
Probabilistic ML algorithm

**Naïve Bayes and Maximum
Likelihood**

There is uncertainty in the world

- Rather than only following the most likely explanation for a given situation, keeping an open mind and considering other possible explanations has proven to be essential in systems that have to work in a real-world environment
- The language of describing uncertainty, that of probability theory, tremendously simplify reasoning in such worlds.
- **Many machine learning problems can be formulated in probabilistic terms.** A probabilistic formulation is the following:
 - the target of a ML system is to learn a classification function $f(x): x \rightarrow y$
 - Given an unseen instance x , assign a category label y to this instance using the learned function $f()$.
 - The corresponding probability formulation is: learn a **probability density function** $p(y/x)$ which is the *conditioned probability of class label y given the observation of x .*
 - A label can be selected based on $\operatorname{argmax}_y p(y/x)$

Example: handwriting recognition



Is this “a” or “9”??

The classifier may return non-zero probabilities for all options

BASIC PROBABILITY NOTIONS YOU NEED

Axioms of Probability Theory

- All probabilities between 0 and 1

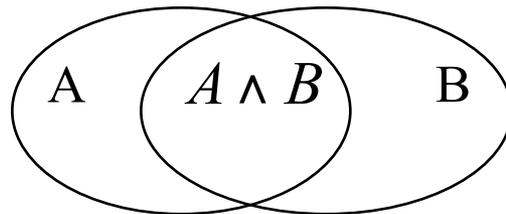
$$0 \leq P(A) \leq 1$$

- True proposition has probability 1, false has probability 0.

$$P(\text{true}) = 1 \quad P(\text{false}) = 0.$$

- The probability of disjunction is:

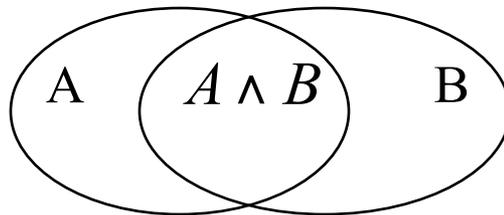
$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$



Conditional Probability

- $P(A | B)$ is the probability of A given B
- Assumes that B is all and only information known.
- Defined by:

$$P(A | B) = \frac{P(A \wedge B)}{P(B)}$$



Statistical Independence

- A and B are *independent* iff:

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

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

These two constraints are logically equivalent

- Therefore, if A and B are independent:

$$P(A | B) = \frac{P(A \wedge B)}{P(B)} = P(A)$$

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

Univariate and Multivariate distributions

- **Univariate** distribution is when there is only one random variable, e.g. if instance vectors \mathbf{x} in \mathbf{X} are described by just one feature, or when there is one classification function $C(\mathbf{x})=Y$
- **Multivariate** if many random variables are involved: e.g. $\mathbf{x}:(x_1 \dots x_d)$. Now any feature of \mathbf{x} can be described by a random variable and we can estimate $P(x_k=v_{jk})$ where v_{jk} $j=1..m_k$ are the possible values for feature k

Probabilistic Classification

- Let Y be the random variable for the class C which takes values $\{y_1, y_2, \dots, y_m\}$ ($|C|=m$ possible classifications for our instances).
- Let X be a random variable describing an instance consisting of a vector of values for d features $\langle X_1, X_2, \dots, X_d \rangle$, let v_{jk} be a possible value for feature X_k .
- For our classification task, we need to compute the (multivariate) conditional probabilities:
$$P(Y=y_i | X=\mathbf{x}(X_1, \dots, X_d)) \text{ for } i=1 \dots m$$

(e.g. $P(Y=\text{positive} / \mathbf{x}=\langle \text{color}=\text{blue}, \text{shape}=\text{circle} \rangle)$)
- E.g. the objective is to classify a new unseen \mathbf{x} by estimating the probability of each possible classification y_i , **given** the observation of feature values of the instance to be classified
- To estimate $P(Y=y_i | X=\mathbf{x})$ we use a learning set D of pairs $(\mathbf{x}_i, C(\mathbf{x}_i))$
- Remember: i is index of class values, k is index of features, j is index of feature values

How can we compute $P(y=y_i/X=x_k)$??

- Example: $x:(\text{color}=\text{red}, \text{shape}=\text{circle})$ $C: y_1=\text{positive}, y_2=\text{negative}$
- We need to compute:

$$P(\text{positive} \mid \text{red} \wedge \text{circle}) = \frac{P(\text{positive} \wedge \text{red} \wedge \text{circle})}{P(\text{red} \wedge \text{circle})}$$

$$P(\text{negative} \mid \text{red} \wedge \text{circle}) = \frac{P(\text{negative} \wedge \text{red} \wedge \text{circle})}{P(\text{red} \wedge \text{circle})}$$

- So we need to compute **joint probabilities**
- The **joint probability distribution** for a set of random INDEPENDENT variables, X_1, \dots, X_d gives the probability of every combination of values (a d -dimensional array with k values if all d variables are discrete with k values, and furthermore $\sum_{k=1..d} P(X_k=v_{jk})=1$)

Joint probability tables

Class=positive

Pr(shape=circle,
color=blue, C=+)

	circle	square
red	0.20	0.02
blue	0.02	0.01

Class=negative

	circle	square
red	0.05	0.30
blue	0.20	0.20

- The probability of all possible conjunctions (assignments of values to some subset of variables) can be estimated from the learning set, by summing the appropriate subset of values from the joint distribution.

