

Machine Learning

Classification Methods

Bayesian Classification, Nearest
Neighbor, Ensemble Methods

Bayesian Classification: Why?



- A statistical classifier: performs *probabilistic prediction*, *i.e.*, predicts class membership probabilities
- Foundation: Based on Bayes' Theorem.
- Performance: A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers
- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct — prior knowledge can be combined with observed data

Bayes' Rule



$$p(h | d) = \frac{P(d | h)P(h)}{P(d)}$$

Understanding Bayes' rule

d = data

h = hypothesis (model)

- rearranging

$$p(h | d)P(d) = P(d | h)P(h)$$

$$P(d, h) = P(d, h)$$

the same joint probability

on both sides

Who is who in Bayes' rule

$P(h)$: prior belief (probability of hypothesis h before seeing any data)

$P(d | h)$: likelihood (probability of the data if the hypothesis h is true)

$P(d) = \sum_h P(d | h)P(h)$: data evidence (marginal probability of the data)

$P(h | d)$: posterior (probability of hypothesis h after having seen the data d)



Example of Bayes Theorem

- Given:
 - A doctor knows that meningitis causes stiff neck 50% of the time
 - Prior probability of any patient having meningitis is $1/50,000$
 - Prior probability of any patient having stiff neck is $1/20$
- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

Choosing Hypotheses



- *Maximum Likelihood* hypothesis:

$$h_{ML} = \arg \max_{h \in H} P(d | h)$$

- Generally we want the most probable hypothesis given training data. This is the *maximum a posteriori* hypothesis:

$$h_{MAP} = \arg \max_{h \in H} P(h | d)$$

- Useful observation: it does not depend on the denominator $P(d)$

Bayesian Classifiers



- Consider each attribute and class label as random variables
- Given a record with attributes (A_1, A_2, \dots, A_n)
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
- Can we estimate $P(C | A_1, A_2, \dots, A_n)$ directly from data?

Bayesian Classifiers



- Approach:
 - compute the posterior probability $P(C \mid A_1, A_2, \dots, A_n)$ for all values of C using the Bayes theorem

$$P(C \mid A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n \mid C) P(C)}{P(A_1 A_2 \dots A_n)}$$

- Choose value of C that maximizes
 $P(C \mid A_1, A_2, \dots, A_n)$
- Equivalent to choosing value of C that maximizes
 $P(A_1, A_2, \dots, A_n \mid C) P(C)$
- How to estimate $P(A_1, A_2, \dots, A_n \mid C)$?



Naïve Bayes Classifier

- Assume independence among attributes A_i when class is given:
 - $P(A_1, A_2, \dots, A_n | C) = P(A_1 | C) P(A_2 | C) \dots P(A_n | C)$
 - Can estimate $P(A_i | C)$ for all A_i and C .
 - This is a simplifying assumption which may be violated in reality
- The Bayesian classifier that uses the Naïve Bayes assumption and computes the MAP hypothesis is called Naïve Bayes classifier

$$c_{Naive\ Bayes} = \arg \max_c P(c)P(\mathbf{x} | c) = \arg \max_c P(c) \prod_i P(a_i | c)$$

How to Estimate Probabilities from Data?



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Class: $P(C) = N_c / N$

- e.g., $P(\text{No}) = 7/10$,
 $P(\text{Yes}) = 3/10$

- For discrete attributes:

$$P(A_i | C_k) = |A_{ik}| / N_{C_k}$$

- where $|A_{ik}|$ is number of instances having attribute A_i and belongs to class C_k
- Examples:

$$P(\text{Status}=\text{Married} | \text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes} | \text{Yes})=0$$

How to Estimate Probabilities from Data?



- For continuous attributes:
 - **Discretize** the range into bins
 - one ordinal attribute per bin
 - violates independence assumption
 - **Two-way split:** $(A < v)$ or $(A > v)$
 - choose only one of the two splits as new attribute
 - **Probability density estimation:**
 - Assume attribute follows a normal distribution
 - Use data to estimate parameters of distribution (e.g., mean and standard deviation)
 - Once probability distribution is known, can use it to estimate the conditional probability $P(A_i | c)$

How to Estimate Probabilities from Data?



Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:

$$P(A_i | c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each (A_i, c_i) pair
- For (Income, Class=No):
 - If Class=No
 - sample mean = 110
 - sample variance = 2975

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi(54.54)}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

Naïve Bayesian Classifier: Training Dataset



Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

New Data:

X = (age <=30,

Income = medium,

Student = yes

Credit_rating = Fair)

age	income	student	credit_rating	comp
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Naïve Bayesian Classifier: An Example



Given X (age=youth, income=medium, student=yes, credit=fair)

Maximize $P(X|C_i)P(C_i)$, for $i=1,2$

First step: Compute $P(C)$ The prior probability of each class can be computed based on the training tuples:

$$P(\text{buys_computer=yes})=9/14=0.643$$

$$P(\text{buys_computer=no})=5/14=0.357$$

Naïve Bayesian Classifier: An Example



Given X (age=youth, income=medium, student=yes, credit=fair)

Maximize $P(X|C_i)P(C_i)$, for $i=1,2$

Second step: compute $P(X|C_i)$

$$\begin{aligned} P(X|\text{buys_computer=yes}) &= P(\text{age=youth}|\text{buys_computer=yes}) \times \\ &\quad P(\text{income=medium}|\text{buys_computer=yes}) \times \\ &\quad P(\text{student=yes}|\text{buys_computer=yes}) \times \\ &\quad P(\text{credit_rating=fair}|\text{buys_computer=yes}) \\ &= 0.044 \end{aligned}$$

$$P(\text{age=youth}|\text{buys_computer=yes}) = 0.222$$

$$P(\text{income=medium}|\text{buys_computer=yes}) = 0.444$$

$$P(\text{student=yes}|\text{buys_computer=yes}) = 6/9 = 0.667$$

$$P(\text{credit_rating=fair}|\text{buys_computer=yes}) = 6/9 = 0.667$$

Naïve Bayesian Classifier: An Example



Given X (age=youth, income=medium, student=yes, credit=fair)

Maximize $P(X|C_i)P(C_i)$, for $i=1,2$

Second step: compute $P(X|C_i)$

$$\begin{aligned} P(X|\text{buys_computer=no}) &= P(\text{age=youth}|\text{buys_computer=no}) \times \\ &\quad P(\text{income=medium}|\text{buys_computer=no}) \times \\ &\quad P(\text{student=yes}|\text{buys_computer=no}) \times \\ &\quad P(\text{credit_rating=fair}|\text{buys_computer=no}) \\ &= 0.019 \end{aligned}$$

$$P(\text{age=youth}|\text{buys_computer=no}) = 3/5 = 0.666$$

$$P(\text{income=medium}|\text{buys_computer=no}) = 2/5 = 0.400$$

$$P(\text{student=yes}|\text{buys_computer=no}) = 1/5 = 0.200$$

$$P(\text{credit_rating=fair}|\text{buys_computer=no}) = 2/5 = 0.400$$

Naïve Bayesian Classifier: An Example



Given X (age=youth, income=medium, student=yes, credit=fair)

Maximize $P(X|C_i)P(C_i)$, for $i=1,2$

We have computed in the first and second steps:

$$P(\text{buys_computer}=\text{yes})=9/14=0.643$$

$$P(\text{buys_computer}=\text{no})=5/14=0.357$$

$$P(X|\text{buys_computer}=\text{yes})= 0.044$$

$$P(X|\text{buys_computer}=\text{no})= 0.019$$

Third step: compute $P(X|C_i)P(C_i)$ for each class

$$P(X|\text{buys_computer}=\text{yes})P(\text{buys_computer}=\text{yes})=0.044 \times 0.643=0.028$$

$$P(X|\text{buys_computer}=\text{no})P(\text{buys_computer}=\text{no})=0.019 \times 0.357=0.007$$

The naïve Bayesian Classifier predicts **X belongs to class (“buys_computer = yes”)**

Example



Training set :

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

k

Given a Test Record:

$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$

Example of Naïve Bayes Classifier



Given a Test Record:

$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$

naive Bayes Classifier:

$$P(\text{Refund}=\text{Yes}|\text{No}) = 3/7$$

$$P(\text{Refund}=\text{No}|\text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$$

$$P(\text{Refund}=\text{No}|\text{Yes}) = 1$$

$$P(\text{Marital Status}=\text{Single}|\text{No}) = 2/7$$

$$P(\text{Marital Status}=\text{Divorced}|\text{No}) = 1/7$$

$$P(\text{Marital Status}=\text{Married}|\text{No}) = 4/7$$

$$P(\text{Marital Status}=\text{Single}|\text{Yes}) = 2/7$$

$$P(\text{Marital Status}=\text{Divorced}|\text{Yes}) = 1/7$$

$$P(\text{Marital Status}=\text{Married}|\text{Yes}) = 0$$

For taxable income:

