

Bayesian inference for time series

Wen Yu

Departamento de Control Automático
CINVESTAV-IPN

The probability of the variable A takes a is

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

Alternatives \rightarrow "add"

$$P(A = a_1 \text{ or } A = a_2) = P(A = a_1) + P(A = a_2)$$

Normalisation

$$\sum_{\text{all possible } a} P(A = a) = 1$$

Joint probability

$P(A = a, B = b)$ the probability that both $A = a$ and $B = b$ occur

Conditional probability

$P(A = a \mid B = b)$ the probability that $A = a$ occurs given the knowledge B

Product rule

$$\begin{aligned} P(A = a, B = b) \\ &= P(A = a) P(B = b \mid A = a) \\ &= P(B = b) P(A = a \mid B = b) \end{aligned}$$

Independence, iff A and B are independent:

$$\begin{aligned} P(A = a \mid B = b) &= P(A = a) \\ P(B = b \mid A = a) &= P(B = b) \\ P(A = a, B = b) &= P(A = a) P(B = b) \end{aligned}$$

Continuous variables, "Sum" \rightarrow "Integral", the probability that X lies between x and $(x + dx)$ is $p(x)dx$, $p(x)$ is a probability density function

$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} p(x) dx$$
$$\int_{-\infty}^{+\infty} p(x) dx = 1$$
$$\int_{-\infty}^{+\infty} p(x, y) dy = p(x)$$

Expectations: averages over a time series X

$$E[A] = \sum_a P(A = a) A$$

$$E[x] = \int_{-\infty}^{+\infty} xp(x)dx$$

If X and Y are independent

$$E[x \times y] = E[x] \times E[y]$$

Bayes' theorem, Bayes' law, Bayes' rule

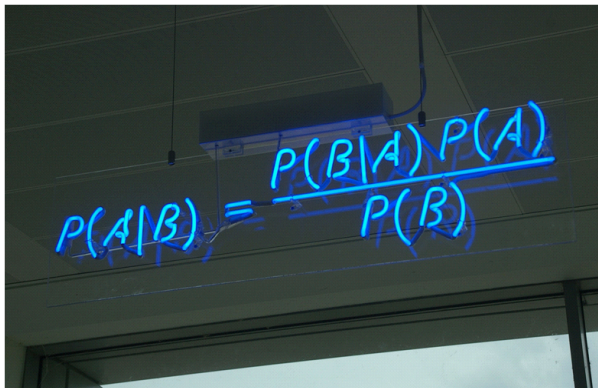
A photograph of a blue neon sign mounted on a ceiling. The sign displays the mathematical formula for Bayes' theorem: $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$. The neon is bright blue and the background is dark.
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Figure: Blue neon at the offices of Autonomy in Cambridge

Bayes' theorem

10 red balls and 5 blue balls in one bag

- "Forward" probability: pick up 1 ball, the probability of red ball is $\frac{10}{15}$, the probability of blue ball is $\frac{5}{15}$.
- "Reverse" probability: there are 15 red and 15 blue balls in one bag, pick up 2 balls, the probability of the color are similar
- "Reverse" probability: there are two bags: 10 red + 5 blue, 5 red + 10 blue. Randomly select one bag, pick up 12 balls and they are 8 red + 4 blue, the probability of bag A ?

Bayes' theorem

Relates current probability to prior probability (degree of belief to account for evidence)

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

where

A and B are events,

$P(A)$ and $P(B)$ are the probabilities of A and B without regard to one other

$P(A | B)$ is the conditional probability (probability of A given that B is true)

$P(B | A)$ is the probability of B given that A is true.

Bayes' theorem

- $P(A)$ is the probability of cause A before effect B is known, it is called the *prior probability* of A .
- $P(A | B)$ is the probability of cause A after effect B is known, and it is called the *posterior probability* of A .
- $P(B | A)$ is the likelihood of B when A

Proof.

