

National Technical University of Ukraine  
“Igor Sikorsky Kyiv Polytechnic Institute”  
Institute of Physics and Technology

## **Lecture 3**

# Generative Models for Real-World Systems. Bayesian Concept

Dmytro Progonov,

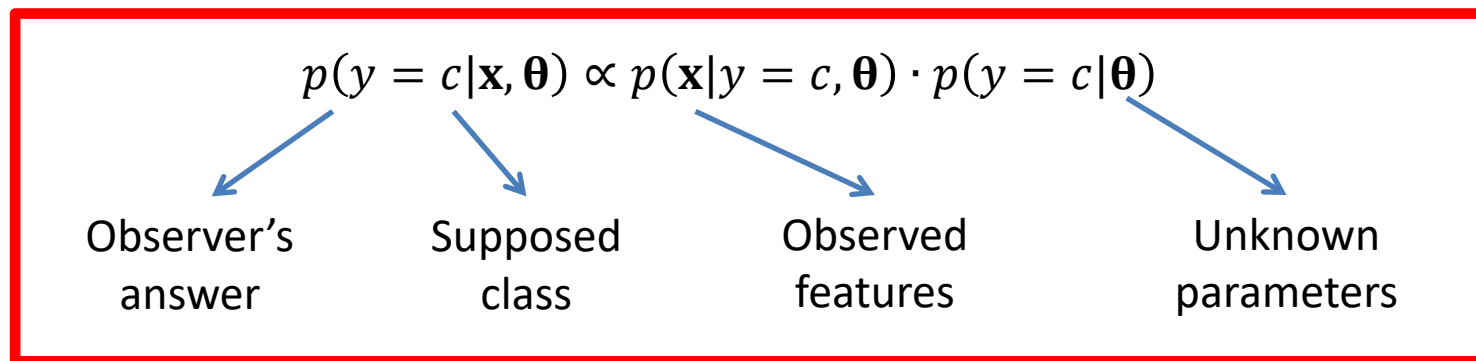
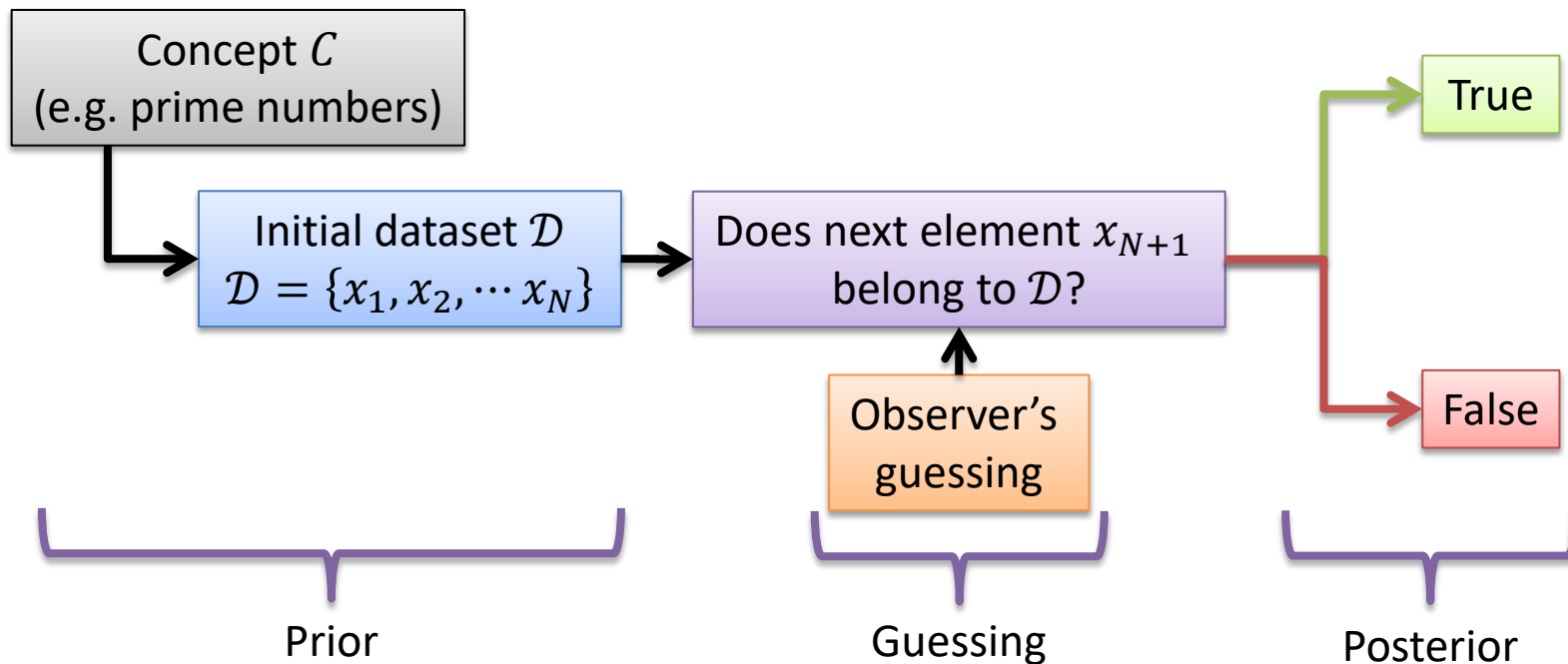
PhD, Associate Professor,

Department Of Physics and Information Security Systems

# Content

- Bayesian concept learning;
- Prior and posterior processing;
- Naïve Bayes classifier;
- Bayesian approach applications.

