

Lecture 4

Generative Models for Discrete Data - Part 1

Luigi Freda

ALCOR Lab
DIAG
University of Rome "La Sapienza"

October 6, 2017

- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - Number Game
 - Likelihood
 - Prior
 - Posterior
 - Posterior Predictive

- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - Number Game
 - Likelihood
 - Prior
 - Posterior
 - Posterior Predictive

Bayesian Approach

General Paradigm

- when making inference about quantities we'll adopt a **Bayesian approach**
- consider the problem of estimating a parameter vector θ of a certain model starting from a dataset \mathcal{D}
- *before* observing the data, we capture our assumptions about θ in the form of a **prior** probability distribution $p(\theta)$
- a probability distribution $p(\mathcal{D}|\theta)$, aka **likelihood function**, expresses how probable the observed dataset \mathcal{D} is for a given parameter setting θ
- Bayes' theorem allows us to express the **posterior** probability distribution

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

which represent the uncertainty about θ *after* we have observed \mathcal{D}

Bayesian Approach

General Paradigma

- we can restate the Bayes' theorem

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

in simple words as

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

- the term $p(\mathcal{D})$ can be seen as a normalization constant which ensure a valid probability on the left-hand side which integrates to one
- in fact, one can integrate/sum both side of the above equation w.r.t. θ

$$p(\mathcal{D}) = \int_{-\infty}^{\infty} p(\mathcal{D}|\theta)p(\theta)d\theta \quad \text{for a continuous RV}\theta$$

$$p(\mathcal{D}) = \sum_j p(\mathcal{D}|\theta_j)p(\theta_j) \quad \text{for a discrete RV}\theta$$

Bayesian vs Frequentist

Likelihood Function

the likelihood function $p(\mathcal{D}|\theta)$ plays a central role and may be differently used

frequentist paradigm

(parameter fixed, data random)

- θ is considered a fixed parameter
- an estimator δ is designed to determine the value of θ to some data, so $\hat{\theta} = \delta(\mathcal{D})$
- uncertainty on θ estimate are given by considering the distribution of possible datasets \mathcal{D}

Bayesian paradigm

(data fixed, parameter random)

- there is only one dataset \mathcal{D} , namely the one observed
- uncertainty in the parameter θ is expressed by the prior $p(\theta)$
- Bayes rule is used to compute the posterior

- Maximum A Posteriori estimate (**MAP**)

$$\theta_{MAP} = \arg \max_{\theta} p(\mathcal{D}|\theta)p(\theta) = \arg \max_{\theta} [\log(p(\mathcal{D}|\theta)) + \log(p(\theta))]$$

which is the mode of the posterior $p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$

- Maximum Likelihood Estimate (**MLE**)

$$\theta_{MLE} = \arg \max_{\theta} p(\mathcal{D}|\theta) = \arg \max_{\theta} \log(p(\mathcal{D}|\theta))$$

which is the mode of the likelihood function $p(\mathcal{D}|\theta)$

- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - **Number Game**
 - Likelihood
 - Prior
 - Posterior
 - Posterior Predictive

Concept Learning

- **concept learning** can be thought of as the process of learning the meaning of a word
- psychological research result: people can learn from positive examples alone¹
- we consider the case of having as input only a sequence of *positive examples* x_i of a word/concept C
- the result of the process can be represented as a binary classification problem

$$f(x) = \begin{cases} 1 & \text{if } x \text{ represents } C \\ 0 & \text{otherwise} \end{cases}$$

- once $f(x)$ is learnt, it can be used for classifying future instances \tilde{x}
- **problem:** how can we learn the binary classifier $f(x)$ which can be used on future data?

¹as an example, consider a child as he/she learns to understand the meaning of the word "dog" being just shown dogs

Number Game

An Example of Concept Learning

number game

- choose a simple arithmetical concept C such as "*prime number*", "*even number*", "*a number between 1 and 10*"
- we are given a set of positive examples $\mathcal{D} = \{x_1, \dots, x_N\}$ drawn from C
- for simplicity we assume x_i is an integer between 1 and 100
- **problem:** (binary classification) given a new test case \tilde{x} , estimate if $\tilde{x} \in C$ or not

how can we solve this problem in a machine?