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \quad P(B | A) = \frac{P(B \cap A)}{P(A)}$$
$$P(B) P(A | B) = P(A) P(B | A) = P(A \cap B)$$
$$P(A | B) = \frac{P(A) P(B | A)}{P(B)}$$



Bayes' theorem

$$\begin{aligned}P(A | B) &= \frac{P(B|A)P(A)}{P(B)} \\&= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\neg A)P(\neg A)} \\&= \frac{P(B|A_i)P(A_i)}{\sum_j P(B|A_j)P(A_j)}\end{aligned}$$

Bayes' theorem- application

Example

A patient comes into the doctor's office. A blood test for cancer is given and the test result is POS. The test is 95% accurate. Only 0.8% of the people in the U.S. have this form of cancer. What is probability of having cancer of this patient $P(\text{cancer} | \text{POS})$?

$$P(\text{cancer}) = 0.008, P(\text{POS} | \text{cancer}) = 0.95$$

$$P(\text{POS}) = P(\text{POS}, \text{cancer}) + P(\text{POS}, \neg\text{cancer})$$

$$= P(\text{POS} | \text{cancer}) P(\text{cancer}) + P(\text{POS} | \neg\text{cancer}) P(\neg\text{cancer})$$

$$= 0.95 \times 0.008 + 0.05 \times 0.992$$

The probability of having cancer is

$$P(\text{cancer} | \text{POS}) = \frac{P(\text{POS} | \text{cancer}) P(\text{cancer})}{P(\text{POS})} = 0.133$$

This patient has a 13% chance of having cancer.

Bayes' theorem- application

Spelling check

"lates" is "late" or " latest "? we need probability, not rules

Given a written word W , we need to find a correct word C in all correct words,

$$\max P(C | W)$$

By Bayes' theorem

$$\max P(C | W) = \max \frac{P(W | C) P(C)}{P(W)}$$

Since we can write any word, $P(W)$ is the same always

$$\max P(C | W) = \max P(W | C) P(C)$$

Bayes' theorem- application

- $P(C)$ is the probability of all English words. It is the language model. For example $P("the")$ is very high, $P("teh")$ is very low.
- $P(W | C)$ is the probability: the person want input C , however he types W . It is error model.
- Language Model: $P(C)$ can be obtained by extracting a big text file
- Error Model: edit distance, in 80-95% spelling error have "edit distance"=1

Bayes' theorem- continuous time

For continuous random variables X and Y , Bayes' Theorem is formulated in terms of densities:

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

where

$P(x, y)$ is the joint probability distribution of X and Y

$P(x)$ is the prior probability of X , $P(x | y)$ is the likelihood of Y when $X = x$

$P(y | x)$ is the posterior probability of Y when $X = x$

$$P(y | x) = \frac{P(x | y) P(y)}{\int_{-\infty}^{\infty} P(x | y) P(y) dy}$$

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

$$\text{posterior} = \frac{\text{prior} \times \text{linkhood}}{\text{evidence}}$$

Naive Bayesian

For input-output model

$$y = f(x) = P(Y | X) = \frac{P(X | Y_i) P(Y_i)}{\sum_i P(X | Y_i) P(Y_i)}$$

where $X = [x_1 \cdots x_n]$, Y_i has i values c_i

Assume that each feature x_i is conditionally independent of every other feature

$$P(X | Y = c_i) = P(x_1 \cdots x_n | Y = c_i) = \prod_j P(X = x_j | Y = c_i)$$

So

$$P(Y = c_i | X) = \frac{\prod_j P(X = x_j | Y = c_i) P(Y = c_i)}{\sum_i P(Y = c_i) \prod_j P(X = x_j | Y = c_i)}$$

We want to the best value of Y , i.e.,

$$y = f(x) = \arg \max_{c_i} \frac{\prod_j P(X = x_j | Y = c_i) P(Y = c_i)}{\sum_i P(Y = c_i) \prod_j P(X = x_j | Y = c_i)}$$

the denominator for all c_i is the same

$$y = f(x) = \arg \max_{c_i} P(Y = c_i) \prod_j P(X = x_j | Y = c_i)$$

Naive Bayes Classifier: select the most likely classification V given the attribute values $x_1 \cdots x_n$

$$V = \max_{a_j \in A} P(a_j | x_1 \cdots x_n) = \max_{a_j \in A} P(a_j) \prod_{i=1}^n P(x_i | a_j)$$

It is equivalent to the joint probability model, because

$P(A | B) = \frac{P(B|A)P(A)}{P(B)}$ the denominator does not depend on A , it is effectively constant.