If class=No: sample mean=110

sample variance=2975

If class=Yes: sample mean=90

sample variance=25

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No})$
 $\times P(\text{Married}|\text{Class}=\text{No})$
 $\times P(\text{Income}=120\text{K}|\text{Class}=\text{No})$
 $= 4/7 \times 4/7 \times 0.0072 = 0.0024$
- $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{Class}=\text{Yes})$
 $\times P(\text{Married}|\text{Class}=\text{Yes})$
 $\times P(\text{Income}=120\text{K}|\text{Class}=\text{Yes})$
 $= 1 \times 0 \times 1.2 \times 10^{-9} = 0$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$

$\Rightarrow \text{Class} = \text{No}$

Avoiding the 0-Probability Problem



- If one of the conditional probability is zero, then the entire expression becomes zero
- Probability estimation:

$$\text{Original : } P(A_i | C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace : } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

$$\text{m - estimate : } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m}$$

c: number of classes

p: prior probability

m: parameter

Naïve Bayes (Summary)



- Advantage
 - Robust to isolated noise points
 - Handle missing values by ignoring the instance during probability estimate calculations
 - Robust to irrelevant attributes
- Disadvantage
 - Assumption: class conditional independence, which may cause loss of accuracy
 - Independence assumption may not hold for some attribute. Practically, dependencies exist among variables
 - Use other techniques such as Bayesian Belief Networks (BBN)

Remember

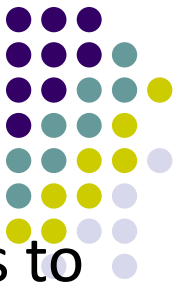


- Bayes' rule can be turned into a classifier
- Maximum A Posteriori (MAP) hypothesis estimation incorporates prior knowledge; Max Likelihood (ML) doesn't
- Naive Bayes Classifier is a simple but effective Bayesian classifier for vector data (i.e. data with several attributes) that assumes that attributes are independent given the class.
- Bayesian classification is a generative approach to classification

Classification Paradigms

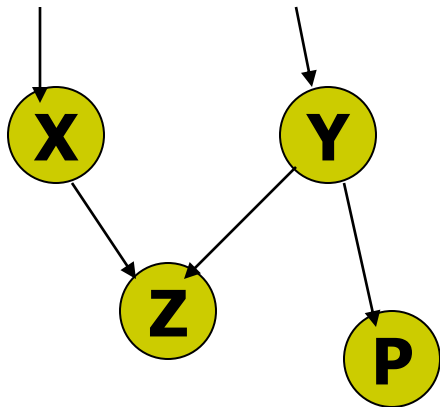


- In fact, we can categorize three fundamental approaches to classification:
- **Generative models:** Model $p(x|C_k)$ and $P(C_k)$ separately and use the Bayes theorem to find the posterior probabilities $P(C_k|x)$
 - E.g. Naive Bayes, Gaussian Mixture Models, Hidden Markov Models,...
- **Discriminative models:**
 - Determine $P(C_k|x)$ directly and use in decision
 - E.g. Linear discriminant analysis, SVMs, NNs,...
- Find a **discriminant function** f that maps x onto a class label directly without calculating probabilities