$$P(\text{red} \wedge \text{circle}) = P(\text{red} \wedge \text{circle} \wedge \text{positive}) + P(\text{red} \wedge \text{circle} \wedge \text{negative}) = 0.20 + 0.05 = 0.25$$

- If all joint probabilities can be estimated, all conditional probabilities can be calculated.

$$P(\text{positive} | \text{red} \wedge \text{circle}) = \frac{P(\text{positive} \wedge \text{red} \wedge \text{circle})}{P(\text{red} \wedge \text{circle})} = \frac{0.20}{0.25} = 0.80$$

Example

- Consider this dataset:
- X1(red,circle), positive
- X2(red,square),negative
- X3(blue,circle),positive
- X4(red, circle),negative
- X5(red, circle), positive

$$\Pr(\text{positive}/\text{red}\&\text{circle})=\Pr(\text{positive}\&\text{red}\&\text{circle})/\Pr(\text{red}\&\text{circle}) = (2/5)/(3/5)=2/3$$

Probabilistic Classification (2)

- However, given no other assumptions, this requires tables **giving the probability of each category for each possible instance (combination of feature values) in the instance space**, which is **impossible to accurately estimate** from a reasonably-sized training set D .
- E.g. $P(Y=y_i/X_1=v_{1j}, X_2=v_{2j}, \dots, X_d=v_{dj})$ for all y_i and v_{kj}
Assuming that Y and all X_i are binary valued, and we have d features, we need 2^d entries to specify $P(Y=\text{pos} \mid X=x_k)$ for each of the 2^d possible x_k in D since:
 - $P(Y=\text{neg} \mid X=x) = 1 - P(Y=1 \mid X=x)$
 - Compared to $2^{d+1} - 1$ entries for the joint distribution $P(Y, X_1, X_2, \dots, X_d)$

Example: dimension of joint probability tables

- $X:(X_1, X_2 \dots X_4)$, $X_i: \{0, 1\}$ $Y: \{0, 1\}$ (all binary values)
- # of features: $d=4$, # of values for the class Y : $m=2$)
- Consider instance $x_k:(0, 1, 0, 0)$
- Need to estimate $\Pr(Y=0/(X_1=0, X_2=1, X_3=0, X_4=0))$
- If $P(Y=0/(0, 1, 0, 0)) > P(Y=1/(0, 1, 0, 0)) = (1 - P(Y=0/(0, 1, 0, 0)))$
then class is 0, else class is 1
- Overall, **2^4 estimates** are needed for our probabilistic classifier (the number of possible combinations for feature values)
- *For large d and n this is not feasible, even with large learning sets. Simply, we do not have all the necessary evidence!!*

Summary

- We are given an unclassified instance \mathbf{x} , which is represented by d features (boolean, multi-valued or continuous)
- We have a multi-valued classification function $C(\mathbf{x})$ with possible values $y_1 \dots y_m$
- Target of the probabilistic classifier is to classify \mathbf{x} based on: **$y^* = \operatorname{argmax}_{y_i} \Pr(y_i/\mathbf{x})$** (the most likely classification, given the specific combination of feature values in \mathbf{x})
- To do so, we need to estimate $\Pr(y_i/\mathbf{x})$ for all y_i , using evidence provided by previously seen examples $\langle \mathbf{x}_i, y_i \rangle$ in D
- Quite likely, the specific combination of feature values of \mathbf{x} has not seen in the learning set: even if values are boolean, there are 2^d possible combinations! So how do we go about?

Solution: Maximum Likelihood Learning

- Let M be a probabilistic formulation (MODEL) of our classification task $p(y_i/x_k)$
- Suppose that we know exactly the structure of M e.g. we know that:

$$P(y = y_i / x(x_1 = v_1, \dots, x_d = v_d)) = y_i \times (v_1)^{\alpha_1} \times \dots \times (v_d)^{\alpha_d} / Z = f(y_k, v_1, \dots, v_d, \alpha_1 \dots \alpha_d)$$

$$\Theta = \alpha_1 \dots \alpha_d$$

Stated more intuitively: learn the values of model parameters that MAXIMIZE the likelihood of observing the **evidence** provided by the learning set

◆ Goal: After observing several examples $x_1 \dots x_N$, estimate the model parameters, Θ , that generated the observed data.

Maximum Likelihood Estimation (MLE)

- ◆ The **likelihood** of the observed data, given the model parameters Θ , is the **conditional probability** that the model, M , with parameters Θ , produces the observations $x_1, \dots, x_{|D|}$.

$$L(\Theta) = P(x_1 \dots x_{|D|} | \Theta, M) = P(D | \Theta, M)$$

- ◆ In MLE we seek the model parameters, Θ , that **maximize the likelihood of observing our EVIDENCE E , represented by the learning set D .**

We reversed the problem..

- Rather than finding the class value y^* that maximized the conditional probability of any class value given an observed instance \mathbf{x} (the $\operatorname{argmax}_y P(y/\mathbf{x})$) we now try to find the model parameters Θ that maximize the probability of having observed a particular dataset D (=our evidence):

$$\operatorname{argmax}_{\Theta}(P(D/\Theta))$$

- $P(D/\Theta)$ is called **likelihood function** and is also denoted as $L_D(D/\Theta)$

What are these “parameters”?

