



# BAYESIAN NETWORK

Bor-shen Lin

[bslin@cs.ntust.edu.tw](mailto:bslin@cs.ntust.edu.tw)

<http://www.cs.ntust.edu.tw/~bslin>

# OUTLINES

- Random Variables and Distribution
- Expectation
- Joint distribution and conditional probability
- Bayesian Theory
- Bayesian Classifier
- Bayesian Network



# TYPE OF RANDOM VARIABLES

- Discrete variables
  - Coin/dice tossing
  - Lottery
  - Blood type
- Continuous variable
  - Temperature
  - Time of leaving home
  - Sound
  - Image

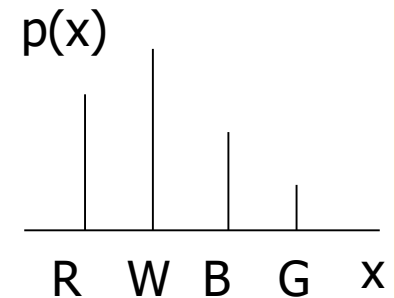


# DISCRETE RANDOM VARIABLES

- A basket: 3 red(R) balls, 4 white(W) balls, 2 blue(B) balls, 1 green(G) ball.
- Random variable X: the color of a ball drawn randomly from the bag
- *Probability Weight Function* (p.w.f.)

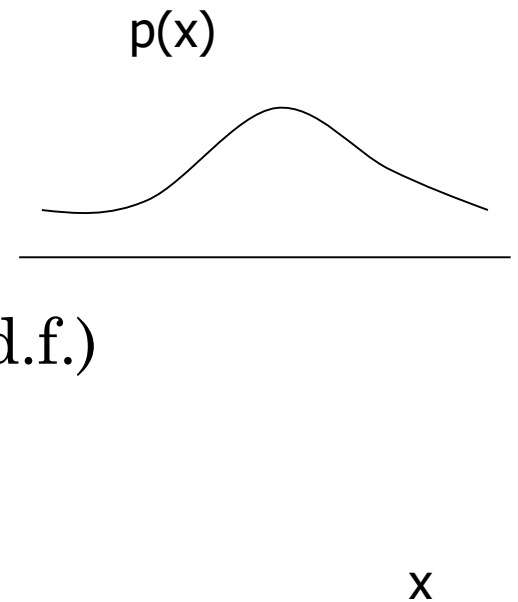
$$p(X=x) = p(x)$$

- $p(X = R) = 3/10$ ,  $p(X = W) = 4/10$ ,  
 $p(X = B) = 2/10$ ,  $p(X = G) = 1/10$
- $\sum_x p(x) = 1.0$



# CONTINUOUS RANDOM VARIABLES

- Example: temperature  $X$
- *Probability Density Function* (p.d.f.)
  - $p(X = x) = p(x)$
  - $\int p(x)dx = 1$
- *Cumulated distribution function* (c.d.f.)
$$P(x) = \int_{-\infty}^x p(x)dx$$
- Example: height 、 weight 、 score



# EXPECTATION

## ○ Continuous Random Variable

- $\mu = E(X) \equiv \int_{-\infty}^{\infty} xp(x)dx$
- $var(X) = E((X - \mu)^2) \equiv \int_{-\infty}^{\infty} (x - \mu)^2 p(x)dx$
- $\sigma \equiv \sqrt{var(X)}$

## ○ Discrete Random Variable

- $\mu = E(X) \equiv \sum_{-\infty}^{\infty} xp(x)$
- $var(X) = E((X - \mu)^2) \equiv \sum_{-\infty}^{\infty} (x - \mu)^2 p(x)$
- $\sigma \equiv \sqrt{var(X)}$

## ○ Numerical only

- might not be able to compute the expectations
- Examples: blood type 、 job



# EXPECTATION

- For any  $f(x)$ , we may take the expectation

- $E(f(X)) \equiv \int_{-\infty}^{\infty} f(x)p(x)dx$

- Entropy

- Average information

- $$I(X) \equiv E\left(\log\left(\frac{1}{p(X)}\right)\right) = \int_{-\infty}^{\infty} \log\left(\frac{1}{p(x)}\right)p(x)dx$$
$$= -\int_{-\infty}^{\infty} \log(p(x))p(x)dx.$$