Maximum-likelihood estimation (MLE)

For independent and identically distributed samples, $x_1 \cdots x_n$, coming from unknown probability density function P_0 . Assume P_0 belongs to a certain family of distributions $\{P(x|\theta), \theta \in \Theta\}$. where θ is parameters for this family. The parametric model is

$$P_0 = P(x|\theta_0)$$

The joint density function for all observations

$$P(x_1 \cdots x_n | \theta) = \prod_{i=1}^n P(x_i | \theta)$$

When θ can be changed and x_i is fixed, the $P(x_1 \cdots x_n | \theta)$ is called likelihood,

$$L(\theta; x_1 \cdots x_n) = P(x_1 \cdots x_n | \theta)$$

where $x_1 \cdots x_n$ are observations

Maximum-likelihood estimation (MLE)

In practice, we use the log-likelihood:

$$\ln L(\theta; x_1 \cdots x_n) = \sum_{i=1}^n \ln P(x_i \cdots x_n | \theta), \text{ or } \hat{L} = \frac{1}{n} \ln L$$

\hat{L} estimates the expected log-likelihood of a single observation in the model.

MLE

$$\theta^* = \max_{\theta} \hat{L}$$

Bayes' theorem:

$$P(\theta; | x_1 \cdots x_n) = \frac{P(x_1 \cdots x_n | \theta) P(\theta)}{P(x_1 \cdots x_n)}$$

Bayesian estimator is obtained by maximizing $P(x_1 \cdots x_n | \theta) P(\theta)$ with respect to θ . If we further assume that the prior $P(\theta)$ is a uniform distribution, the Bayesian estimator is obtained by maximizing the likelihood function $P(x_1 \cdots x_n | \theta)$.

Thus the Bayesian estimator coincides with MLE for a uniform prior distribution $P(\theta)$.

For the normal distribution $N(\mu, \sigma^2)$, the probability density function is

$$P(x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

the likelihood of n independent identically distributed samples is

$$L(\theta; x_1 \cdots x_n) = P(x_1 \cdots x_n | \mu, \sigma^2) = \prod_{i=1}^n P(x_i | \mu, \sigma^2) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right]$$

$$\ln L = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

Because

$$\exp \left[-\frac{\sum (x - \mu)^2}{2\sigma^2} \right] = \exp \left[-\frac{\sum (x - \bar{x})^2 + n \sum (\bar{x} - \mu)^2}{2\sigma^2} \right]$$

$$\frac{\partial}{\partial \mu} \ln L = 0$$

$$\mu = \frac{1}{n} \sum x_i$$

$$\frac{\partial}{\partial \sigma} \ln L = 0$$

$$\sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

If $x_1 \cdots x_n$ have the same probability density function, such as Bernoulli distribution

$$p(x, \theta) = \begin{cases} (1 + \theta) x^\theta & 0 < x < 1 \\ 0 & \textit{otherwise} \end{cases}$$

Construct the likelihood function

$$L = \begin{cases} (1 + \theta)^n \prod_{i=1}^n x_i^\theta & 0 < x < 1 \\ 0 & \textit{otherwise} \end{cases}$$

and

$$\ln L = \begin{cases} n \ln(1 + \theta) + \theta \sum \ln x_i & 0 < x_i < 1 \\ 0 & \textit{otherwise} \end{cases}$$

$$\frac{\partial}{\partial \theta} \ln L = \frac{n}{1 + \theta} + \sum \ln x_i = 0$$

MLE of θ is

$$\hat{\theta} = \frac{n}{\sum \ln x_i} - 1$$

Normal model

A random variable X is said to be normally distributed with mean θ and variance σ^2 , if the density of X is

$$p(x | \theta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2\sigma^2} (x - \theta)^2 \right] \quad (1)$$

If $(x_1 \dots x_n | \theta, \sigma^2) \sim i.i.d. normal(\theta, \sigma^2)$, then the sampling density is

$$\begin{aligned} p(x_1 \dots x_n | \theta, \sigma^2) &= \prod p(x_i | \theta, \sigma^2) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left[-\frac{1}{2\sigma^2} \sum (x_i - \theta) \right] \end{aligned} \quad (2)$$

The problem of the model training is to estimate the two parameters, θ and σ^2 , from the data $x_1 \dots x_n$.

