

Probabilistic Methods in Artificial Intelligence

Hadar Tal

Hebrew University of Jerusalem, Israel

May 24, 2024

1 Probability Review

Definition 1.1 (Probability Space)

A probability space is a triple (Ω, \mathcal{F}, P) where:

1. Ω is the sample space
2. \mathcal{F} is a σ -algebra of subsets of Ω
3. P is a probability measure on \mathcal{F} such that $P(\Omega) = 1$

Definition 1.2 (Joint Probability)

The joint probability of two events A and B is:

$$P(A, B) := P(A \cap B)$$

Definition 1.3 (Random Variable)

A random variable X is a function $X : \Omega \rightarrow \mathbb{R}$.

$$\text{Val}(X) = \text{Image}(X) = \{x \in \mathbb{R} : \exists \omega \in \Omega \text{ s.t. } X(\omega) = x\}$$

Definition 1.4 (Probability Mass Function (PMF))

The probability mass function of a random variable X is:

$$P(X = x) := P(\{\omega \in \Omega : X(\omega) = x\})$$

Definition 1.5 (Joint Distribution)

A joint distribution over a set of RVs $\mathcal{X} = \{X_1, X_2, \dots, X_n\}$ is a probability distribution $P_{\mathcal{X}} : \text{Val}(X_1) \times \text{Val}(X_2) \times \dots \times \text{Val}(X_n) \rightarrow [0, 1]$ defined by:

$$\forall x_1, \dots, x_n : x_i \in \text{Val}(X_i) \quad P_{\mathcal{X}}(x_1, x_2, \dots, x_n) := P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

Proposition 1.1 (Law of Total Probability)

For X, Y random variables, we can write:

$$P(X) = \sum_{y \in \text{Val}(Y)} P(X, Y = y)$$

Definition 1.6 (Conditional distribution)

For X, Y RVs, and for any $y \in \text{Val}(Y)$ where $P(Y = y) > 0$ the conditional distribution of X given $Y=y$ is:

$$P(X|y) := \frac{P_{X,Y}(X = x, Y = y)}{P_Y(Y = y)}$$

Proposition 1.2 (Chain Rule)

For any set of random variables X_1, X_2, \dots, X_n :

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots P(X_n|X_1, X_2, \dots, X_{n-1})$$

Proposition 1.3 (Bayes' Rule)

For any two random variables H, E :

$$P(H = h|E = e) = \frac{P(E = e|H = h)P(H = h)}{P(E = e)}$$

where we often call:

- $P(H = h)$ the **prior** probability
- $P(H = h|E = e)$ the **posterior** probability in light of evidence $E = e$
- $P(E = e|H = h)$ the **likelihood** of the evidence $E = e$ given the hypothesis $H = h$

Definition 1.7 (Marginal Independence)

Let P be a probability distribution over a set of random variables \mathcal{X} and let $X, Y \in \mathcal{X}$. We say that X is independent of Y , denoted $P \models X \perp Y$, if

$$P(X|Y) = P(X)$$

Definition 1.8 (Conditional Independence)

Let P be a probability distribution over a set of random variables \mathcal{X} and let $X, Y, Z \in \mathcal{X}$. We say that X is independent of Y given Z , denoted $P \models X \perp Y|Z$, if

$$P(X|Y, Z) = P(X|Z)$$

Lemma 1.1 (Equivalent Definitions of Conditional Independence)

Let P be a probability distribution over a set of random variables \mathcal{X} and let $X, Y, Z \in \mathcal{X}$. The following are equivalent:

1. $P \models X \perp Y|Z$
2. $P(X, Y|Z) = P(X|Z)P(Y|Z)$
3. $P(X, Y, Z) = P(X|Z)P(Y, Z)$
4. $\exists f, g : P(X, Y, Z) = f(X, Z)g(Y, Z)$

Theorem 1.1 (Properties of Conditional Independence)

Let P be a probability distribution over a set of random variables \mathcal{X} and let $X, Y, Z, W \in \mathcal{X}$. The following hold:

1. **Symmetry** - $(X \perp Y|Z) \implies (Y \perp X|Z)$
2. **Decomposition** - $(X \perp Y, W|Z) \implies (X \perp Y|Z) \wedge (X \perp W|Z)$
3. **Weak Union** - $(X \perp Y, W|Z) \implies (X \perp Y|W, Z)$
4. **Contraction** - $(X \perp Y|Z) \wedge (X \perp W|Y, Z) \implies (X \perp Y, W|Z)$
5. **Intersection** - For strictly positive distributions,

$$(X \perp Y|W, Z) \wedge (X \perp W|Y, Z) \implies (X \perp Y, W|Z)$$

2 Bayesian Networks

Definition 2.1 (Bayesian Network)

A Bayesian Network B is:

1. A directed acyclic graph (DAG) $G = (V, E)$
2. A set of conditional probability distributions $P_i(X_i|Pa(X_i))$ for each node X_i in the graph

the network defines a probability distribution:

$$P_B(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P_i(X_i|Pa(X_i))$$

Theorem 2.1 (Bayesian Network defines a probability distribution)

For any Bayesian Network B , $P_B(X_1, X_2, \dots, X_n)$ is a joint probability distribution over the variables X_1, X_2, \dots, X_n .

Definition 2.2 ($I_{LM}(G)$)

The **Local Markov Independencies Set** of a Bayesian Network B is the set of all independencies that hold in the network:

$$I_{LM}(G) = \{(X_i \perp ND(X_i)|Pa(X_i)) \quad : i \in |V|\}$$

Definition 2.3 (I-map)

A DAG G is an I-map of a distribution P if all independencies assumptions of G hold in P :

$$I_{LM}(G) \subseteq I(P)$$

Theorem 2.2 (If G is an I-map of P , then P factorizes according to G)

If G is an I-map of P , then

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i|Pa(X_i))$$

Theorem 2.3 (Independencies in P_B)

For P_B it holds for all i that

1. $X_i \perp ND(X_i)|Pa(X_i) \quad (I_{LM}(G))$
2. $P_B(X_i|ND(X_i)) = P_i(X_i|Pa(X_i))$

Definition 2.4 (Minimal I-map)

A DAG G is a minimal I-map of a distribution P if

1. G is an I-map of P
2. If $G' \subset G$ then G' is not an I-map of P