$P(P_1=T)$	$P(P_1=F)$
---	------------	------------

T	0.95	0.05
---	------	------

F	0.05	0.95
---	------	------

The Person 'P₁' node can be 'true' or 'false' (i.e. may have called the person 'gfg' or not). It has a Parent node, The alarm 'A', which can be 'true' or 'false' (i.e. may have rung or may not have rung, upon burglary 'B' or fire 'F').

Person 'P₂':-

A	$P(P_2=T)$	$P(P_2=F)$
---	------------	------------

T	0.80	0.20
---	------	------

F	0.01	0.99
---	------	------

The Person 'P₂' node can be 'true' or 'false'
 (i.e. may have called the person 'Fry' or not)
 It has a parent node, The Alarm 'A' which
 can be 'true' (or) 'false' (i.e. may have rung
 & may not have rung, upon burglary B & fire F).

Solution:- Consider the observed probabilistic scan
 with respect to the question - $P(P_1, P_2, A, \sim B, \sim F)$
 we need to get the probability of 'P₁'. we find
 it with regard to its parent node - alarm
 'A'.

alarm 'A' To get the probability of 'P₂' we
 bind it with regard to its Parent node - alarm
 'A'.

we find the probability of alarm 'A' node with
 regard to ' $\sim B$ ' & ' $\sim F$ ' since burglary 'B' and
 fire 'F' are Parent nodes of alarm 'A'

From the observed probabilistic scan, we can

$$\begin{aligned}
 & P(P_1, P_2, A, \sim B, \sim F) \\
 & = P(P_1/A) * P(P_2/A) * P(A/\sim B, \sim F) * P(\sim B) * \\
 & \quad P(\sim F) \\
 & = 0.95 * 0.30 * 0.001 * 0.999 * 0.998 \\
 & = 0.00075
 \end{aligned}$$