Bayesian inference for time series

Exponential model

The inverse exponential distribution

$$p(\beta) = \frac{1}{\beta} e^{(-X/\beta)} \quad (3)$$

where $X = [x_1 \dots x_n]$, $\beta > 0$. This distribution can be calculated with the cumulative of the previous distribution with $\lambda = 1/\beta$,

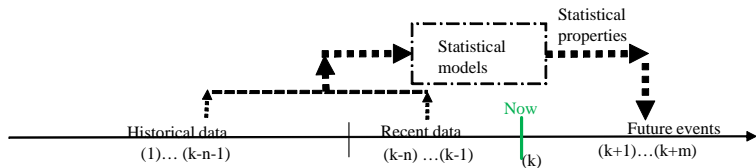
$$p(\lambda) = \lambda e^{-\lambda X} \quad (4)$$

The mean and variance of the exponential distribution can be represented as

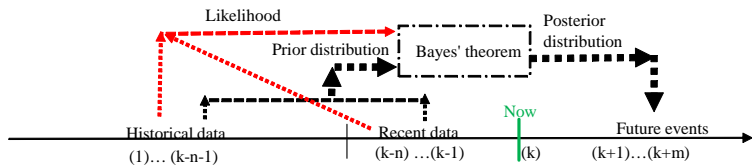
$$\mu = \beta, \quad \sigma^2 = \beta^2 \quad (5)$$

The problem of the model training is to estimate the parameter β from the data $x_1 \dots x_n$.

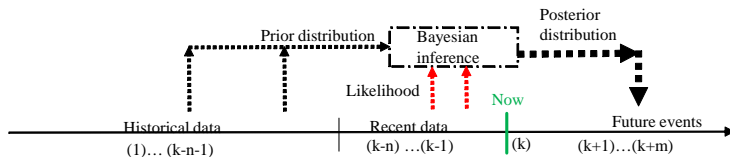
Bayesian inference for time series



Bayesian inference for time series



Bayesian inference for time series



Bayes' theorem in model form

Bayes' theorem in model form is written as:

$$\begin{aligned} p(\textit{parameter} \mid \textit{data}) &= \frac{p(\textit{data} \mid \textit{parameter})p(\textit{parameter})}{p(\textit{data})} \\ &= \frac{p(x \mid \theta)p(\theta)}{p(x)} \end{aligned} \tag{6}$$

- $p(\textit{parameter})$ is the prior distribution. It represents our beliefs about the true value of the parameters,
- $p(\textit{data} \mid \textit{parameter})$ is the likelihood distribution
- $p(\textit{parameter} \mid \textit{data})$ is the posterior distribution. This is the distribution representing the parameter values after we have calculated everything taking the observed data into account.

Bayes' theorem in model form

The simplified Bayes' theorem is

$$\begin{aligned} p(\textit{parameter} \mid \textit{data}) \\ \propto p(\textit{parameter}) \times p(\textit{data} \mid \textit{parameter}) \end{aligned} \quad (7)$$

- the posterior distribution can be obtained from the prior distribution and the likelihood function.
- the prior distribution and the likelihood function are approximated by Monte Carlo method.

MCMC or Gibbs sampling

The Monte Carlo procedure allows to approximate the distribution $p(x)$ in $[a, b]$ by sampling random variables.

Because

$$\begin{aligned} \int_a^b f(x) dx &= \int_a^b \frac{f(x)}{p(x)} p(x) dx \\ &= \int_a^b F(x) p(x) dx = E[F(x)] \end{aligned} \quad (8)$$

where $F(x) = \frac{f(x)}{p(x)}$.

We can sample $x_1 \cdots x_n$ to obtain the distribution $p(x)$ as

$$p(x_i) \approx \frac{x_i}{\sum_{i=1}^n x_i} \quad (9)$$

So

$$E[F(x)] \approx \frac{1}{n} \sum_{i=1}^n \frac{f(x_i)}{p(x_i)}$$