

BAYESIAN NETWORKS & STATISTICAL GENETICS

LECTURE 2.

Timo Koski

12.04.2010



X is a (**discrete**) random variable that assumes values in \mathcal{X} and Y is a (**discrete**) random variable that assumes values in \mathcal{Y} .

Random Variables \mathcal{X} and \mathcal{Y} are two discrete *state spaces*, whose generic elements are called *values* or *instantiations* and denoted by x_i and y_j , respectively.

$$\mathcal{X} = \{x_1, \dots, x_L\}, \mathcal{Y} = \{y_1, \dots, y_J\}.$$

$|\mathcal{X}|$ ($:=$ the number of elements in \mathcal{X}) $= L \leq \infty$, $|\mathcal{Y}| = J \leq \infty$. Unless otherwise stated the alphabets considered here are finite.

A two dimensional *joint (simultaneous) probability distribution* (simultan sannolikhetsfördelning) is a probability defined on the alphabet $\mathcal{X} \times \mathcal{Y}$

$$p(x_i, y_j) := P(X = x_i, Y = y_j). \quad (1)$$

$$p(x_i, y_j) \geq 0, \quad (2)$$

$$\sum_{i=1}^L \sum_{j=1}^J p(x_i, y_j) = 1. \quad (3)$$

Marginal distribution for X:

$$p(x_i) = \sum_{j=1}^J p(x_i, y_j). \quad (4)$$

Marginal distribution for Y:

$$p(y_j) = \sum_{i=1}^L p(x_i, y_j). \quad (5)$$

These notions can be extended to define the joint (simultaneous) probability distribution of n random variables and the marginal distributions of any subset thereof.

SIMULTANEOUS DISTRIBUTION AS A TABLE

X/Y	y_1	y_2	y_3
x_1	0.05	0.10	0.05
x_2	0.15	0.00	0.25
x_3	0.10	0.20	0.10

For example

$$p(X = x_2, Y = y_3) = 0.25$$



MARGINAL DISTRIBUTION

X/Y	y_1	y_2	y_3
x_1	0.05	0.10	0.05
x_2	0.15	0.00	0.25
x_3	0.10	0.20	0.10

$$p(X = x_1) = 0.05 + 0.10 + 0.05 = 0.20$$

$$p(X = x_2) = 0.15 + 0.00 + 0.25 = 0.40$$

$$p(X = x_3) = 0.10 + 0.20 + 0.10 = 0.40$$



The conditional probability for $X = x_i$ given $Y = y_j$ is

$$p(x_i | y_j) := \frac{p(x_i, y_j)}{p(y_j)}. \quad (6)$$

The conditional probability for $Y = y_j$ given $X = x_i$ is

$$p(y_j | x_i) := \frac{p(x_i, y_j)}{p(x_i)}. \quad (7)$$

Here we assume $p(y_j) > 0$ and $p(x_i) > 0$.

In other words

$$p(y_j | x_i) = \frac{\text{prob. for the event } \{X = x_i, Y = y_j\}}{\text{prob. for the event } \{X = x_i\}}.$$



Hence

$$\sum_{i=1}^L p(x_i | y_j) = 1, \sum_{j=1}^J p(y_j | x_i) = 1.$$

for every j and i , respectively.



In the table above

$$p(y_1|x_1) = \frac{p(x_1, y_1)}{p(x_1)} = \frac{0.05}{0.20} = \frac{5}{20}$$

$$p(y_2|x_1) = \frac{p(x_1, y_2)}{p(x_1)} = \frac{0.10}{0.20} = \frac{1}{2}$$

$$p(y_3|x_1) = \frac{p(x_1, y_3)}{p(x_1)} = \frac{0.05}{0.20} = \frac{5}{20}$$

$$\frac{5}{20} + \frac{1}{2} + \frac{5}{20} = 1$$

Next

$$P_X(A) := \sum_{x_i \in A} p(x_i) \quad (8)$$

is the probability of the event that X assumes a value in A , a subset of \mathcal{X} . Then one easily establishes the complement rule

$$P_X(A^c) = 1 - P_X(A), \quad (9)$$

where A^c is the complement of A , i.e., those outcomes which do not lie in A .



$$P_X(A \cup B) = P_X(A) + P_X(B) - P_X(A \cap B), \quad (10)$$

is immediate. If $A \cap B = \emptyset$, then $P_X(A \cap B) = 0$.

