

Bayesian Methods & Association Rule Mining

$$P(f | TD) = \frac{P(TD | f)P(f)}{P(TD)}$$

Lektor Dr.techn. Alexander K. Seewald
Österreichisches Forschungsinstitut
für Artificial Intelligence

Bayesian Methods & Bayes Theorem

Bayesian Methods provide the basis for probabilistic reasoning:

- Theoretical framework for machine learning, classification, knowledge representation and analysis.
- Allows to integrate uncertain and partial domain knowledge
- Direct modeling of uncertainty
- Easily handles noisy and incomplete data sets with MVs

One cornerstone of Bayesian Methods is Bayes' Rule:

$$P(f | TD) = \frac{P(TD | f)P(f)}{P(TD)}$$

The probability of function/model f given training data TD is equal to the probability of TD given f multiplied by the (prior) probability of f divided by the (prior) probability of TD . All classification methods can be seen as estimating Bayes' Rule, with different techniques to estimate $P(TD|f)$.

Maximum Likelihood

Among all f from Concept Space CS , choose f which has highest probability given training data TD . This is the **maximum a posteriori (MAP)** model:

$$f_{BEST} = \arg \max_{f \in CS} P(f | TD) = \arg \max_{f \in CS} \frac{P(TD | f)P(f)}{P(TD)} = \arg \max_{f \in CS} P(TD | f)P(f)$$

In some cases, we can assume every $f \in CS$ to be equally probable. In that case, the above simplifies to the **maximum likelihood (ML)** model:

$$f_{ML} = \arg \max_{f \in CS} P(TD | f)$$

Example: Cancer diagnosis

Two hypotheses: cancer, \neg cancer ($|CS|=2$)

Diagnostic test for cancer with two outcomes: $\oplus = +$, $\emptyset = -$

Known prob.: $P(\text{cancer}) = 0.008$ $P(\neg\text{cancer}) = 0.992$

(*sensitivity*) $P(\oplus | \text{cancer}) = 0.98$ $P(\emptyset | \text{cancer}) = 0.02$

$P(\oplus | \neg\text{cancer}) = 0.03$ $P(\emptyset | \neg\text{cancer}) = 0.97$ (*specificity*)

$P(\text{cancer} | \oplus) = P(\oplus | \text{cancer})P(\text{cancer}) = 0.98 * 0.008 = 0.0078$ (21%)

$P(\neg\text{cancer} | \oplus) = P(\oplus | \neg\text{cancer})P(\neg\text{cancer}) = 0.03 * 0.992 = 0.0298$ (**79%,MAP**)

Basic Probability Formulas

Product rule : probability of a conjunction of events A and B (= A AND B)

$$P(A \wedge B) = P(A | B)P(B) = P(B | A)P(A)$$

Sum rule : probability of a disjunction of two events A and B (= A OR B)

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

Bayes Theorem : relating posterior and prior probabilities

$$P(f | TD) = \frac{P(TD | f)P(f)}{P(TD)}$$

Theorem of total probability : if events A_1, \dots, A_n are mutually exclusive

$$\text{with } \sum_{i=1}^n P(A_i) = 1 \Rightarrow P(B) = \sum_{i=1}^n P(B | A_i)P(A_i)$$

Bayes Optimal Classifier

So far, we have searched for the best $f \in CS$. However, it is possible to do better by combining all functions/models in CS , weighted by posterior probabilities:

$$Bayes_{optimal} = \arg \max_{Cl_j \in Class} \sum_{f_i \in CS} P(Cl_j | f_i) P(f_i | TD)$$

where $Class$ is the set of possible classes, TD = training data, CS = concept space.

This is the **Bayes Optimal Classifier**.

Advantages

- Optimality: No other classification method with same CS and same prior knowledge can outperform this method (on average). This method maximizes the probability that the new instance is classified correctly, given the available data TD , concept space CS and posterior prob. over all the hypotheses f .
- Predictions correspond to a function/model not in CS – more general than CS !

Disadvantages

- Needs to sum over all possible functions f . Very costly and often intractable.
- Probabilities are usually unknown, and some of them are very hard to estimate.