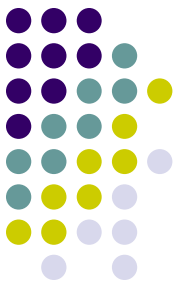


Bayesian Belief Networks

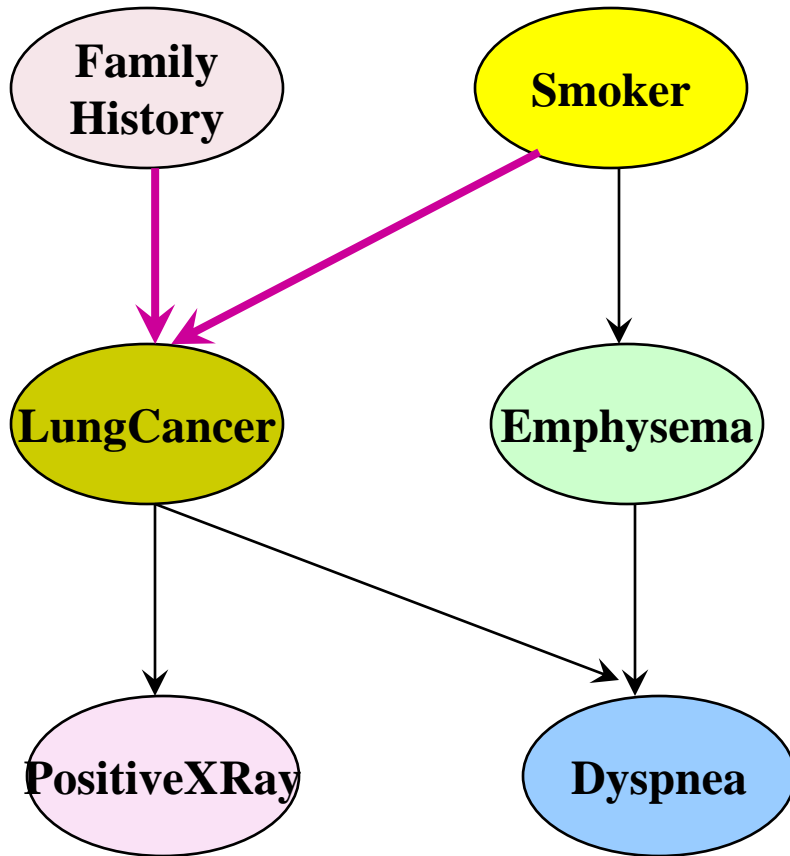
- Bayesian belief network allows a *subset* of the variables to be conditionally independent
- A graphical model of causal relationships
 - Represents dependency among the variables
 - Gives a specification of joint probability distribution



- ❑ **Nodes:** random variables
- ❑ **Links:** dependency
- ❑ X and Y are the parents of Z, and Y is the parent of P
- ❑ No dependency between Z and P
- ❑ Has no loops or cycles



Bayesian Belief Network: An Example



The **conditional probability table (CPT)** for variable LungCancer:

	(FH, S)	(FH, ~S)	(~FH, S)	(~FH, ~S)
LC	0.8	0.5	0.7	0.1
~LC	0.2	0.5	0.3	0.9

CPT shows the conditional probability for each possible combination of its parents

Derivation of the probability of a particular combination of values of **X**, from CPT:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{Parents}(Y_i))$$

Bayesian Belief Networks

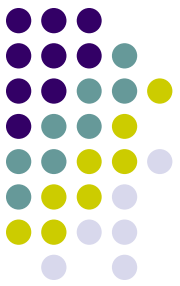
Bayesian network through an example by creating a directed acyclic graph:



- Example: Harry installed a new burglar alarm at his home to detect burglary. The alarm reliably responds at detecting a burglary but also responds for minor earthquakes. Harry has two neighbors David and Sophia, who have taken a responsibility to inform Harry at work when they hear the alarm. David always calls Harry when he hears the alarm, but sometimes he got confused with the phone ringing and calls at that time too. On the other hand, Sophia likes to listen to high music, so sometimes she misses to hear the alarm. Here we would like to compute the probability of Burglary Alarm.
- Problem: Calculate the probability that alarm has sounded, but there is neither a burglary, nor an earthquake occurred, and David and Sophia both called the Harry.