The conditional probability for $X = x_i$ given $X \in A$ is denoted by $P_X(x_i | A)$ and given by

$$P_X(x_i | A) = \begin{cases} \frac{P_X(x_i)}{P_X(A)} & \text{if } x_i \in A \\ 0 & \text{otherwise.} \end{cases} \quad (11)$$

Law of Total Probability

$$P(X \in A) = \sum_{j=1}^J P(X \in A \mid Y = y_j) p(Y = y_j) \quad (*)$$

$$P(Y \in B) = \sum_{i=1}^L P(Y \in B \mid X = x_i) p(X = x_i)$$

X and Y are *independent* random variables if and only if

$$p(x_i, y_j) = p(x_i) \cdot p(y_j) \quad (12)$$

for all pairs (x_i, y_j) in $\mathcal{X} \times \mathcal{Y}$. In other words all events $\{X = x_i\}$ and $\{Y = y_j\}$ are to be independent.

Independence

X/Y	y_1	y_2	y_3
x_1	0.05	0.10	0.05
x_2	0.15	0.00	0.25
x_3	0.10	0.20	0.10

$$p(X = x_1) = 0.05 + 0.10 + 0.05 = 0.20$$

$$p(Y = y_3) = 0.05 + 0.25 + 0.10 = 0.40$$

$$p(X = x_1) \cdot p(Y = y_3) = 0.08 \neq 0.05 = p(X = x_1, Y = y_3)$$

We say that X_1, X_2, \dots, X_n are **independent** random variables if and only if the joint distribution

$$p(x_{i_1}, x_{i_2}, \dots, x_{i_n}) = P(X_1 = x_{i_1}, X_2 = x_{i_2}, \dots, X_n = x_{i_n}) \quad (13)$$

equals

$$= p_{X_1}(x_{i_1}) \cdot p_{X_2}(x_{i_2}) \cdots p_{X_n}(x_{i_n}) \quad (14)$$

for every $x_{i_1}, x_{i_2}, \dots, x_{i_n} \in \mathcal{X}^n$.

Let Z be a (discrete) random variable that assumes values in $\mathcal{Z} = \{z_k\}_{k=1}^K$. If $p(z_k) > 0$,

$$p(x_i, y_j | z_k) = \frac{p(x_i, y_j, z_k)}{p(z_k)}.$$

Then we obtain as an identity

$$p(x_i, y_j | z_k) = \frac{p(x_i, y_j, z_k)}{p(y_j, z_k)} \cdot \frac{p(y_j, z_k)}{p(z_k)},$$

and again by definition of conditional probability

$$p(x_i \mid y_j, z_k) \cdot p(y_j \mid z_k).$$



Chain Rule So,

$$p(x_i, y_j | z_k) = \frac{p(x_i, y_j, z_k)}{p(y_j, z_k)} \cdot \frac{p(y_j, z_k)}{p(z_k)}$$

In other words,

$$p(x_i, y_j | z_k) = p(x_i | y_j, z_k) \cdot p(y_j | z_k). \quad (15)$$



A generalization

$$p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i \mid X_1, \dots, X_{i-1})$$

$$p(X_1 \mid X_0) = p(X_0).$$

Conditional Independence (IRRELEVANCE) The random variables X and Y are called *conditionally independent* given Z if

$$p(x_i, y_j | z_k) = p(x_i | z_k) \cdot p(y_j | z_k) \quad (16)$$

for all triples $(z_k, x_i, y_j) \in \mathcal{Z} \times \mathcal{X} \times \mathcal{Y}$. We write this as

$$X \perp Y | Z. \quad (17)$$

Y is irrelevant for X given Z , and X is irrelevant for Y given Z .

