



Lecture Slides for

INTRODUCTION TO

Machine Learning

ETHEM ALPAYDIN
© The MIT Press, 2004

Edited for CS536 Fall 05- Rutgers University
Ahmed Elgammal

CHAPTER 3:

Bayesian Decision Theory

Probability and Inference

- Result of tossing a coin is $\in \{\text{Heads}, \text{Tails}\}$
- Random var $X \in \{1, 0\}$
Bernoulli: $P\{X=1\} = p_o^X(1 - p_o)^{(1-X)}$
- Sample: $\mathbf{X} = \{x^t\}_{t=1}^N$
Estimation: $p_o = \#\{\text{Heads}\} / \#\{\text{Tosses}\} = \sum_t x^t / N$
- Prediction of next toss:
Heads if $p_o > 1/2$, Tails otherwise

Classification

- Credit scoring: Inputs are income and savings.
Output is low-risk vs high-risk
- Input: $\mathbf{x} = [x_1, x_2]^T$, Output: $\mathbf{C} \in \{0, 1\}$
- Prediction:

$$\text{choose } \begin{cases} \mathbf{C} = 1 & \text{if } P(\mathbf{C} = 1 | x_1, x_2) > 0.5 \\ \mathbf{C} = 0 & \text{otherwise} \end{cases}$$

or equivalently

$$\text{choose } \begin{cases} \mathbf{C} = 1 & \text{if } P(\mathbf{C} = 1 | x_1, x_2) > P(\mathbf{C} = 0 | x_1, x_2) \\ \mathbf{C} = 0 & \text{otherwise} \end{cases}$$

Bayes' Rule

$$\text{posterior} \rightarrow P(\mathbf{C} | \mathbf{x}) = \frac{\overset{\text{prior}}{P(\mathbf{C})} \overset{\text{likelihood}}{p(\mathbf{x} | \mathbf{C})}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

Prior: The knowledge we have as to the value of C before looking at the observation x

$$P(\mathbf{C} = 0) + P(\mathbf{C} = 1) = 1$$

Bayes' Rule

$$\text{posterior} \rightarrow P(\mathbf{C} | \mathbf{x}) = \frac{\overset{\text{prior}}{P(\mathbf{C})} \overset{\text{likelihood}}{p(\mathbf{x} | \mathbf{C})}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

- Likelihood: the conditional probability that an event belonging to C has associated observation value x
- Evidence: the marginal probability that an observation x is seen regardless of C

$$p(\mathbf{x}) = p(\mathbf{x} | \mathbf{C} = 1)P(\mathbf{C} = 1) + p(\mathbf{x} | \mathbf{C} = 0)P(\mathbf{C} = 0)$$

Bayes' Rule

$$\text{posterior} \rightarrow P(\mathbf{C} | \mathbf{x}) = \frac{\overset{\text{prior}}{P(\mathbf{C})} \overset{\text{likelihood}}{p(\mathbf{x} | \mathbf{C})}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

- Combining the prior with what the data tells us, we can calculate the posterior probability $P(\mathbf{C} | \mathbf{x})$ after having seen the observation \mathbf{x}
- The posterior sum up to 1

$$P(\mathbf{C} = 0 | \mathbf{x}) + P(\mathbf{C} = 1 | \mathbf{x}) = 1$$

Bayes' Rule

$$\text{posterior} \rightarrow P(\mathbf{C} | \mathbf{x}) = \frac{\overset{\text{prior}}{P(\mathbf{C})} \overset{\text{likelihood}}{p(\mathbf{x} | \mathbf{C})}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

- Bayesian learning: from training data how to estimate $P(\mathbf{C})$ and $P(\mathbf{x} | \mathbf{C})$

Bayes' Rule

$$\text{posterior} \rightarrow P(\mathbf{C} | \mathbf{x}) = \frac{\overset{\text{prior}}{P(\mathbf{C})} \overset{\text{likelihood}}{p(\mathbf{x} | \mathbf{C})}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

$$P(\mathbf{C} = 0) + P(\mathbf{C} = 1) = 1$$

$$p(\mathbf{x}) = p(\mathbf{x} | \mathbf{C} = 1)P(\mathbf{C} = 1) + p(\mathbf{x} | \mathbf{C} = 0)P(\mathbf{C} = 0)$$

$$p(\mathbf{C} = 0 | \mathbf{x}) + p(\mathbf{C} = 1 | \mathbf{x}) = 1$$

Bayes' Rule: $K > 2$ Classes

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})}$$
$$= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)}$$
$$P(C_i) \geq 0 \text{ and } \sum_{i=1}^K P(C_i) = 1$$