# Bayesian concept learning



# Prior and posterior processing (1/4)

## Prior

*Likelihood*

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

$\mathcal{D}$  – observed data;

$h$  – total number of items;

$N$  – number of sampled (with replacement) items.

### Occam's razor:

The model favors the simplest (smallest) hypothesis consistent with the data.

### Jeffreys-Linley paradox:

Bayes-oriented decision systems will always favor the simpler model, since the probability of the observed data under a complex model with a very diffuse prior will be very small.

## Posterior

$$p(h|\mathcal{D}) = \frac{p(\mathcal{D}|h)p(h)}{\sum_{\hat{h} \in \mathcal{H}} p(\mathcal{D}, \hat{h})} = \frac{p(h)\mathbb{I}(\mathcal{D} \in h)/|h|^N}{\sum_{\hat{h} \in \mathcal{H}} p(\hat{h})\mathbb{I}(\mathcal{D} \in \hat{h})/|\hat{h}|^N}$$

where  $\mathbb{I}(\mathcal{D} \in h)$  is 1 **if and only if** all the data are in extension of the hypothesis  $h$ .

# Prior and posterior processing (2/4)

If we have enough data, the posterior  $p(h|\mathcal{D})$  becomes **Maximum A Posterior Estimate (MAP)**:

$$p(\mathcal{D}|h) \rightarrow \delta_{\hat{h}^{MAP}}(h),$$

$\hat{h}^{MAP} = \operatorname{argmax}_h p(h|\mathcal{D})$  – the posterior mode;

$$\delta_x(A) = \begin{cases} 1, & x \in A \\ 0, & x \notin A \end{cases} \text{ – Dirac measure.}$$

Note that the MAP estimate can be written as:

$$\hat{h}^{MAP} = \operatorname{argmax}_h [p(\mathcal{D}|h)p(h)] = \operatorname{argmax}_h [\log(p(\mathcal{D}|h)) + \log(p(h))].$$

As we get more and more data, the MAP estimate converges towards the **Maximum Likelihood Estimate (MLE)**:

$$\hat{h}^{MLE} = \operatorname{argmax}_h [p(\mathcal{D}|h)] = \operatorname{argmax}_h [\log(p(\mathcal{D}|h))].$$

# Prior and posterior processing (3/4)

The way to test if our beliefs are justified is to use them to predict objectively observable quantities with usage of **posterior predictive distribution**:

$$p(\tilde{x} \in C|\mathcal{D}) = \sum_h p(y = 1|\tilde{x}, h)p(h|\mathcal{D}).$$

When we have a small and/or ambiguous dataset, the posterior  $p(h|\mathcal{D})$  is vague, which induces a broad predictive distribution. However, once we have “figured things out”, the posterior becomes a delta function centered at the MAP estimate. In this case, we can use **plug-in approximation**:

$$p(\tilde{x} \in C|\mathcal{D}) = \sum_h p(\tilde{x}|h)\delta_{\hat{h}}(h) = p(\tilde{x}|\hat{h}).$$

# Prior and posterior processing (4/4)

In general, computing the Maximum a Posterior Estimate  $p(\mathcal{D}|h)$  can be quite difficult. One simple but popular approximation is known as the **Bayesian information criterion** (**BIC**)

$$BIC \triangleq \log p(\mathcal{D}|\hat{\theta}) - \frac{dof(\hat{\theta})}{2} \log N \approx \log p(\mathcal{D}),$$

where  $\hat{\theta}$  – maximum likelihood estimation of used model parameters;  $dof(\hat{\theta})$  – the number of degrees of freedom in used model.

The BIC method is very closely related to the **Minimum Description Length** or **MDL** principle, which characterizes the score of a model in terms of how well it fits the data, minus how complex the model is to define.