CONDITIONAL INDEPENDENCE

There are several equivalent ways of expressing conditional independence. We have for instance

$$X \perp Y|Z \Leftrightarrow p(x_i|y_j, z_k) = p(x_i|z_k).$$

To see this equivalence in one direction we write

$$p(x_i|y_j, z_k) = \frac{p(x_i, y_j, z_k)}{p(y_j, z_k)}$$

and assume $p(z_k) > 0$, so

$$\begin{aligned} &= \frac{p(x_i, y_j, z_k)}{p(z_k)} \frac{p(z_k)}{p(y_j, z_k)} \\ &= \frac{p(x_i, y_j | z_k)}{p(y_j | z_k)}, \end{aligned}$$

and assuming $X \perp Y|Z$ we get

$$= \frac{p(x_i|z_k) \cdot p(y_j|z_k)}{p(y_j | z_k)} = p(x_i|z_k),$$

as claimed.



$$p(X | Y) \cdot p(Y) = p(Y | X) \cdot p(X)$$

we have in a formal way

$$p(X | Y) = \frac{p(Y | X) \cdot p(X)}{p(Y)}.$$

But the marginal distribution $p(Y)$ is by the law of total probability (see (*) above) written as

$$p(y_j) = \sum_{i=1}^L p(y_j | x_i) p(x_i). \quad (18)$$



$$p(x_i | y_j) = \frac{p(y_j | x_i) \cdot p(x_i)}{\sum_{i=1}^L p(y_j | x_i) p(x_i)}. \quad (19)$$

TERMINOLOGY FOR BAYES' RULE

$p(X)$: A **Prior Distribution** on \mathcal{X} .

$p(X | Y)$: A **Posterior Distribution** on \mathcal{X} .

If X and Y are independent, then the prior distribution and posterior distribution are identical and there is no *learning*. Bayes' rule can be seen as just a formal identity derived from certain premises and definitions.



KTH Matematik

Learning and Bayes' Rule Bayes' rule gives a fundamental operation for *up-date of probability distributions* in response to observed information. The rule shows how knowledge about the occurrence of the event $Y = y_j$ is to be used to transform probabilities on \mathcal{X} . Probability is a degree of rational belief, Bayes' rule is a rule for reasoning.

If we learn that the event $Y = y_j$ is true, then we change $p(X)$ to a new probability distribution $p^*(X)$ according to Bayes' Rule.

$$p(X) \mapsto p^*(X) = p(X \mid Y = y_j)$$

So the posterior becomes the new prior.

$$p(X | Y) \propto p(Y | X) \cdot p(X)$$

Posterior \propto **likelihood** \times **prior**



JEFFREY'S RULE: KINEMATICS OF PROBABILITY (1)

Suppose you change your probabilities on \mathcal{Y} from the distribution $p(Y)$ to the distribution $p^*(Y)$. How should this change be propagated to the distribution on \mathcal{X} . R. Jeffrey thinks that Bayes' rule is not the only way. He suggests that $p(X)$ is updated to $p^*(X)$ defined by the rule

$$p^*(x_i) = \sum_{j=1}^J p(x_i | y_j) p^*(y_j), \quad (20)$$

where the assumption is that

$$p(x_i | y_j) = p^*(x_i | y_j).$$



$$p^*(x_i) = \sum_{j=1}^J p(x_i | y_j) p^*(y_j)$$

The argument is that if the event $X = x_i$ is 'not directly affected' by the flow of experience that was involved in $p(Y) \mapsto p^*(Y)$, then we should not use Bayes' rule. What does this mean ?

JEFFREY'S RULE: KINEMATICS OF PROBABILITY (3)

Let us say that e is the evidence that made us do $p(Y) \mapsto p^*(Y)$. Then we set

$$p^*(x_i) = p(x_i | e)$$

and get by Bayes rule and law of total probability

$$\begin{aligned} p(x_i | e) &= \sum_{j=1}^J p(x_i | y_j, e) p(y_j | e) \\ &= \sum_{j=1}^J p(x_i | y_j) p^*(y_j), \end{aligned}$$

if X and e are conditionally independent given Y .