- MAP: Maximal a Posteriori class: pick the class that maximizes the posterior probability

$$\text{choose } C_i \text{ if } P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$$

$$C_{\text{MAP}} \equiv \arg_{C_j} \max P(C_j | \mathbf{x})$$

Bayes' Rule: $K > 2$ Classes

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})}$$
$$= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)}$$

$P(C_i) \geq 0$ and $\sum_{i=1}^K P(C_i) = 1$

choose C_i if $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

$C_{\text{MAP}} \equiv \arg_{C_j} \max P(C_j | \mathbf{x}) = P(\mathbf{x} | C_j)P(C_j)$

$C_{\text{ML}} \equiv \arg_{C_j} \max P(\mathbf{x} | C_j)$

- ML: Maximal Likelihood Class: Special case; assume all class priors $P(C_j)$ are equal

11

Lecture Notes for E. Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Losses and Risks

- It may be the case that decisions are not equally good or costly.
- Actions: α_i
- Loss of α_i when the state is C_k : λ_{ik}
- Expected risk (Duda and Hart, 1973)

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x})$$

choose α_i if $R(\alpha_i | \mathbf{x}) = \min_k R(\alpha_k | \mathbf{x})$

12

Lecture Notes for E. Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Losses and Risks: 0/1 Loss

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$$

$$\begin{aligned} R(\alpha_i | \mathbf{x}) &= \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x}) \\ &= \sum_{k \neq i} P(C_k | \mathbf{x}) \\ &= 1 - P(C_i | \mathbf{x}) \end{aligned}$$

For minimum risk, choose the most probable class

13

Losses and Risks: Reject

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ \lambda & \text{if } i = K + 1, \quad 0 < \lambda < 1 \\ 1 & \text{otherwise} \end{cases}$$

$$R(\alpha_{K+1} | \mathbf{x}) = \sum_{k=1}^K \lambda P(C_k | \mathbf{x}) = \lambda$$

$$R(\alpha_i | \mathbf{x}) = \sum_{k \neq i} P(C_k | \mathbf{x}) = 1 - P(C_i | \mathbf{x})$$

choose C_i if $P(C_i | \mathbf{x}) > P(C_k | \mathbf{x}) \forall k \neq i$ and $P(C_i | \mathbf{x}) > 1 - \lambda$
reject otherwise

14

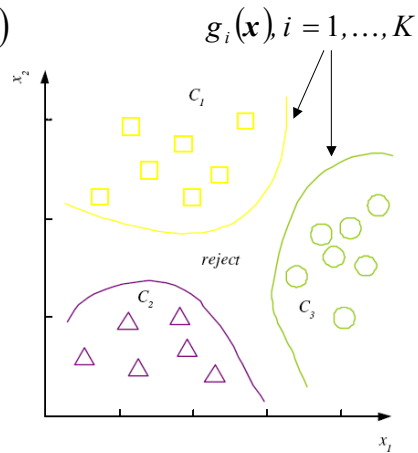
Discriminant Functions

choose C_i if $g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$

$$g_i(\mathbf{x}) = \begin{cases} -R(\alpha_i | \mathbf{x}) \\ P(C_i | \mathbf{x}) \\ p(\mathbf{x} | C_i)P(C_i) \end{cases}$$

K decision regions R_1, \dots, R_K

$$R_i = \{\mathbf{x} | g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$$



15

Lecture Notes for E. Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

$K=2$ Classes

- Dichotomizer ($K=2$) vs Polychotomizer ($K>2$)

- $g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$

$$\text{choose } \begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

- Log odds:

$$\log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})}$$

16

Lecture Notes for E. Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Utility Theory

- Prob of state k given evidence \mathbf{x} : $P(S_k|\mathbf{x})$
- Utility of α_i when state is k : U_{ik}
- Expected utility:

$$EU(\alpha_i | \mathbf{x}) = \sum_k U_{ik} P(S_k | \mathbf{x})$$

$$\text{Choose } \alpha_i \text{ if } EU(\alpha_i | \mathbf{x}) = \max_j EU(\alpha_j | \mathbf{x})$$