- The definition of Θ is quite general.
- It can be a set of coefficients in a probabilistic formulation :

$$p(\mathbf{y}/\mathbf{x})=f(\alpha_1, \alpha_2, \dots \alpha_d, \mathbf{x}) \quad \Theta = (\alpha_1, \alpha_2, \dots \alpha_d)$$

- Or a set of probabilities:
- $P(\mathbf{y}/\mathbf{x})=P(X_1=v_1) \wedge P(X_2=v_2) \wedge \dots P(X_d=v_d)$
- $\Theta = (P(X_1=v_1), P(X_2=v_2), \dots P(X_d=v_d))$

Notation: \mathbf{p} is a **probability distribution** and is a general formulation since it applies also to continuous random variables. \mathbf{P} is the **probability** of observing specific **values**.

Statistical Parameter Fitting (general definition for multivariate case)

- Consider instances $D: \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ such that
 - The set of values that \mathbf{x} can take is known
 - Each \mathbf{x}_i is sampled from the same distribution
 - Each \mathbf{x}_i is sampled independently of the rest

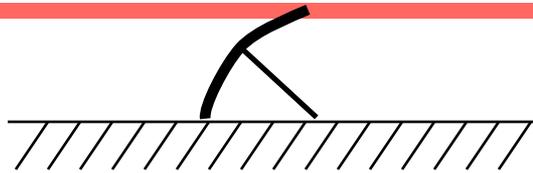
} i.i.d. Samples
- ◆ The task is to find a vector of parameters Θ that have generated the given data D . This vector parameter Θ can be used to predict future data.

Maximum Likelihood Estimation (MLE)

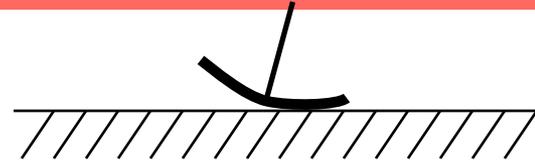
- ◆ In MLE we seek the model parameters, Θ , that **maximize the likelihood**: it is thus an OPTIMIZATION problem (as for in SVM!)
- ◆ The MLE principle is applicable in a wide variety of ML problems, from speech recognition, to natural language processing, computational biology, etc.
- ◆ We will start with the simplest example: Estimating the **bias of a** thumbtack (we got bored with coins!).



Example: Binomial Experiment



Head



Tail

- When tossing the thumbtack, it can land in one of two positions: Head (**H**) or Tail (**T**)
 - ◆ We denote by θ the (unknown) probability $P(\mathbf{H})$.

Estimation task:

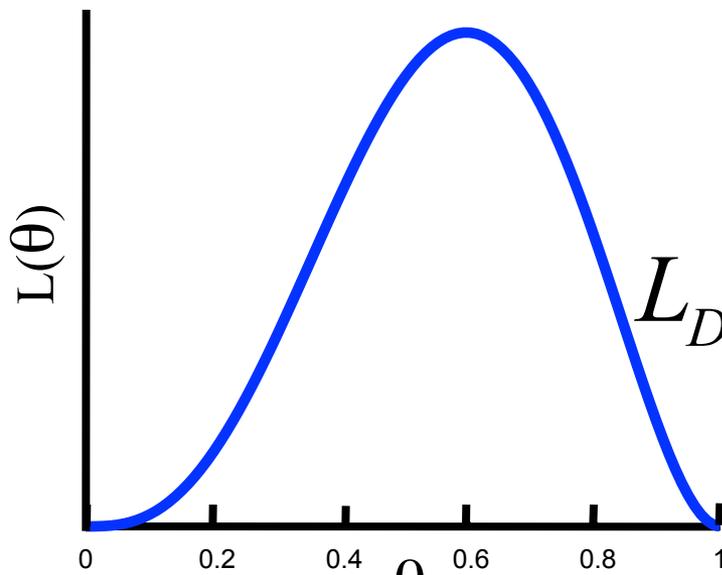
- ◆ Given a sequence of toss samples $x_1..x_m$ (our evidence D) we want to estimate the probabilities $P(\mathbf{H})=\theta$ and $P(\mathbf{T}) = 1 - \theta$
- ◆ θ is the model parameter
- ◆ In this case the problem is univariate (only one stochastic variable)

The Likelihood Function

- How good is a particular θ ? It depends on **how likely it is to generate the observed data**

$$L_D(\theta) = P(D | \theta) = \prod_{j=1..m} P((x_j, y_j) | \theta)$$

- Where are the observed data? In the thumbtack example, we can toss the thumbtack several times and observe a sequence of results
- Ex: the likelihood for the sequence H, T, T, H, H is:



$$L_D(\theta) = \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta$$

Diagram showing arrows pointing from the sequence H, T, T, H, H to the corresponding terms in the likelihood equation: H to θ , T to $(1 - \theta)$, T to $(1 - \theta)$, H to θ , and H to θ .