Gibbs Algorithm & Naïve Bayes

Gibbs Algorithm is a more efficient but less optimal classifier. Under certain conditions it has at most twice the error rate of the optimal Bayes classifier. However, it is still very inefficient.

Gibbs Algorithm

1. Choose a function f from CS at random, according to posterior probability distribution over CS (i.e. $P(f | TD)$)
2. Use f to predict the classification of next instance \mathbf{x} .

A very efficient classifier is obtained by assuming the attributes to be conditionally independent ($P(A_i/A_j)=P(A_i)$ for $\forall i,j$). This is **Naïve Bayes**:

$$Bayes_{naïve} = \arg \max_{Cl_j \in Class} P(Cl_j) \prod_{\forall i} P(a_i | Cl_j)$$

$P(Cl_j)$ and $P(a_i|Cl_j)$ can be efficiently estimated from training data TD by counting. This is a commonly used classifier in machine learning, and works reasonably well even when the conditional independence assumption is violated.

Example: Weather dataset

Classify weather dataset with Naïve Bayes

- Estimate $P(Cl_j)$: $P(yes)=9/14$, $P(no)=5/14$
- Estimate $P(a_i|Cl_j)$, i.e. probability of attribute i having value a_i , given class of Cl_j :

Outlook	Play?	
	yes	no
overcast	4	0
rainy	3	2
sunny	2	3

Windy	Play?	
	yes	no
true	3	3
false	6	2

$$P(\text{outlook}=\text{overcast} \mid \text{yes})=4/9 \quad P(\text{w.}=t \mid \text{yes})=3/9$$

$$P(\text{outlook}=\text{rainy} \mid \text{yes})= 3/9 \quad P(\text{w.}=f \mid \text{yes})=6/9$$

$$P(\text{outlook}=\text{sunny} \mid \text{yes})= 2/9$$

$$P(\text{outlook}=\text{overcast} \mid \text{no}) = 0/5 \quad P(\text{w.}=t \mid \text{no}) =3/5$$

$$P(\text{outlook}=\text{rainy} \mid \text{no}) = 2/5 \quad P(\text{w.}=f \mid \text{no}) =2/5$$

$$P(\text{outlook}=\text{sunny} \mid \text{no}) = 3/5$$

Problem: Estimates may be zero $\Rightarrow P(\text{no} \mid \text{outlook}=\text{overcast})$ would always be 0.

\Rightarrow **Laplace correction:** Use $(a+1)/(b+2)$ instead of a/b , e.g. $1/7$ instead of $0/5$.

Outlook	T	H	Windy	Play?
overcast	64°F	65%	true	yes
overcast	72°F	90%	true	yes
overcast	81°F	75%	false	yes
overcast	83°F	86%	false	yes
rainy	68°F	80%	false	yes
rainy	70°F	96%	false	yes
rainy	75°F	80%	false	yes
rainy	65°F	70%	true	no
rainy	71°F	91%	true	no
sunny	69°F	70%	false	yes
sunny	75°F	70%	true	yes
sunny	72°F	95%	false	no
sunny	80°F	90%	true	no
sunny	85°F	85%	false	no

Example: Weather dataset (2)

- For quantitative / numerical variables, probabilities cannot be determined by counting. Probability density functions must be defined. Common assumption: Values are normally distributed. Then, arithmetic mean and standard deviation define a normal probability density function as follows:

$$P(A_i = x | Cl_j) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

where μ and σ^2 are chosen depending on Cl_j and A_i
 E.g. for $P(\text{temp.}=x | \text{no})$ use $\mu=74.6$ and $\sigma^2=62.3$;
 for $P(\text{hum.}=x | \text{yes})$ use $\mu=79.1$ and $\sigma^2=104.4$ etc..