A very similar expression of BIC / MDL is called the **Akaike information criterion** or **AIC**

$$AIC(m, \mathcal{D}) \triangleq \log p(\mathcal{D}|\hat{\theta}_{MLE}) - dof(m).$$

# Examples. Beta-binomial model (1/3)

Suppose  $X_i \sim \text{Ber}(\theta)$ , where  $X_i = 1$  represents “heads”,  $X_i = 0$  represents “tails”, and  $\theta \in [0; 1]$  is the rate parameter (probability of head). If the data are iid, the likelihood has the form:

$$p(\mathcal{D}|\theta) = \theta^{N_1} (1 - \theta)^{N_0}.$$

where we have  $N_1 = \sum_{i=1}^N \mathbb{I}(x_i = 1)$  heads and  $N_0 = \sum_{i=1}^N \mathbb{I}(x_i = 0)$  tails,  $N = N_0 + N_1$  is observed trials. In this case we have  $N_1 \sim \text{Bin}(N, \theta)$ , which has following pdf:

$$\text{Bin}(k|n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}.$$

Since binomial coefficients  $\binom{n}{k}$  is a constant independent of  $\theta$ , the **likelihood of the binomial sampling model** is the same as the likelihood for the Bernoulli model – any inference we have about  $\theta$  will be same whether we observe the counts  $\mathcal{D} = (N_1, N)$  or sequence of trials  $\mathcal{D} = \{x_1, \dots, x_N\}$ .



# Examples. Beta-binomial model (2/3)

To make the math easier, it would be convenient if the prior had the same form as the likelihood:

$$p(\theta) \propto \theta^{\gamma_1}(1 - \theta)^{\gamma_0}$$

for some prior parameters  $\gamma_1$  and  $\gamma_2$ . Then we could easily evaluate the posterior by simply adding the exponents

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta) = \theta^{N_1}(1 - \theta)^{N_0}\theta^{\gamma_1}(1 - \theta)^{\gamma_0} = \theta^{N_1+\gamma_1}(1 - \theta)^{N_0+\gamma_0}$$

When the prior and the posterior have the same form, we say that the prior is a **conjugate prior** for the corresponding likelihood. In case of the Bernoulli, the conjugate prior is the beta distribution:

$$\text{Beta}(\theta|a, b) \propto \theta^{a-1}(1 - \theta)^{b-1}$$

Then posterior is

$$p(\theta|\mathcal{D}) \propto \text{Bin}(N_1|\theta, N_0 + N_1)\text{Beta}(\theta|a, b) \propto \text{Beta}(\theta|N_1 + a, N_0 + b)$$

# Examples. Beta-binomial model (3/3)

Consider predicting the probability of heads in a single future trial under a  $\text{Beta}(a, b)$  posterior:

$$p(\tilde{x} = 1|\mathcal{D}) = \int_0^1 p(x = 1|\theta)p(\theta|\mathcal{D})d\theta = \int_0^1 \theta \text{Beta}(\theta|a, b)d\theta = \mathbb{E}[\theta|\mathcal{D}] = \frac{a}{a+b}.$$

**Zero count (sparse data) problem** – occurs when estimating counts from small amount of data;

**Black swan paradox** – problem of how to draw general conclusion about the future from specific observation from the past.

Suppose now we were interested in predicting the number of heads,  $x$ , in  $M$  future trials:

$$p(x|\mathcal{D}, M) = \int_0^1 \text{Bin}(x|\theta, M) \text{Beta}(\theta|a, b) d\theta = \\ \binom{M}{x} \frac{1}{B(a, b)} \int_0^1 \theta^x (1 - \theta)^{M-x} \theta^{a-1} (1 - \theta)^{b-1} d\theta = \binom{M}{x} \frac{B(x + a, M - x + b)}{B(a, b)}.$$

# Naïve Bayes classifier

Let us classify vector of discrete-valued features  $\mathbf{x} \in \{1, \dots, K\}^D$ , where  $K$  is number of values for each feature,  $D$  is the number of features. The simplest approach is to assume the features are conditionally independent given the class label:

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|y = c, \boldsymbol{\theta}_{jc}).$$

The model is called “naïve” since we do not expect the features to be independent, even conditional on the class label. One reason of successful application of naïve Bayes classifier is that the model is quite simple, and hence it is relatively immune to overfitting.

Type of feature	Recommended type of class-conditional density
Real-valued	$p(\mathbf{x} y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \mathcal{N}(x_j \mu_{jc}, \sigma_{jc}^2)$
Binary	$p(\mathbf{x} y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_j \mu_{jc})$
Categorical	$p(\mathbf{x} y = c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Multinoulli}(x_j \boldsymbol{\mu}_{jc})$

# Bayesian approach applications

**Hierarchical Bayes** (multi-level model) based on putting a prior on used prior:

$$\eta \rightarrow \theta \rightarrow \mathcal{D}.$$

**Empirical Bayes** violates the principle that the prior should be chosen independently of the data:

$$\hat{\eta} = \operatorname{argmax}_{\eta} p(\mathcal{D}|\eta) = \operatorname{argmax}_{\eta} \left[ \int p(\mathcal{D}|\theta)p(\theta|\eta)d\theta \right].$$

Method	Definition
Maximum likelihood	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)$
MAP estimation	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)p(\theta \eta)$
Empirical Bayes	$\hat{\eta} = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)$
Full Bayes	$p(\theta, \eta \mathcal{D}) \propto p(\mathcal{D} \theta)p(\theta \eta)p(\eta)$

# Conclusion

- Concept of Bayesian learning was considered;
- Methods for processing the prior and posterior information were presented;
- Practical applications of Bayesian approach were shown.