Value of Information

- Expected utility using \mathbf{x} only

$$EU(\mathbf{x}) = \max_i \sum_k U_{ik} P(S_k | \mathbf{x})$$

- Expected utility using \mathbf{x} and new feature z

$$EU(\mathbf{x}, z) = \max_i \sum_k U_{ik} P(S_k | \mathbf{x}, z)$$

- z is useful if $EU(\mathbf{x}, z) > EU(\mathbf{x})$

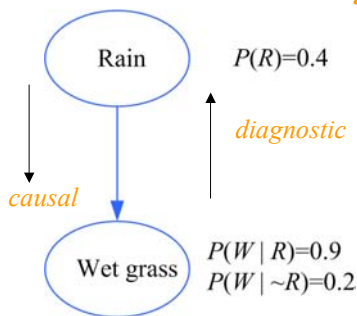
Bayesian Networks

- graphical models, probabilistic networks, belief networks
- **Nodes** are hypotheses (random vars) and the prob corresponds to our belief in the truth of the hypothesis
- **Arcs** are direct direct influences between hypotheses
- The **structure** is represented as a directed acyclic graph (DAG)
- The **parameters** are the conditional probs in the arcs
- (Pearl, 1988, 2000; Jensen, 1996; Lauritzen, 1996)

19

Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Causes and Bayes' Rule



Diagnostic inference:
 Knowing that the grass is wet,
 what is the probability that rain is
 the cause?

$$\begin{aligned}
 P(R|W) &= \frac{P(W|R)P(R)}{P(W)} \\
 &= \frac{P(W|R)P(R)}{P(W|R)P(R) + P(W|\sim R)P(\sim R)} \\
 &= \frac{0.9 \times 0.4}{0.9 \times 0.4 + 0.2 \times 0.6} = 0.75
 \end{aligned}$$

*Notice: knowing the grass is wet
 increases the probability of rain
 from 0.4 to 0.75*

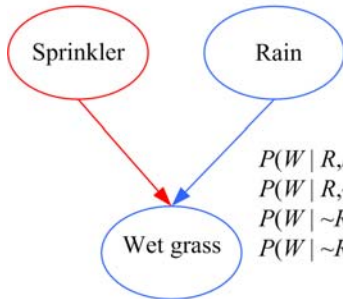
20

Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Causal vs Diagnostic Inference

$$P(S)=0.2$$

$$P(R)=0.4$$



$$\begin{aligned} P(W | R, S) &= 0.95 \\ P(W | R, \sim S) &= 0.90 \\ P(W | \sim R, S) &= 0.90 \\ P(W | \sim R, \sim S) &= 0.10 \end{aligned}$$

Causal inference: If the sprinkler is on, what is the probability that the grass is wet?

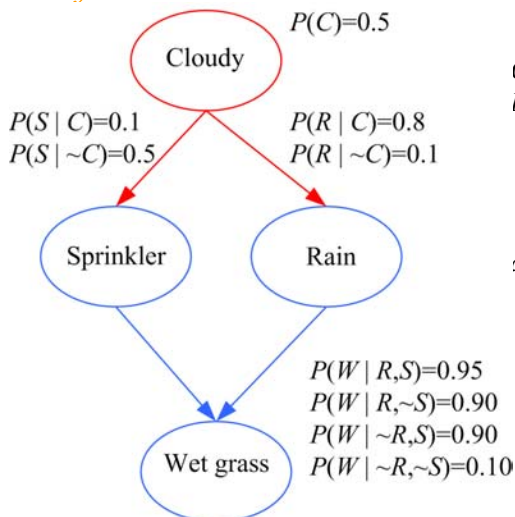
$$\begin{aligned} P(W|S) &= P(W|R,S) P(R|S) + P(W|\sim R,S) P(\sim R|S) \\ &= P(W|R,S) P(R) + P(W|\sim R,S) P(\sim R) \\ &= 0.95 \cdot 0.4 + 0.9 \cdot 0.6 = 0.92 \end{aligned}$$

Diagnostic inference: If the grass is wet, what is the probability that the sprinkler is on? $P(S|W) = 0.35 > 0.2 P(S)$

$P(S|R,W) = 0.21$ *Explaining away:* Knowing that it has rained decreases the probability that the sprinkler is on.