	<i>Play?</i>	
	<i>yes</i>	<i>no</i>
Temp.		
μ	73.0	74.6
σ^2	38.0	62.3

	<i>Play?</i>	
	<i>yes</i>	<i>no</i>
Hum.		
μ	79.1	86.2
σ^2	104.4	94.7

Outlook	T	H	Windy	Play?
overcast	64°F	65%	true	<i>yes</i>
rainy	68°F	80%	false	<i>yes</i>
sunny	69°F	70%	false	<i>yes</i>
rainy	70°F	96%	false	<i>yes</i>
overcast	72°F	90%	true	<i>yes</i>
sunny	75°F	70%	true	<i>yes</i>
rainy	75°F	80%	false	<i>yes</i>
overcast	81°F	75%	false	<i>yes</i>
overcast	83°F	86%	false	<i>yes</i>
rainy	65°F	70%	true	<i>no</i>
rainy	71°F	91%	true	<i>no</i>
sunny	72°F	95%	false	<i>no</i>
sunny	80°F	90%	true	<i>no</i>
sunny	85°F	85%	false	<i>no</i>

Bayesian Belief Networks

Naïve Bayes: Assumes conditional independence of all attribute. If this is true, then it outputs the optimal Bayes classification. However, in many cases this assumption is overly restrictive.

⇒ **Bayesian Networks:** Allow arbitrary conditional dependence. Dependency information can be learned from training data, or specified as background knowledge. Usually visualized as directed acyclic graph (DAG)

E.g. given a burglary, what is the prob. that John calls?

$$P(J | B) = ?$$

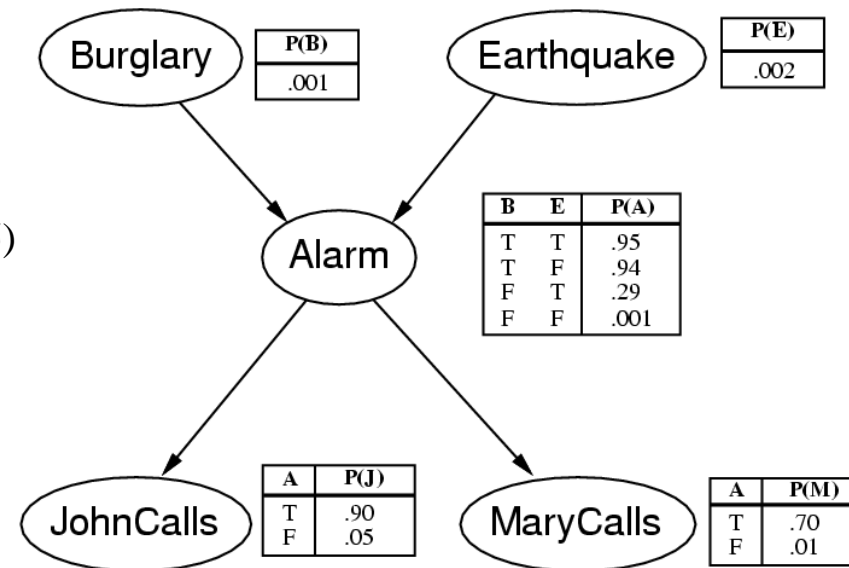
$$P(A | B) = P(B)P(\neg E)(0.94) + P(B)P(E)(0.95)$$

$$P(A | B) = 1(0.998)(0.94) + 1(0.002)(0.95)$$

$$P(A | B) = 0.94$$

$$P(J | B) = P(A | B)(0.9) + P(\neg A | B)(0.05)$$

$$P(J | B) = (0.94)(0.9) + (0.06)(0.05) \\ = 0.85$$



Exact computation of complex queries is NP-hard(!)

Association Rule Mining

Different learning methodology (not prediction)

- Find rules of the form $A \Rightarrow B$ where A and B are itemsets, (i.e. sets of *items*) and $A \cap B = \emptyset$.
- An item is a specific object which may either be part of an itemset or not.
- Rules can related to any attribute; there is no target variable (*unsupervised* as in Reinforcement Learning)
- Usually applied to binary data (*market basket analysis*); Nominal variables can be transformed in the usual way into one binary attribute per value. Each binary attribute corresponds to an item with $\text{attr}_i = 1 \Leftrightarrow \text{Item } i \text{ is present}$.