# JOINT DISTRIBUTION

## ○ Joint distribution : $p(x,y)$

- $\iint p(x, y) dx dy = 1.0$
- $p(x) \equiv \int p(x, y) dy \rightarrow \int p(x) dx = 1$
- $p(y) \equiv \int p(x, y) dx \rightarrow \int p(y) dy = 1$

## ○ Conditional probability.

- $p(x | y) \equiv p(x, y) / p(y)$ 
  - *p. d. f. of  $X$  when  $y$  is given*
  - $\int p(x | y) dx = \int [p(x, y) / p(y)] dx = 1$
- $p(x) \equiv \int p(x, y) dy = \int p(x | y) p(y) dy$
- $p(y) \equiv \int p(x, y) dx$





# CONDITIONAL MEAN

- $E(X | y) = \int x \cdot p(x | y) dx \neq E(X)$ 
  - Expectation of  $X$  given  $Y = y$
  - $E(X | y)$  is a function of  $y$  (varies w.r.t.  $y$ )
- $E(X) = \int x p(x) dx = \iint x p(x, y) dx dy$ 
$$= \iint [x p(x | y) dx] p(y) dy$$
$$= \int E(X | y) p(y) dy$$
- $E(X)$  is not the function of  $y$  (the weighting average of  $E(X | y)$  on all  $y$ 's)



# STATISTICAL INDEPENDENCE

- $p(x, y) = p(x) p(y)$   
→ X, Y are statistically independent

- $p(x|y) = \frac{p(x,y)}{p(y)} = p(x)$

the probability weight of  $x$  is NOT influenced by  $y$

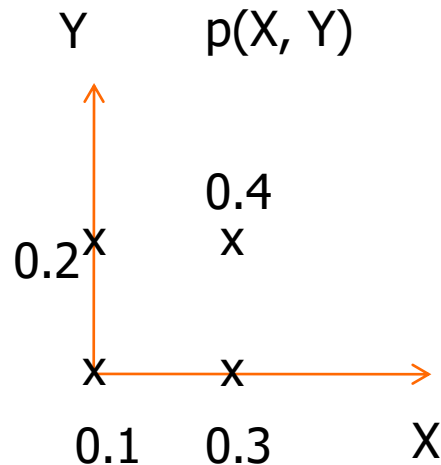
- Example

- Coin tossing result X : Head/Tail
- Dice tossing result Y: point 1,2,..., 6
- When X and Y are *statistically independent*

$$\begin{aligned} p(X = \text{Head}, Y = 4) &= p(X = \text{Head})p(Y = 4) \\ &= (1/2) * (1/6) \end{aligned}$$



# EXAMPLE OF STATISTICALLY DEPENDENT



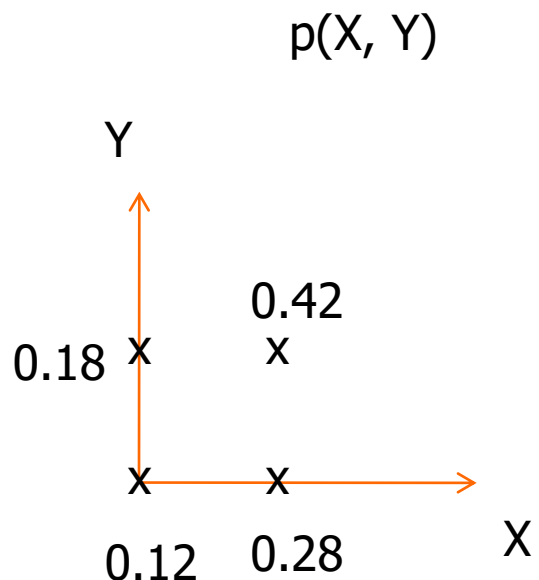
X, Y are statistically  
dependent or  
independent?

$$\begin{aligned}p(X = 0, Y = 0) &= 0.1 \\p(X = 0, Y = 1) &= 0.2 \\p(X = 1, Y = 0) &= 0.3 \\p(X = 1, Y = 1) &= 0.4\end{aligned}$$