Sufficient Statistics

- To compute the likelihood in the thumbtack example we only require N_H and N_T (the number of heads and the number of tails in a sequence of tosses)

$$L_D(\theta) = \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta = \theta^{N_H} \cdot (1 - \theta)^{N_T}$$

- N_H and N_T are **sufficient statistics** for the binomial distribution
- A sufficient statistic is a function whose value contains all the information needed to compute any estimate of the parameter

Maximum Likelihood Estimation

MLE Principle:

Choose parameters that maximize the likelihood function

- $\Theta = \alpha_1 \dots \alpha_n$
- $\Theta_{optimal} = \operatorname{argmax}_{\Theta} (P(D/\Theta))$
- MLE is one of the **most commonly used estimators** in statistics
- One usually maximizes the **log-likelihood function**, defined as $l_D(\theta) = \ln L_D(\theta)$

Example: MLE in **Binomial Data** (i.e., Y is boolean so we only need $P(Y=1/X)$ since $P(Y=0/X)=1-P(Y=1/X)$)

$$l_D(\theta) = N_H \log \theta + N_T \log(1 - \theta)$$

Taking derivative and equating it to 0 we get

$$\frac{N_H}{\theta} = \frac{N_T}{1 - \theta} \Rightarrow \hat{\theta} = \frac{N_H}{N_H + N_T}$$

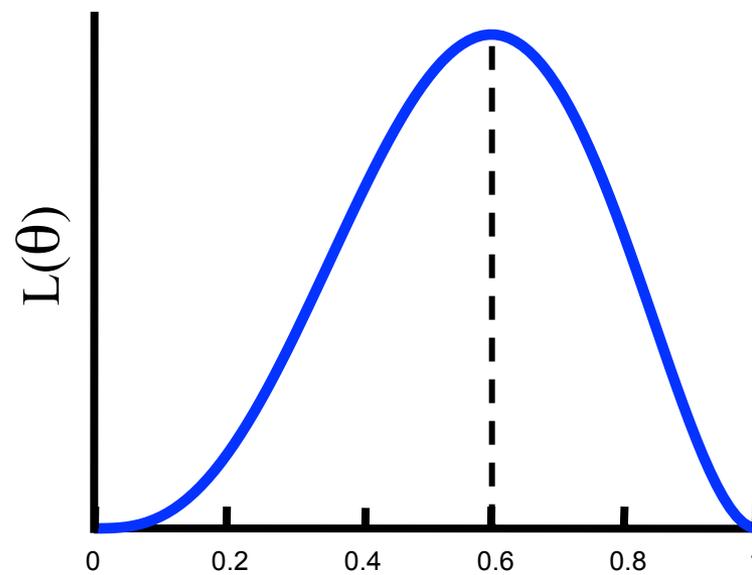
Remember, to maximize minimize a function you need to take the derivative

(which coincides with what one would expect,

Example:

$$(N_H, N_T) = (3, 2)$$

MLE estimate is $3/5 = 0.6$



From Binomial to Multinomial

- Now suppose Y can have the values $1, 2, \dots, K$ (For example a die has $K=6$ sides)
- We want to learn the parameters $\theta_1, \theta_2, \dots, \theta_n$ (the vector Θ of probabilities $\theta_i = P(Y=y_i)$)

Sufficient statistics:

◆ N_1, N_2, \dots, N_K - the number of times each outcome is observed

◆ The optimization problem is: maximize the log

of:
$$L_D(\theta) = \prod_{i=1}^K \theta_i^{N_k} \text{ such that: } \sum_i \theta_i = 1 \text{ and } \theta_i \geq 0 \quad \forall i$$

Maximizing log-likelihood with constraints

We consider the log-likelihood and we represent the optimization problem with a Lagrangian (again!! Remember SVM)

$$l(\alpha, \Theta) = \sum N^i \log \theta_i - \alpha \left(\sum \theta_i - 1 \right)$$

$$\frac{dl(\alpha, \Theta)}{d\theta_i} = 0$$

$$\frac{\partial(i \log(x))}{\partial x} = \frac{i}{x}$$

$$\frac{N^i}{\theta_i} - \alpha = 0 \Rightarrow \theta_i = \frac{N^i}{\alpha}$$

$$\sum \frac{N^i}{\alpha} = 1 \Rightarrow \alpha = \sum N^i$$

$$\hat{\theta}_i = \frac{N^i}{\sum_j N^j}$$

Another example: proteine sequences

- Let $x_1x_2\dots x_n$ be a protein sequence
- We want to learn the parameters $\theta_1, \theta_2, \dots, \theta_{20}$ corresponding to the probabilities of the 20 amino acids
- N_1, N_2, \dots, N_{20} - the number of times each amino acid is observed in the sequence

$$L_D(q) = \prod_{i=1}^{20} \theta_i^{N_i}$$

Likelihood function:

MLE: $\theta_i = \frac{N_i}{n} \quad n = \sum_{i=1}^{20} N_i$

NAIVE BAYES CLASSIFIER: A ML ALGORITHM BASED ON MLE PRINCIPLE

Naive bayes

- We are given an evidence represented by an annotated learning set D of classified instances (\mathbf{x}_i, y_i) .
- We would like to estimate the probabilities $p(Y=y_i/\mathbf{x})$ for **unobserved instances**, given the evidence D
- In probability calculus, often estimating $P(a/b)$ is easier than estimating $P(b/a)$ so we need the Bayes theorem to invert conditional probabilities

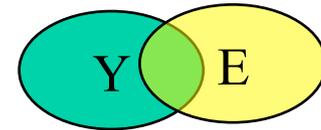
Bayes Theorem

- Y=classification
- E= evidence of data

$$P(Y | E) = \frac{P(E | Y)P(Y)}{P(E)}$$

Simple proof from definition of conditional probability:

$$P(Y | E) = \frac{P(Y \wedge E)}{P(E)} \quad (\text{Def. cond. prob.})$$



$$P(E | Y) = \frac{P(Y \wedge E)}{P(Y)} \quad (\text{Def. cond. prob.})$$

$$P(Y \wedge E) = P(E | Y)P(Y)$$

QED:
$$P(Y | E) = \frac{P(E | Y)P(Y)}{P(E)}$$

Naïve Bayes Model (1)