Association Rule Mining: Defs.

- T is the set of all possible items. When $A, B \subseteq T$, $A \cap B = \emptyset$: $A \Rightarrow B$ is an association rule.

E.g. supermarket: Each distinct product is an item $\in T$. Each transaction Tr_i corresponds to the products bought by one shopper, and thus to a specific itemset. $TD = \{Tr_i\}$.

- Rules of the form $A \Rightarrow B$ tell us about correlations between itemsets, e.g. $(milk, beer) \Rightarrow (diapers, babyfood)$
- Rules are characterized by *Support* (how common is the rule?), and *Confidence* (how well does rule $A \Rightarrow B$ hold?)
- Algorithms to find **all** rules with given minimum *Support* and *Confidence* exist, and are space- and time-efficient.

Support & Confidence

$Support(A \Rightarrow B) = P(A \cup B)$ [$P(A)$ for Borgelt's Implementation]

$$Confidence(A \Rightarrow B) = \frac{P(A \cup B)}{P(A)}$$

$$P(I) = \frac{1}{|TD|} \sum_{Tr \in TD} \prod_{i \in I} (Tr_i = 1) \qquad Lift(A \Rightarrow B) = \frac{P(A \cup B)}{P(A)P(B)}$$

Itemset I is frequent $\Leftrightarrow Support(I) \geq \text{minimum support}$

Adding items: Support can only monotonically decrease ($P(I \cup X) \leq P(I)$), since we add restrictions to the itemset.

\Rightarrow If an itemset I is frequent, all its subsets must also be frequent. If any subset has lower support, then I cannot be frequent. This is an efficient pruning criterion

The Apriori Algorithm

Computes the set of all frequent itemsets.

```
L1 = {set of frequent items (=frequent itemsets of size 1)}  
for (k=1; Lk != ∅; k++) {  
    Ck+1 = {A ∪ B | A, B ∈ Lk, |A ∪ B| = k+1,  
              ∀ X ⊂ (A ∪ B) ⇒ X ∈ (∪ Lk)}  
    Lk+1 = {C | C ∈ Ck+1 && Support(C) > MinSupp }  
}  
return ∪ Lk;
```

Afterwards, compute all possible rules (partitions) of frequent itemsets and output those with min. confidence.

Open-source Apriori implementation by C. Borgelt

<http://fuzzy.cs.uni-magdeburg.de/~borgelt/apriori.html>

Also available in WEKA, Associate Tab.

Rules from Weather Dataset (nominal)

Minimum Support = 0.3, Minimum Confidence = 0.7

- **outlook=overcast \Rightarrow play=yes**
supp: 0.28, conf: 1.0, lift: 1.56
- **temperature=cool \Rightarrow humidity=normal**
supp: 0.28, conf: 1.0, lift: 2.0
- **humidity=normal, windy=false \Rightarrow play=yes**
supp: 0.28, conf: 1.0, lift: 1.56
- **humidity=normal \Rightarrow play=yes**
supp: 0.43, conf: 0.86, lift: 1.33
- **play=no \Rightarrow humidity=high**
supp: 0.28, conf: 0.8, lift: 1.6

...

For large datasets, outputs a large set of rules (>1000), so understanding rules is more challenging than mining.

Ex.: Computing Support and Confidence

- **humidity=normal \Rightarrow play=yes**

confidence=6/7=0.86 (6x correct, 1x wrong)

					Humidity	Windy
Outlook	Temp.				CLASS	
sunny	hot	high	false		Don't Play	
sunny	hot	high	true		Don't Play	
overcast	hot	high	false		Play	
rain	mild	high	false		Play	
rain	cool	normal	false	→	Play	
rain	cool	normal	true	→	Don't Play	
overcast	cool	normal	true	→	Play	
sunny	mild	high	false		Don't Play	
sunny	cool	normal	false	→	Play	
rain	mild	normal	false	→	Play	
sunny	mild	normal	true	→	Play	
overcast	mild	high	true		Play	
overcast	hot	normal	false	→	Play	
rain	mild	high	true		Don't Play	

support=6/14=0.43