But, are we permitted to write

$$\begin{aligned} p^*(x_i) &= p(x_i \mid e) \\ &= \frac{p(x_i, e)}{p(e)}, \end{aligned}$$

as $p(x_i, e)$ was not specified, if e was not a part of our knowledge base. E.g., e may not have been anticipated. Hence Jeffrey's rule seems more generally valid than Bayes' rule.

JEFFREY'S RULE: KINEMATICS OF PROBABILITY (5)

But even if $p(x_i, e)$ was not specified as a numerical quantity, we may still be permitted to apply conditional independence of X and e given Y by qualitative judgement.

Lesson: We shall specify conditional independencies instead of numerical joint distributions.



Consider X with values $\mathcal{X} = \{0, 1\}$ and $0 < \theta < 1$ with the probability table

p	$x = 1$	$x = 0$
$p(x)$	θ	$1 - \theta$

then we call X a Bernoulli random variable with the *probability of success* θ . We write

$$X \in Be(\theta).$$

We refer to θ as the *parameter* of the Bernoulli distribution p .



If X_1, X_2, \dots, X_n are independent and $X_i \in Be(\theta)$, then

$$\begin{aligned} p(1, 1, 0, 1, 0, 1, 1) &= \\ &= \theta \cdot \theta \cdot (1 - \theta) \cdot \theta \cdot (1 - \theta) \cdot \theta \cdot \theta \\ &= \theta^5 \cdot (1 - \theta)^2. \end{aligned}$$

A SEQUENCE OF FLIPS OF A THUMBTAACK

If we throw a thumbtack in the air, it will come to rest either on its point (0) or on its head (1). Suppose we flip the thumbtack n times (fixing n in advance), making sure that the physical properties of the thumbtack and the conditions under which it is flipped remain stable over time. We let \mathbf{x} denote the sequence of outcomes of the flips

$$\mathbf{x} = x_{i_1} x_{i_2} \dots x_{i_n}, x_{i_l} \in \{0, 1\}.$$



As our model we take the bits in \mathbf{x} to be outcomes of $X_i \in \text{Be}(\theta)$ conditionally independent given $\Theta = \theta$.

$$X_i \perp X_j \mid \Theta \quad \text{for all } i \neq j$$

Not only are pairs independent, but all subsets of X_{i_1}, \dots, X_{i_k} . In subjective probability the parameters of a probability model are regarded as random variables.

Hence

$$\begin{aligned} P(\mathbf{x} \mid \Theta = \theta) &= \prod_{l=1}^n \theta^{x_{i_l}} \cdot (1 - \theta)^{1 - x_{i_l}} = \\ &= \theta^{\sum_{l=1}^n x_{i_l}} \cdot (1 - \theta)^{n - \sum_{l=1}^n x_{i_l}} = \theta^k \cdot (1 - \theta)^{n - k}, \end{aligned}$$

if $\sum_{l=1}^n x_{i_l} = k$.

Find the model that is in some sense best for \mathbf{x} . In the thumbtack example we understand this as follows. We have observed n outcomes of flips of a thumbtack and wish to determine which of the values θ that best describes this set of flips.

LEARNING ABOUT PROBABILITIES: BAYES' RULE FOR PARAMETERS

$$p(\Theta | X) = \frac{p(X | \Theta) \cdot p(\Theta)}{p(X)}.$$

$p(\Theta | X)$ and $p(\Theta)$ are probability densities.



LEARNING ABOUT PROBABILITIES: BAYES' RULE FOR PARAMETERS

Θ is given a probability density function $f_{\Theta}(\theta)$, called the *prior density*.

$$f_{\Theta}(\theta) \geq 0, 0 \leq \theta \leq 1$$

and $f_{\Theta}(\theta) = 0$ elsewhere, and

$$\int_0^1 f_{\Theta}(\theta) d\theta = 1.$$

Also $P(a < \Theta \leq b) = \int_a^b f_{\Theta}(\theta) d\theta$.

$$f_{\Theta|\mathbf{x}}(\theta | \mathbf{x}) = \frac{P(\mathbf{x} | \Theta = \theta) \cdot f_{\Theta}(\theta)}{\int_0^1 P(\mathbf{x} | \Theta = \theta) \cdot f_{\Theta}(\theta) d\theta}, 0 \leq \theta \leq 1 \quad (21)$$

and zero elsewhere. Due to the standardization $f_{\Theta|\mathbf{x}}(\theta | \mathbf{x})$ is a probability density for Θ .