- in a Bayesian approach we are going to estimate the full **posterior predictive distribution** $p(\tilde{x} \in C | \mathcal{D})$

Number Game

Hypothesis Space and Likelihood

let's restate the problem

- assume we have a finite **hypothesis space** of concepts \mathcal{H}
- \mathcal{H} collects all the interesting/reasonable hypotheses $h_i \in \mathcal{H}$ which can represent C (e.g, $h_1 =$ "odd numbers", $h_2 =$ "even numbers", $h_{two} =$ "powers of two", $h_{end6} =$ "all numbers ending in 6", etc)
- we observe a set of positive examples $\mathcal{D} = \{x_1, \dots, x_N\}$ drawn from C
- C is represented by an unknown hypothesis $h \in \mathcal{H}$ (h is a discrete RV)
- the **version space** for a given dataset \mathcal{D} is the subset of hypotheses in \mathcal{H} which are consistent with \mathcal{D}

as the numbers x_i are observed, we would like to

- 1 estimate the distribution of h
- 2 be able to infer the most probable hypothesis at each step
- 3 continuously shrink the version space by using Bayesian inference
- 4 predict if a new \tilde{x} belongs to the concept "represented" by \mathcal{D}

- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - Number Game
 - Likelihood
 - Prior
 - Posterior
 - Posterior Predictive

Number Game

Likelihood

- assume the numbers x_i are uniformly sampled from the extension of an hypothesis h (i.e., the set of numbers that belong to it)

$$p(x_i|h) = \frac{1}{|h|}$$

where $|h|$ denotes the size of the hypothesis h

- assuming N independent samples one has

$$p(\mathcal{D}|h) = \prod_{i=1}^N p(x_i|h) = \left[\frac{1}{|h|} \right]^N$$

- **Occam's razor** (*lex parsimoniae*): among competing hypotheses, the one with the fewest assumptions should be selected
- **size principle**: the model should favor the simplest (smallest) hypothesis consistent with data (equivalent to Occam's razor and implemented by the above $p(\mathcal{D}|h)$ definition)

Number Game

Likelihood

- let's check how this works

$$p(\mathcal{D}|h) = \left[\frac{1}{|h|} \right]^N$$

assume $\mathcal{D} = \{16\}$, in this case

- $p(\mathcal{D}|h_{two}) = 1/6$ (powers of two)
- $p(\mathcal{D}|h_{even}) = 1/50$ (even numbers)
- $p(\mathcal{D}|h_{end6}) = 1/10$ (number ending with 6)

- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - Number Game
 - Likelihood
 - **Prior**
 - Posterior
 - Posterior Predictive

Number Game

Prior

- **prior** is the mechanism by which background knowledge can be brought to bear a problem
- in this case the prior $p(h)$ assigns a **subjective** probability to each hypothesis $h \in \mathcal{H}$ (i.e., $p(h)$ is the probability that h actually represents the concept C)
- the subjectivity of the prior is controversial but is quite useful since it allows rapid learning by using previously refined knowledge

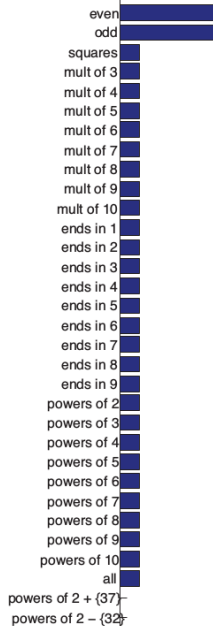
- suppose $\mathcal{D} = \{16, 8, 2, 64\}$, a child and a math professor will certainly reach different answers since they presumably start from **different priors** and **different hypothesis spaces**

Number Game

Prior

in this example we use **simple prior** $p(h)$ which

- basically puts a **uniform prior** on 30 simple arithmetical concepts h_i
- puts more prior weight on the concepts of even and odd numbers
- adds some "unnatural concepts" with low prior weights (for making things more interesting)



- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - Number Game
 - Likelihood
 - Prior
 - **Posterior**
 - Posterior Predictive

Number Game

Posterior

- by using Bayes' theorem

$$\begin{aligned} p(h|\mathcal{D}) &= \frac{p(\mathcal{D}|h)p(h)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|h)p(h)}{\sum_j p(\mathcal{D}, h_j)} = \frac{p(\mathcal{D}|h)p(h)}{\sum_j p(\mathcal{D}|h_j)p(h_j)} = \\ &= \frac{p(h)\mathbb{I}(\mathcal{D} \in h)/|h|^N}{\sum_j p(h_j)\mathbb{I}(\mathcal{D} \in h_j)/|h_j|^N} \end{aligned}$$