21

Bayesian Networks: Causes



$$P(C)=0.5$$

$$\begin{aligned} P(S | C) &= 0.1 \\ P(S | \sim C) &= 0.5 \end{aligned}$$

$$\begin{aligned} P(R | C) &= 0.8 \\ P(R | \sim C) &= 0.1 \end{aligned}$$

$$\begin{aligned} P(W | R, S) &= 0.95 \\ P(W | R, \sim S) &= 0.90 \\ P(W | \sim R, S) &= 0.90 \\ P(W | \sim R, \sim S) &= 0.10 \end{aligned}$$

Causal inference:

$$\begin{aligned} P(W|C) &= P(W|R,S) P(R,S|C) + P(W|\sim R,S) P(\sim R,S|C) + \\ &P(W|R,\sim S) P(R,\sim S|C) + P(W|\sim R,\sim S) P(\sim R,\sim S|C) \end{aligned}$$

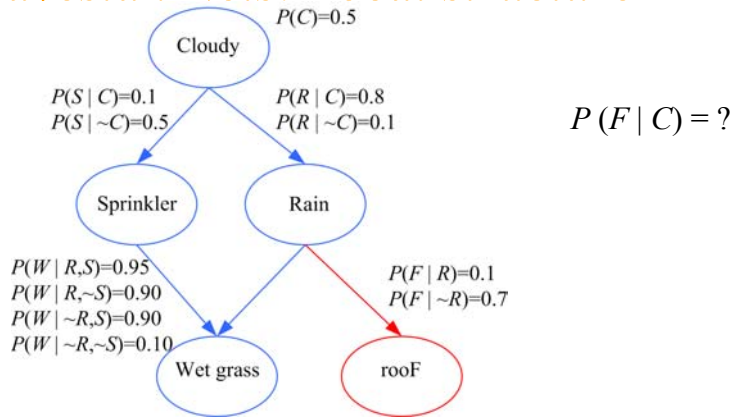
and use the fact that

$$P(R,S|C) = P(R|C) P(S|C)$$

Diagnostic: $P(C|W) = ?$

22

Bayesian Nets: Local structure



$$P(C, S, R, W, F) = P(C)P(S | C)P(R | C)P(W | S, R)P(F | R)$$

$$P(X_1, \dots, X_d) = \prod_{i=1}^d P(X_i | \text{parents}(X_i))$$

Bayesian Networks: Inference

$$P(C, S, R, W, F) = P(C)P(S|C)P(R|C)P(W|R, S)P(F|R)$$

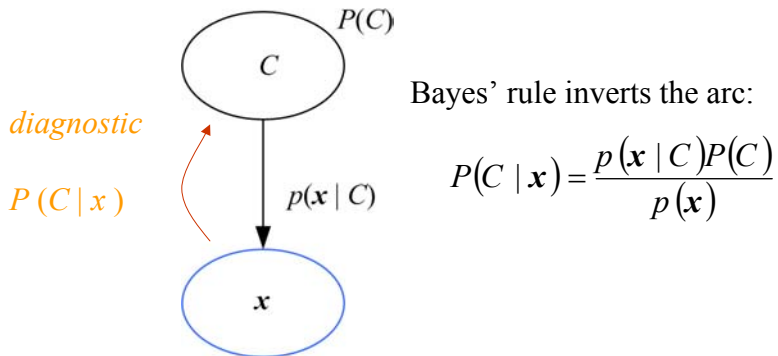
$$P(C, F) = \sum_S \sum_R \sum_W P(C, S, R, W, F)$$

$$P(F|C) = P(C, F) / P(C) \quad \text{Not efficient!}$$

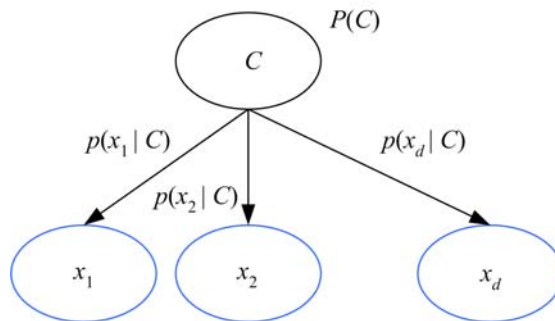
Belief propagation (Pearl, 1988)

Junction trees (Lauritzen and Spiegelhalter, 1988)

Bayesian Networks: Classification



Naive Bayes' Classifier



$$p(\mathbf{x}|C) = p(x_1|C) p(x_2|C) \dots p(x_d|C)$$