1.  $p(X=0)=0.1+0.2 = 0.3$   
 $p(X=1) = 0.3+0.4=0.7$
2.  $p(X=0 \mid Y = 0)=0.1/0.4=0.25$   
 $p(X=1 \mid Y = 0)=0.3/0.4=0.75$
3.  $p(X=0 \mid Y = 1)=0.2/0.6=0.333$



# EXAMPLE OF STATISTIC INDEPENDENCE



$X, Y$  are statistically independent

$$p(X = 1) = 0.42 + 0.28 = 0.7$$

$$p(X = 0) = 0.18 + 0.12 = 0.3$$

$$p(X = 1|Y = 1) = \frac{p(X=1, Y=1)}{p(Y=1)} = \frac{0.42}{0.18+0.42} = 0.7$$

$$p(X = 0|Y = 1) = \frac{0.18}{0.18 + 0.42} = 0.3$$

# BAYESIAN THEORY

- $p(y|x) = \frac{p(x,y)}{p(x)} = \frac{p(y)p(x|y)}{\sum_y p(y)p(x|y)}$
- Could be used for inference between observable phenomena( $x$ ) and unseen cause ( $y$ )
  - X: phenomena/events (e.g.: having fever, cough,...)
    - Could have multiple dimensions!
    - X could be continuous/discrete/composite!
  - Y: cause(e.g. having cold, inflammation, plague, ...)
    - Classification/detection: Y is discrete variable
    - Estimation: Y is continuous variable



# CORRELATION COEFFICIENT

$$\rho_{xy} \equiv \frac{E[(X - E(X))(Y - E(Y))]}{[\text{Var}(X) \cdot \text{Var}(Y)]^{\frac{1}{2}}}$$

- If X and Y are statistically independent  
 $\rightarrow \rho_{xy} = 0$
- $\rho_{xy} = 0$  does not guarantee statistical independence



# CONDITIONAL PROBABILITY FOR MEDICAL DIAGNOSIS

- F: fever (**symptom**), C: having a cold (**cause**)
  - $F \rightarrow C$  or  $C \rightarrow F$  does not hold in general
- $p(F | C) = 0.8$  having fever in case of having cold  
 $p(F) = 0.001$  having fever  
 $p(C) = 0.0001$  having cold  
 $p(C | F) = p(C)p(F | C)/p(F) = 0.008$
- Though it is highly probable that one who has cold also have fever (0.8), it can be concluded that one must have fever because he/she caught cold (only 0.008)
- $p(C) = 0.0001 \rightarrow p(C | F) = 0.008$ 

The event F increases the probability of having cold (from 0.0001 to 0.008), but the probability is not high.



# INFERENCE WITH CONDITIONAL PROBABILITY

- $p(F) = 0.001$  (F: one has fever)  
 $p(C) = 0.0001$  (C: have a cold)  
 $p(P) = 0.0000000001$  (P: plague)  
 $p(F | C) = 0.8$  (having fever when having a cold)  
 $p(F | P) = 0.99$  (have fever in case getting plague)  
 $p(F | P) > p(F | C)$  ? (ML detection)
- With F (one has fever)  $\rightarrow$  guess C or P? (disease)  
$$\frac{p(C | F)}{p(P | F)} = \frac{p(C)p(F | C)}{p(P)p(F | P)}$$
$$= (0.0001 * 0.8) / (0.0000000001 * 0.99) = 84210$$
$$\rightarrow p(C | F) \gg p(P | F) \text{ (C is more likely)}$$
$$\rightarrow \text{it is more reasonable to guess C!}$$
- Without F:  $p(C)/p(P) = 100000$





# SIMPLE BAYESIAN LEARNING (1)

- Suppose we have a set of hypotheses  $H_1 \dots H_n$ .
- For each  $H_i$ ,  $p(H_i|E) = \frac{p(H_i)p(E|H_i)}{p(E)}$ 
  - $p(H_i|E)$  is used to represent the probability that some hypothesis,  $H$ , is true, given evidence  $E$ .
  - Hence, given a piece of evidence, a learner can determine which is the most likely explanation by finding the hypothesis that has the highest **posterior probability**.
- Maximum a Posterior (MAP) detection
$$i^* = \operatorname{argmax}_i \{p(H_i|E)\}$$