- in simple words: **posterior** \propto **likelihood** \times **prior**

$$p(h|\mathcal{D}) \propto p(\mathcal{D}|h)p(h) = p(h)\mathbb{I}(\mathcal{D} \in h)/|h|^N$$

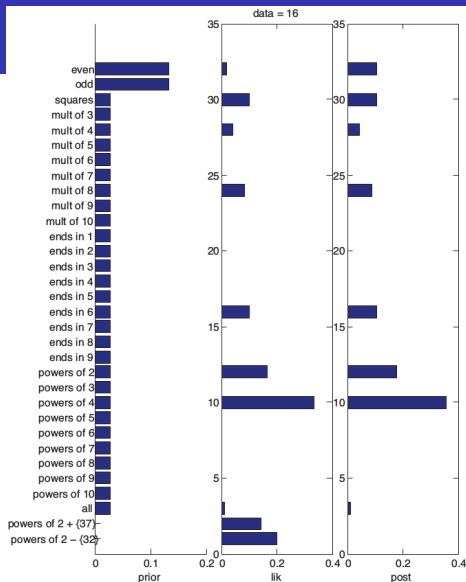
- recall that the denominator $p(\mathcal{D})$ can be seen as a normalization constant which ensure a valid probability on the left-hand side which integrates to one

N.B.: $p(\mathcal{D}, h_j)$ is the probability of observing \mathcal{D} and having h_j equal to the "true" hypothesis representing the concept C

Number Game

Posterior

assume $\mathcal{D} = \{16\}$

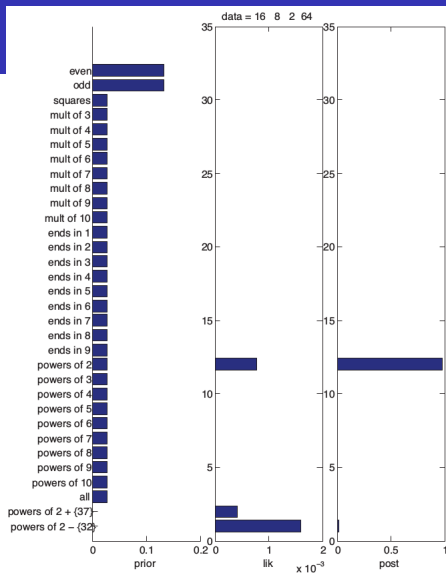


$$p(h) \quad p(\mathcal{D}|h) \quad \propto \quad p(h|\mathcal{D})$$

Number Game

Posterior

assume $\mathcal{D} = \{16, 8, 2, 64\}$



$$p(h) \quad p(\mathcal{D}|h) \quad \propto \quad p(h|\mathcal{D})$$

Number Game

MAP

- in general when we have **enough data** the **posterior** becomes peaked on a single concept, i.e.

$$p(h|\mathcal{D}) \rightarrow \delta_{h_{MAP}}(h)$$

where the **MAP** is

$$h_{MAP} = \arg \max_h p(h|\mathcal{D}) = \arg \max_h p(\mathcal{D}|h)p(h) = \arg \max_h [\log(p(\mathcal{D}|h)) + \log(p(h))]$$

- in this case we can extract and use the **MAP** as a good representative estimate

Number Game

MLE

- note that the prior $p(h)$ stays constant
- since in our case the likelihood depends exponentially on N , i.e.

$$p(\mathcal{D}|h) = \mathbb{I}(\mathcal{D} \in h)/|h|^N$$

we have that the **MAP** estimate converges toward the **MLE**

$$h_{MLE} = \arg \max_h p(\mathcal{D}|h) = \arg \max_h \log(p(\mathcal{D}|h))$$

- when N is large enough, the **data overwhelms the prior**, i.e.

$$\lim_{N \rightarrow \infty} h_{MAP} = h_{MLE}$$

- 1 Bayesian Approach
 - General Paradigma
- 2 Bayesian Concept Learning
 - Number Game
 - Likelihood
 - Prior
 - Posterior
 - Posterior Predictive