For each classification value y_i we have (applying Bayes):

$$P(Y = y_i | X = \mathbf{x}) = \frac{P(Y = y_i)P(X = \mathbf{x} | Y = y_i)}{P(X = \mathbf{x})}$$

- $P(Y=y_i)$ and $P(X=\mathbf{x})$ are called **priors** and can be estimated from learning set D since categories are **complete** and **disjoint**

Naive Bayes Model (2)

$$P(Y = y_i | X = \mathbf{x}) = \frac{P(Y = y_i)P(X_1 = v_1, X_2 = v_2, \dots, X_d = v_d | Y = y_i)}{P(X = \mathbf{x})} =$$

$$P(Y = y_i) \frac{\prod_{j=1}^d P(X_j = v_{jk} | Y = y_i)}{P(X = \mathbf{x})}$$

We assume feature values v_{jk} of different features X_j being **statistically independent**. v_{jk} is the k -th value of feature j where $j=1,2,\dots,d$ and $k=1,\dots,K_j$ (if binary features, $k=0$ or 1)

e.g. $P(\mathbf{x}(\text{color}=\text{blue}, \text{shape}=\text{circle}, \text{dimension}=\text{big})) = P(\text{color}=\text{blue})P(\text{shape}=\text{circle})P(\text{dimension}=\text{big})$

and furthermore

$$\sum P(X_j = v_{jk} / Y = y_i) = 1$$

Estimating model parameters Θ

- The parameters are:
- θ^1 : $P(Y=y_i)$ for all i , and
- θ^2 : $P(X_j = v_{jk} | Y = y_i)$ for all i, j and k
- The log-likelihood function is:

$$L(\Theta) = \sum N^i \log \theta_i^1 + \sum \prod N^{jki} \log \theta_{jki}^2$$

- Where N_i is the number of times we see class y_i in dataset, and N^{jki} is the number of times we see an instance with $X_j = v_{jk}$ and $Y=y_i$

Remember: i index of class, j index of features, k index of feature values

Estimating model parameters (2)

The MLE problem is therefore: **maximize** $L(\theta)$ subject to:

$$\sum \theta_i^1 = 1 \quad \theta_i^1 \geq 0; \quad \sum_k \theta_{jki}^2 = 1 \quad \theta_{jki}^2 \geq 0$$

$$L(\Theta) = \sum N^i \log \theta_i^1 + \sum \prod N^{jki} \log \theta_{jki}^2 - \alpha \sum_{i=1..|C|} (\theta_i^1 - 1) - \beta \sum_{k=1..K} (\theta_{jki}^2 - 1)$$

Since we have logarithms:

$$\sum \prod N^{jki} \log \theta_{jki}^2 = \sum \sum N^{jki} \log \theta_{jki}^2$$

Estimating model parameters (3)

To maximize the likelihood we need to find values that take to zero the derivative

$$\frac{\partial L(\alpha, \beta, \Theta)}{\partial \theta_{jki}^2} = \frac{N^{jki}}{\theta_{jki}^2} - \beta = 0 \rightarrow \theta_{jki}^2 = \frac{N^{jki}}{\beta}$$

and since $\sum_{k=1..K} \theta_{jki}^2 = 1$

we obtain $\theta_{jki}^2 = \frac{N^{jki}}{\sum_k \theta_{jki}^2} = \frac{N^{jki}}{N^i}$

I.e. the θ^2 can be estimated as the ratio between the number of times feature X_j takes value v_{jk} when $Y=y_i$ and the total number of examples in D for which $Y=y_i$

How do we predict the category of an unseen instance with Naive Bayes?

$$P(Y = y_i | X = x_k) = P(Y = y_i) \prod_{j=1..d} P(v_{jk} | Y = y_i) / \prod_{j=1..d} P(v_{jk}) = \theta_i^1 \prod_{j=1..d} \theta_{ikj}^2 / \prod_{j=1..d} P(v_{jk})$$

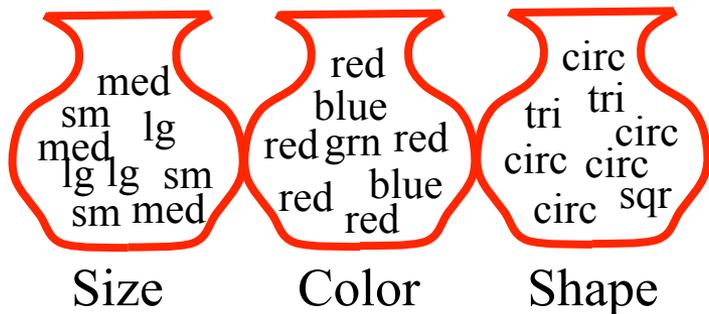
$$y_i = \underset{k}{\operatorname{argmax}} \left(\theta_i^1 \prod_{j=1..d} \theta_{ikj}^2 \right)$$

Note that since the denominator is common to all conditional probabilities, it does not affect the ranking.