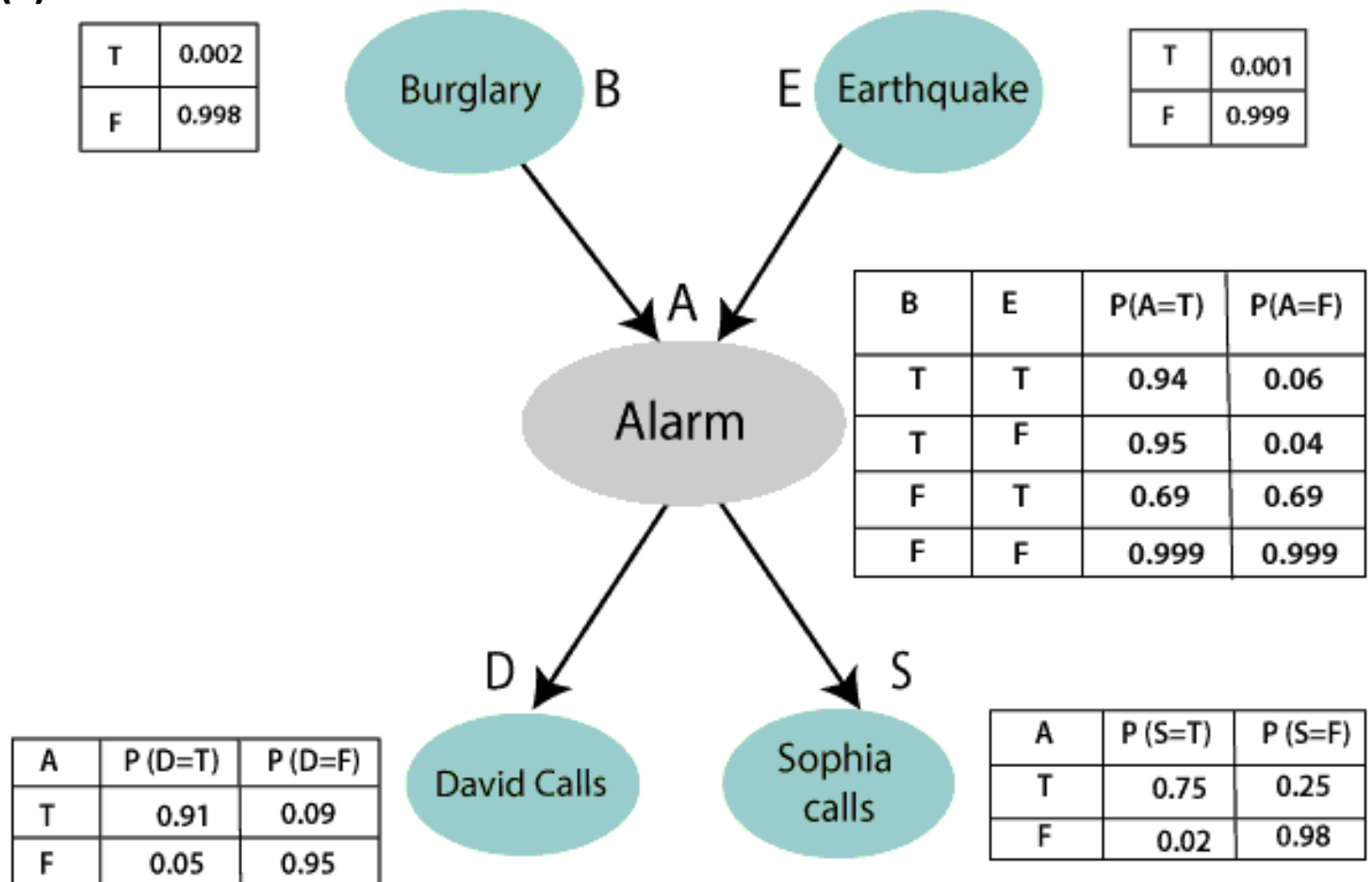


- The Bayesian network for the above problem is given below. The network structure is showing that burglary and earthquake is the parent node of the alarm and directly affecting the probability of alarm's going off, but David and Sophia's calls depend on alarm probability.
- The network is representing that our assumptions do not directly perceive the burglary and also do not notice the minor earthquake, and they also not confer before calling.
- The conditional distributions for each node are given as conditional probabilities table or CPT.
- Each row in the CPT must be sum to 1 because all the entries in the table represent an exhaustive set of cases for the variable.
- In CPT, a boolean variable with k boolean parents contains 2^k probabilities. Hence, if there are two parents, then CPT will contain 4 probability values



List of all events occurring in this network:

- Burglary (B)
- Earthquake(E)
- Alarm(A)
- David Calls(D)
- Sophia calls(S)





$$P(S, D, A, \neg B, \neg E) = P(S|A) * P(D|A) * P(A|\neg B \wedge \neg E) * P(\neg B) * P(\neg E).$$

$$= 0.75 * 0.91 * 0.001 * 0.998 * 0.999$$

$$= \mathbf{0.00068045}.$$