# SIMPLE BAYESIAN CONCEPT LEARNING (2)

- Since  $P(E)$  is independent of  $H_i$  it will have the same value for each hypothesis.
- Hence, it can be ignored, and we can find the hypothesis with the highest value of:

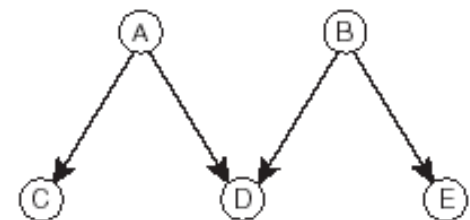
$$P(E|H_i) \cdot P(H_i)$$

- We can simplify this further if all the hypotheses are **equally likely**, in which case we simply seek the hypothesis with the highest value of  $P(E|H_i)$ .
- This is the likelihood of  $E$  given  $H_i$ .
- Maximum Likelihood (ML) detection



# BAYESIAN BELIEF NETWORKS (1)

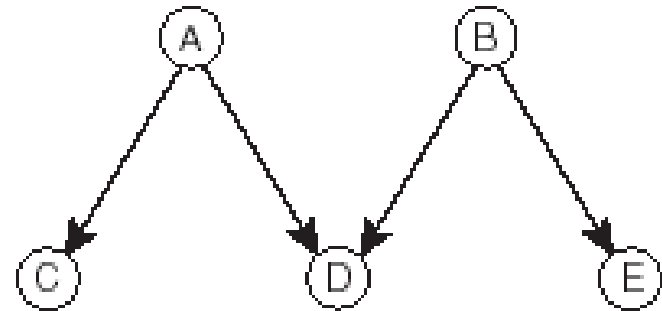
- A belief network assumes some **statistical dependencies** between a set of variables.
  - Such assumptions can simplify the computations of complicated conditional probabilities
- Two variables A and B are statistically independent if the likelihood that A will occur has nothing to do with whether B occurs.
- Example : C and D are dependent on A; D and E are dependent on B.
- The Bayesian belief network has the conditional probabilities associated with each link.



# BAYESIAN BELIEF NETWORKS (2)

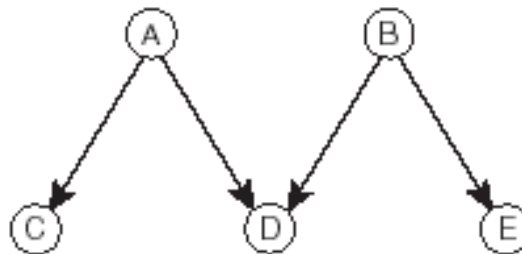
- A complete set of probabilities for this BBN

- $P(A) = 0.1$
- $P(B) = 0.7$
- $P(C|A) = 0.2$
- $P(C|\neg A) = 0.4$
- $P(D|A \wedge B) = 0.5$
- $P(D|A \wedge \neg B) = 0.4$
- $P(D|\neg A \wedge B) = 0.2$
- $P(D|\neg A \wedge \neg B) = 0.0001$
- $P(E|B) = 0.2$
- $P(E|\neg B) = 0.1$



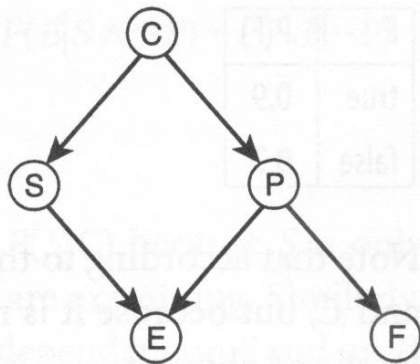
# BAYESIAN BELIEF NETWORKS (3)

- The joint probability
  - $P(A, B, C, D, E) = P(A) \cdot P(B|A) \cdot P(C|A, B) \cdot P(D|A, B, C) \cdot P(E|A, B, C, D)$
- With the assumption of independences, the computation of  $P(A, B, C, D, E)$  can be **simplified**.
  - $P(C|A, B) = P(C|A)$ ,
  - $P(D|A, B, C) = P(D|A, B)$
  - $P(E|A, B, C, D) = P(E|B)$
  - $P(A, B, C, D, E) = P(A) \cdot P(B) \cdot P(C|A) \cdot P(D|A, B) \cdot P(E|B)$



# EXAMPLE OF BN

- C: you go to college
- S: you will study
- P: you will party
- E: you will be successful in exams
- F: you will have fun
- $P(E | F, \neg P, S, C)$  ?