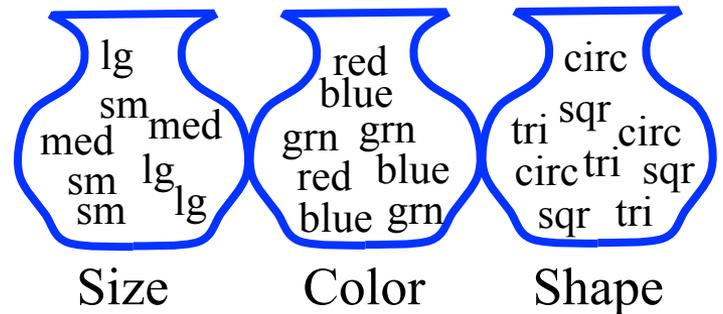
No need to compute it!

Naïve Bayes Generative Model Example

$K=3, |C|=2, d=3$



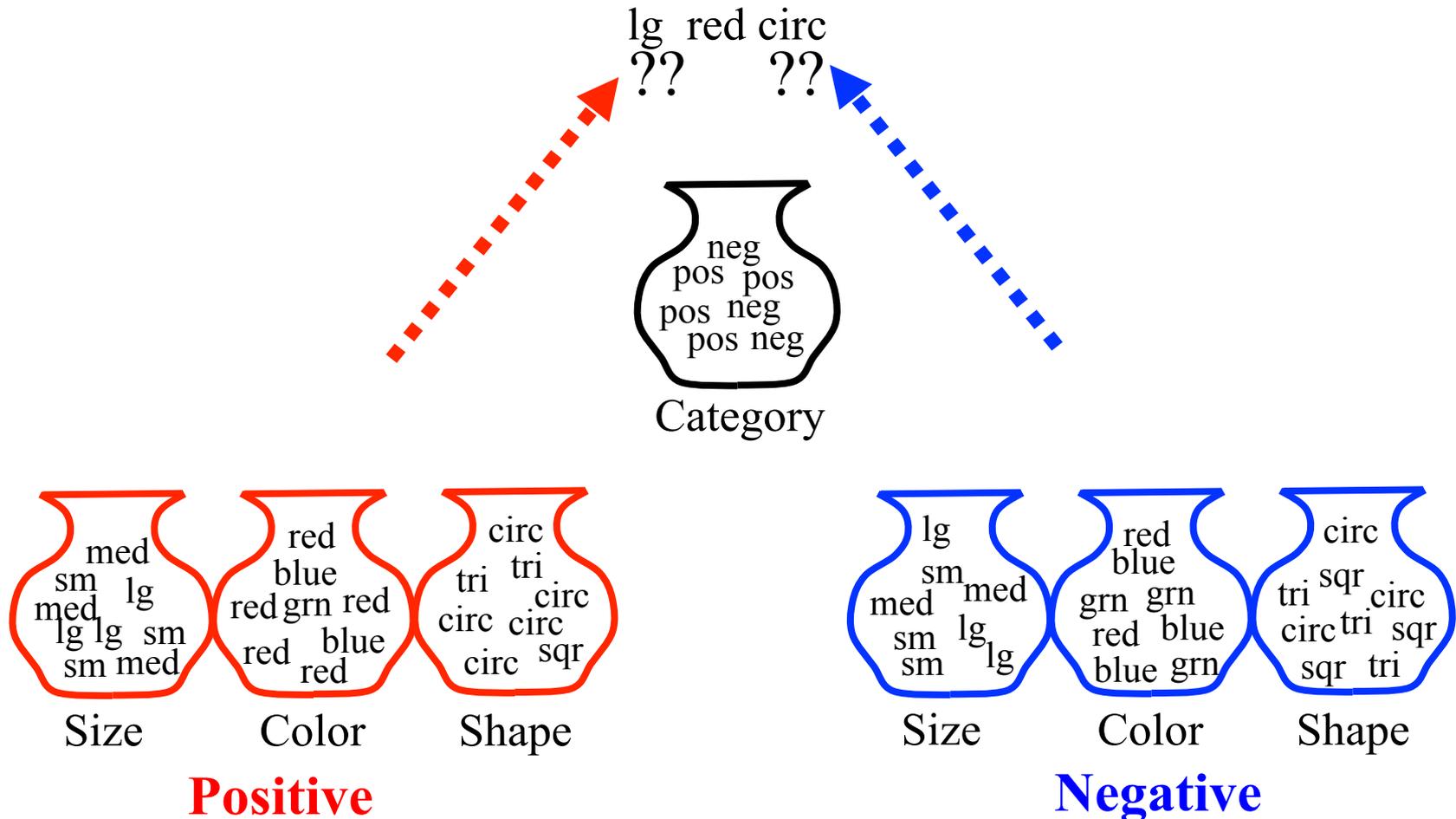
Positive



Negative

Naïve Bayes Inference Problem

I estimate on the learning set the probability of extracting **lg**, **red**, **circ** from the red or blue urns. Whichone is higher?