The posterior $f_{\Theta|\mathbf{x}}(\theta | \mathbf{x})$ expresses our updated belief in the statement that θ is the probability of success given that we have observed \mathbf{x} .

Let us consider the *uniform prior* given by

$$f_{\Theta}(\theta) = \begin{cases} 1 & 0 \leq \theta \leq 1 \\ 0 & \text{elsewhere.} \end{cases}$$

The uniform prior is often interpreted as a representation of complete ignorance. This is a special case of a *Beta density*.

$$\int_0^1 P(\mathbf{x} \mid \Theta = \theta) \cdot f_{\Theta}(\theta) d\theta = \int_0^1 \theta^k \cdot (1 - \theta)^{n-k} d\theta = \frac{k!(n-k)!}{(n+1)!}$$

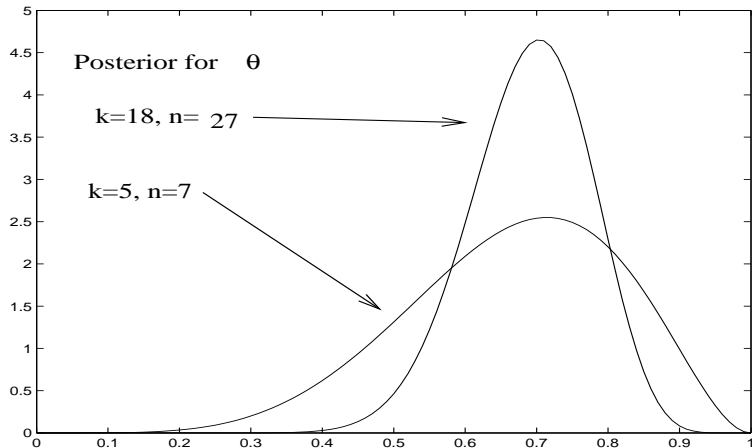
by the *Beta integral*.

POSTERIOR IS A BETA DENSITY

$$f_{\Theta|\mathbf{x}}(\theta | \mathbf{x}) = \begin{cases} \frac{(n+1)!}{k!(n-k)!} \cdot \theta^k (1-\theta)^{n-k} & 0 \leq \theta \leq 1 \\ 0 & \text{elsewhere.} \end{cases} \quad (22)$$



POSTERIOR DENSITIES FOR θ IN $\text{BE}(\theta)$



The **maximum likelihood estimate** MLE, $\hat{\theta}_{\text{ML}}$ of θ , is defined by

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \operatorname{argmax}_{0 \leq \theta \leq 1} P(\mathbf{x} \mid \Theta = \theta) \\ &= \operatorname{argmax}_{0 \leq \theta \leq 1} \theta^k \cdot (1 - \theta)^{n-k}.\end{aligned}$$

The **maximum a posterior estimate** MAP $\hat{\theta}_{\text{MAP}}$ of θ is defined by

$$\hat{\theta}_{\text{MAP}} = \operatorname{argmax}_{0 \leq \theta \leq 1} f_{\Theta|\mathbf{x}}(\theta | \mathbf{x})$$

*Find the parameter value within the model that gives the (training) sequence \mathbf{x} the highest possible probability. The probability $P(\mathbf{x} \mid \Theta = \theta)$ regarded as a function of θ is known as the *likelihood function**

$$L_{\mathbf{x}}(\theta) = P(\mathbf{x} \mid \Theta = \theta).$$

The likelihood function $L_{\mathbf{x}}(\theta)$ thus compares the plausibilities of different models for given \mathbf{x} .

$$-\log L_{\mathbf{x}}(\theta) = -\log P(\mathbf{x} \mid \Theta = \theta).$$

is called the log likelihood function.

Maximization of the likelihood function or the log likelihood function by calculus gives

$$\hat{\theta}_{\text{ML}} = \frac{k}{n}. \quad (23)$$

What is $\hat{\theta}_{\text{MAP}}$ in this case ?