P	P(F)
true	0.9
false	0.7

$p(F|P)$

P(C)
0.2

$p(C)$

C	P(S)
true	0.8
false	0.2

$p(S|C)$

C	P(P)
true	0.6
false	0.5

$p(P|C)$

S	P	P(E)
true	true	0.6
true	false	0.9
false	true	0.1
false	false	0.2

$p(E|S,P)$

# NAÏVE BAYES CLASSIFIER (1)

- A vector of data is classified.
  - $P(c_i|x_1, x_2, \dots, x_n)$
  - The classification with the highest posterior probability is chosen.
  - The hypothesis which has the highest posterior probability is the maximum a posteriori, or MAP hypothesis.
  - In this case, we are looking for the MAP classification.
- Bayes' theorem is used to find the posterior probability

- $$P(c_i|x_1, x_2, \dots, x_n) = \frac{P(c_i)P(x_1, x_2, \dots, x_n|c_i)}{P(x_1, x_2, \dots, x_n)}.$$



## THE NAÏVE BAYES CLASSIFIER (2)

- Since  $P(x_1, \dots, x_n)$  is independent of  $c_i$ , we can eliminate it, and simply aim to find the classification  $c_i$ , for which the following is maximized:

$$P(c_i)P(x_1, x_2, \dots, x_n|c_i).$$

- We now **assume that all the attributes  $x_1, \dots, x_n$  are independent**, so  $P(c_i, x_1, x_2, \dots, x_n)$  can be rewritten as:

$$P(c_i) \prod_{j=1}^n P(x_j|c_i).$$

- The classification for which this is highest is chosen to classify the data.





# THE NAÏVE BAYES CLASSIFIER (3)

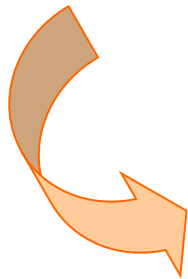
- $i^* = \operatorname{argmax} P(C_i | x, y, z)$   
 $= \operatorname{argmax} P(C_i, x, y, z)$
- $C_1 = A, C_2 = B, C_3 = C$
- For  $x=2, y=3, z=4$ ,  
compare
  - $P(A, x=2, y=3, z=4)$
  - $P(B, x=2, y=3, z=4)$
  - $P(C, x=2, y=3, z=4)$

x	y	z	Classification
2	3	2	A
4	1	4	B
1	3	2	A
2	4	3	A
4	2	4	B
2	1	3	C
1	2	4	A
2	3	3	B
2	2	4	A
3	3	3	C
3	2	1	A
1	2	1	B
2	1	4	A
4	3	4	C
2	2	4	A

# EXAMPLE

Training data

	Buy A	Buy B	Buy C	Buy D
Customer 1	Yes	Yes	No	Yes
Customer 2	Yes	No	Yes	No
...	...	...	...	...



compare

$$P(D|A, B, \bar{C})$$

$$P(\bar{D}|A, B, \bar{C})$$

Recommend

those products

with high ratios



# APPLICATIONS OF BAYESIAN NETWORK

- Bayesian Network
  - A set of random variables with an assumption of statistical dependency
  - Detection and Estimation (ML, MAP)
- Bayesian Classifier and Probabilistic Reasoning
- Markov Process
- N-Gram Language Model (e.g. trigram  $p(W_i|W_{i-2}, W_{i-1}))$ )
- Probabilistic Latent Semantic Indexing (PLSI)
- Gaussian Mixture Model
  - Continuous variables dependent of a state variable
- Hidden Markov Model
  - Hidden state sequence as Markov variables
  - Observation variables dependent of hidden state



# REFERENCES

- Artificial Intelligence Illuminated, Ben Coppin