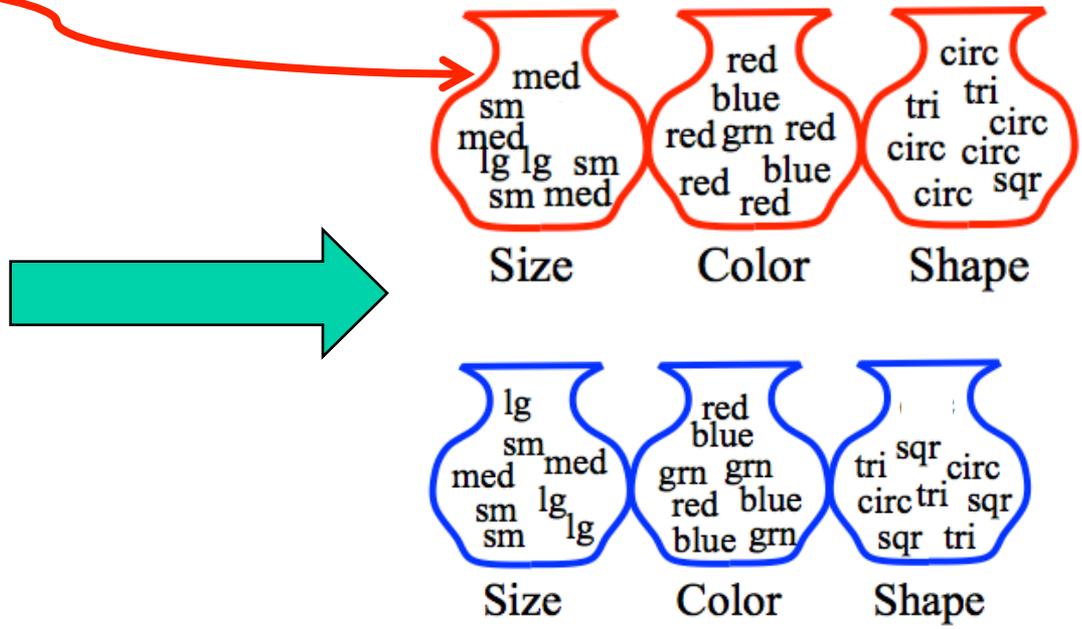
HOW?

We fill urns using the available dataset

D

1	med	red	circ	pos
2	sm	blue	tri	pos
3	med	red	tri	pos
4	lg	grn	circ	pos
5	lg	red	circ	pos
6	sm	blue	circ	pos
7	sm	red	sqr	pos
8	med	red	circ	pos
9	lg	red	sqr	neg
10	sm	blue	tri	neg
11	med	grn	circ	neg
12	med	grn	tri	neg
13	lg	red	circ	neg
14	sm	blue	sqr	neg
15	sm	blue	tri	neg
16	lg	grn	sqr	neg

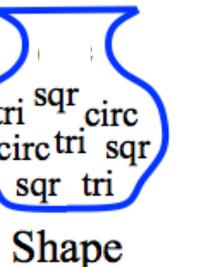
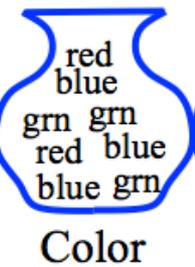
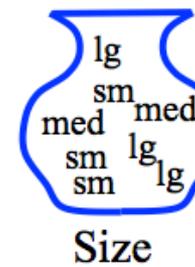
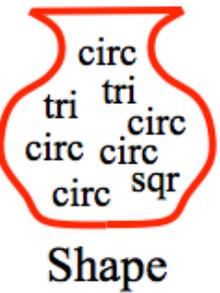
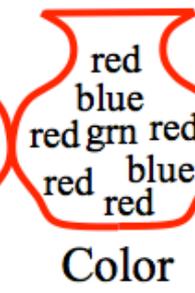
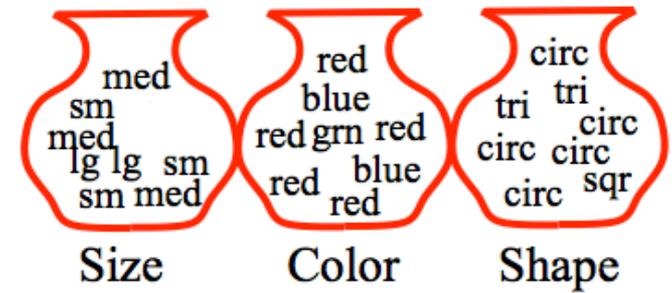
Note that because of independence hypothesis, the specific combination of feature values does not matter!



Now we can compute parameters

$$\theta_{\text{size,small,positive}} = P(\text{size} = \text{small} / C = \text{positive})$$

Probability	Y=positive	Y=negative
P(Y)	0.5	0.5
P(small Y)	3/8	3/8
P(medium Y)	3/8	2/8
P(large Y)	2/8	3/8
P(red Y)	5/8	2/8
P(blue Y)	2/8	3/8
P(green Y)	1/8	3/8
P(square Y)	1/8	3/8
P(triangle Y)	2/8	3/8
P(circle Y)	5/8	2/8



Have 3 small out of 8 instances in red “size” urn
then $P(\text{size}=\text{small}/\text{pos})=3/8=0,375$ (round 4)

Training set D (evidence)

Using parameters, we can classify unseen instances

Probability	Y=positive	Y=negative
P(Y)	0.5	0.5
P(medium Y)	3/8	2/8
P(red Y)	5/8	2/8
P(circle Y)	5/8	2/8

$$y_i = \underset{k}{\operatorname{argmax}}(\theta_i^1 \prod_{j=1..d} \theta_{ikj}^2)$$

Test Instance:
X.<medium ,red, circle>

$$P(\text{positive} | X) = P(\text{positive}) * P(\text{medium} | \text{positive}) * P(\text{red} | \text{positive}) * P(\text{circle} | \text{positive}) / P(X)$$

$$0.5 * 3/8 * 5/8 * 5/8 = 0,073$$

$$P(\text{negative} | X) = P(\text{negative}) * P(\text{medium} | \text{negative}) * P(\text{red} | \text{negative}) * P(\text{circle} | \text{negative}) / P(X)$$

$$0.5 * 2/8 * 2/8 * 2/8 = 0.0078$$

P(positive/X) > P(negative/X) → positive

Note: sum is not 1 since we ignore the denominator of original formulation $\prod P(v_{jk})$