Number Game

Posterior Predictive

- the posterior $p(h|\mathcal{D})$ is our internal **belief state** of the world
- we want to test it by predicting observable quantities
- assume a new instance \tilde{x} arrives
- **posterior predictive distribution**

$$\begin{aligned} p(\tilde{x} \in C|\mathcal{D}) &= p(y = 1|\tilde{x}, \mathcal{D}) = \sum_j p(y = 1, h_j|\tilde{x}, \mathcal{D}) = \sum_j p(y = 1|\tilde{x}, h_j, \mathcal{D})p(h_j|\mathcal{D}) \\ &= \sum_j p(y = 1|\tilde{x}, h_j)p(h_j|\mathcal{D}) \end{aligned}$$

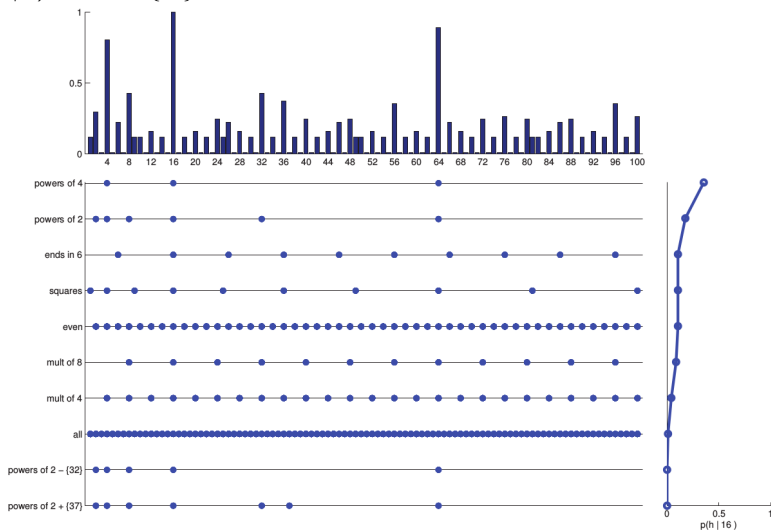
weighted average of the predictions of the individual hypotheses

- this is also called **Bayesian model averaging**
- in this case $p(y = 1|\tilde{x}, h_j) = \mathbb{I}(\tilde{x} \in h_j)$

Number Game

Posterior Predictive

$p(\tilde{x} \in C | \mathcal{D})$ with $\mathcal{D} = \{16\}$



Number Game

Posterior Predictive

- when the dataset is small or ambiguous, the posterior $p(h|\mathcal{D})$ is vague
- this induces a broad predictive distribution
- as noticed, as more data arrives we have that the **posterior** becomes peaked on a single concept, i.e.

$$p(h|\mathcal{D}) \rightarrow \delta_{h_{MAP}}(h)$$

- in this case we can use the **plug-in approximation**

$$p(\tilde{x} \in C|\mathcal{D}) = \sum_j p(y = 1|\tilde{x}, h_j)p(h_j|\mathcal{D}) \simeq \sum_j p(y = 1|\tilde{x}, h_j)\delta_{h_{MAP}}(h_j) =$$

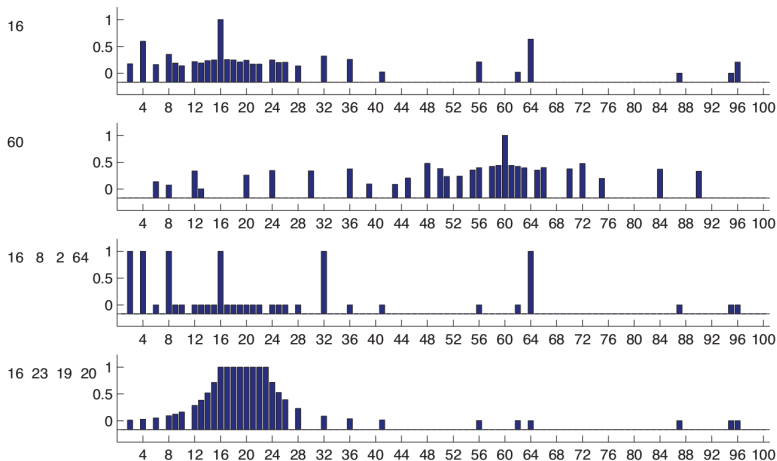
and hence

$$p(\tilde{x} \in C|\mathcal{D}) \simeq p(y = 1|\tilde{x}, h_{MAP})$$

- this approximation can be used at the cost of **under-representing our uncertainty** by losing smooth Bayesian "transitions"

Number Game

Posterior Predictive



- Bayesian approach: we start broad and then narrow down as we learn more
- plug-in approx gets broader or stay the same

- Kevin Murphy's book