Naive summary

Classify any new datum instance $\mathbf{x}_k = (x_1, \dots, x_n)$ as:

$$y_{Naive\ Bayes} = \operatorname{argmax}_i P(y_i) P(\mathbf{x} | y_i) \rightarrow \operatorname{argmax}_i P(y_i) \prod_{j=1..d} P(v_{jk} | y_i)$$

- To do this based on training examples, estimate the parameters from the training examples in D :

- For each target value of the classification variable (hypothesis) y_i

$$\hat{P}(Y = y_j) := \mathbf{estimate} P(y_i)$$

- For each attribute value v_{jk} of each datum instance

$$\hat{P}(x_j = v_{jk} | Y = y_i) := \mathbf{estimate} P(v_{jk} | y_i)$$

Estimating Probabilities

- Normally, as in previous example, probabilities are estimated based on observed frequencies in the training data.
- If D contains N_i examples in category y_i , and N_{jki} of these N_i examples have the k -th value for feature X_j , v_{jk} , then:

$$P(X_j = v_{jk} | Y = y_i) = \frac{N_{jki}}{N_i}$$

- However, estimating such probabilities from small training sets is **error-prone**.
- If -due only to chance- a rare feature, $X_j=v_{jk}$ is never observed in the training data, then $P(X_j=v_{jk} | Y=y_i) = 0$.
- If $X_j=v_{jk}$ then occurs in a test example, X , the result is that $\forall y_k: P(X | Y=y_i) = 0$ and $\forall y_i: P(Y=y_i | X) = 0$ (since individual probability estimates are multiplied)

Probability Estimation Example

Ex	Size	Color	Shape	Category
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative

Test Instance X :
 <medium, red, circle>

Probability	positive	negative
$P(Y)$	0.5	0.5
$P(\text{small} Y)$	0.5	0.5
$P(\text{medium} Y)$	0.0	0.0
$P(\text{large} Y)$	0.5	0.5
$P(\text{red} Y)$	1.0	0.5
$P(\text{blue} Y)$	0.0	0.5
$P(\text{green} Y)$	0.0	0.0
$P(\text{square} Y)$	0.0	0.0
$P(\text{triangle} Y)$	0.0	0.5
$P(\text{circle} Y)$	1.0	0.5

$$P(\text{positive} | X) = 0.5 * 0.0 * 1.0 * 1.0 / P(X) = 0$$

$$P(\text{negative} | X) = 0.5 * 0.0 * 0.5 * 0.5 / P(X) = 0$$

Smoothing

- To account for estimation from small samples, probability estimates are adjusted or *smoothed*.
- Laplace smoothing using an m -estimate assumes that each feature **is given a prior probability, p** , that is assumed to have been previously observed in a “virtual” sample of size m .

$$P(X_j = v_{jk} | Y = y_i) = \frac{N_{jki} + mp}{N_i + m}$$

- For binary features, p is simply assumed to be 0.5, while it can be set to $1/k$ for k -valued features.

Laplace Smoothing Example

- Assume training set contains 10 positive examples, and feature size has 3 values, but 1 (medium) is not observed in D:
 - 4: small
 - 0: medium
 - 6: large
- Estimate parameters as follows (if $m=1, p=1/3$)
 - $P(\text{small} \mid \text{positive}) = (4 + 1/3) / (10 + 1) = \underline{0.394}$
 - $P(\text{medium} \mid \text{positive}) = (0 + 1/3) / (10 + 1) = 0.03$
 - $P(\text{large} \mid \text{positive}) = (6 + 1/3) / (10 + 1) = 0.576$
 - $P(\text{small or medium or large} \mid \text{positive}) = 1.0$

Continuous Attributes

- If X_j is a **continuous** feature rather than a discrete one, need another way to calculate $P(X_j | Y)$.
- Assume that X_j has a **Gaussian** distribution whose mean and variance depends on Y .
- During training, for each combination of a continuous feature X_j and a class value for Y , y_i , estimate a **mean**, μ_{ji} , and standard deviation σ_{ji} based on the values of feature X_j in class y_i in the training data. μ_{ji} is the mean value of X_j observed in instances for which $Y = y_i$ in D
- **During testing**, estimate $P(X_j | Y = y_i)$ for a given example, using the Gaussian distribution defined by μ_{ji} and σ_{ji} .

$$P(X_j = v_{jk} | Y = y_i) = \frac{1}{\sigma_{ji} \sqrt{2\pi}} \exp\left(\frac{-(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

Comments on Naïve Bayes

- Tends to work well despite strong assumption of conditional independence.
- Experiments show it to be quite competitive with other classification methods on standard UCI datasets.
- Although it does not produce accurate probability estimates when its independence assumptions are violated, it may still pick the correct maximum-probability class in many cases.
 - Able to learn conjunctive concepts in any case
- Does not perform any search of the hypothesis space. Directly constructs a hypothesis from parameter estimates that are easily calculated from the training data.
 - Strong bias
- Not guaranteed consistency with training data.
- Typically handles noise well since it does not even focus on completely fitting the training data.