

# Bayesian Models in Machine Learning

Generative models

Lukáš Burget



BAYa lectures, November 2023

# Frequentist vs. Bayesian

- Frequentist point of view:
  - Probability is the frequency of an event occurring in a large (infinite) number of trials
  - E.g. When flipping a coin many times, what is the proportion of heads?
- Bayesian
  - Inferring probabilities for events that have never occurred or believes which are not directly observed
  - Prior beliefs are specified in terms of prior probabilities
  - Taking into account uncertainty (posterior distribution) of the estimated parameters or hidden variables in our probabilistic model.

# Simple classification problem – I.

- Simple example of learning a probabilistic model for maximum a-posteriori classification
  - to introduce classification as a basic problem from machine learning field
  - to understand frequentist's view of "probability" and to show its limitations as compared to the Bayesian approaches
  - to refresh basics from probability theory
- The task is to classify an object (*grenade* or *apple*) given an observation (discrete weight category)
  - It is heavy. Is it grenade or apple?
- Let's have 150 observations as training data
  - Table of observation counts for each class and weight category



1	6	12	15	12	2	2	50
4	22	50	14	6	3	1	100
<i>lightest</i> 0.0 - 0.1	<i>lighter</i> 0.1 - 0.2	<i>light</i> 0.2 - 0.3	<i>middle</i> 0.3 - 0.4	<i>heavy</i> 0.4 - 0.5	<i>heavier</i> 0.5 - 0.6	<i>heaviest</i> 0.6 - 0.7	[kg]

# Simple classification problem – II.

- Let's estimate joint probabilities  $P(\text{class}, \text{observation})$ 
  - normalizing the counts by the total count gives **Maximum likelihood (ML) estimates** (see later):  $P(\text{grenade}, \text{heavy}) = \frac{12}{150}$
  - We need many observations to obtain robust estimates this way.
  - How certain can we be about correctness of these estimates?
- Maximum a-posteriori classification rule:
  - given an observation select the most likely class
  - i.e. select class with highest posterior probability  $P(\text{class}|\text{observation})$
  - ML estimate:  $P(\text{grenade}|\text{heavy}) = \frac{12}{12+6}$



$\frac{1}{150}$	$\frac{6}{150}$	$\frac{12}{150}$	$\frac{15}{150}$	$\frac{12}{150}$	$\frac{2}{150}$	$\frac{2}{150}$	$\frac{50}{150}$
$\frac{4}{150}$	$\frac{22}{150}$	$\frac{50}{150}$	$\frac{14}{150}$	$\frac{6}{150}$	$\frac{3}{150}$	$\frac{1}{150}$	$\frac{100}{150}$
<i>lightest</i> 0.0 - 0.1	<i>lighter</i> 0.1 - 0.2	<i>light</i> 0.2 - 0.3	<i>middle</i> 0.3 - 0.4	<i>heavy</i> 0.4 - 0.5	<i>heavier</i> 0.5 - 0.6	<i>heaviest</i> 0.6 - 0.7	[kg]

# Basic rules of probability theory – I.

Sum rule:

$$P(x) = \sum_y P(x, y)$$

Product rule:

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

Bayes rule:

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

# Basic rules of probability theory – II.

- Sum rule:

$$P(\text{heavy}) = P(\text{grenade, heavy}) + P(\text{apple, heavy}) = \frac{12}{150} + \frac{6}{150} = \frac{18}{150}$$

$$P(\text{grenade}) = \sum_x P(\text{grenade, } x) = \frac{50}{150}$$

- Product rule:

$$P(\text{grenade, heavy}) = P(\text{grenade}|\text{heavy})P(\text{heavy}) = \frac{12}{18} \frac{18}{150} = \frac{12}{150}$$

$$P(\text{grenade, heavy}) = P(\text{heavy}|\text{grenade})P(\text{grenade}) = \frac{12}{50} \frac{50}{150} = \frac{12}{150}$$



$\frac{1}{150}$	$\frac{6}{150}$	$\frac{12}{150}$	$\frac{15}{150}$	$\frac{12}{150}$	$\frac{2}{150}$	$\frac{2}{150}$	$\frac{50}{150}$
$\frac{4}{150}$	$\frac{22}{150}$	$\frac{50}{150}$	$\frac{14}{150}$	$\frac{6}{150}$	$\frac{3}{150}$	$\frac{1}{150}$	$\frac{100}{150}$
<i>lightest</i> 0.0 - 0.1	<i>lighter</i> 0.1 - 0.2	<i>light</i> 0.2 - 0.3	<i>middle</i> 0.3 - 0.4	<i>heavy</i> 0.4 - 0.5	<i>heavier</i> 0.5 - 0.6	<i>heaviest</i> 0.6 - 0.7	[kg]

# Basic rules of probability theory – III.

- Bayes rule:

The diagram shows the Bayes' rule formula with callouts for its components. A light blue box labeled 'Posterior probability' points to the left side of the equation. A light blue box labeled 'Likelihood' points to the numerator. A light blue box labeled 'Prior probability' points to the denominator. A light blue box labeled 'Evidence' points to the denominator.

$$P(\text{grenade}|\text{heavy}) = \frac{P(\text{heavy}|\text{grenade})P(\text{grenade})}{P(\text{heavy})}$$

- The evidence can be evaluated using the sum and product rules in terms of likelihoods and priors:

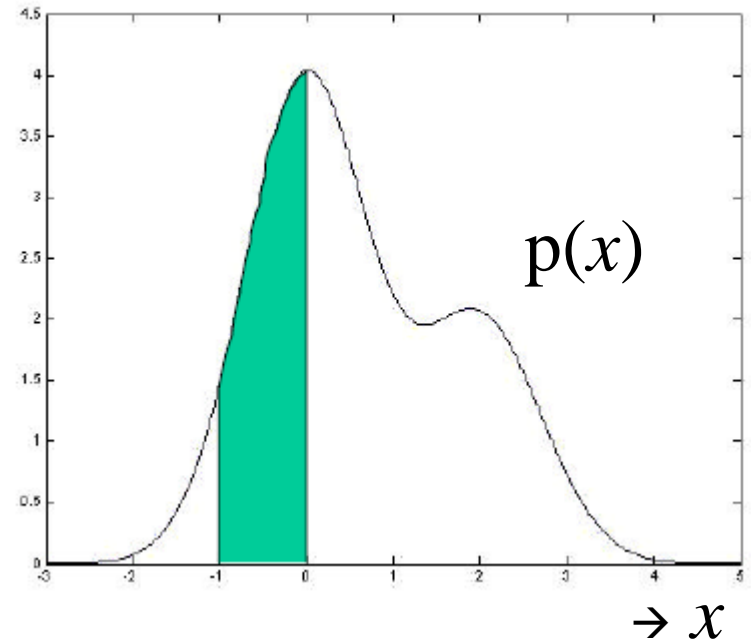
$$P(\text{heavy}) = P(\text{heavy}|\text{grenade})P(\text{grenade}) + P(\text{heavy}|\text{apple})P(\text{apple})$$

- Bayes rule for calculating the class posterior may not seem very useful now, but it will be useful in case continuous valued observations.

# Continuous random variables

- $P(x)$  –probability
- $p(x)$  –probability density function

$$P(x \in (a, b)) = \int_a^b p(x) dx$$

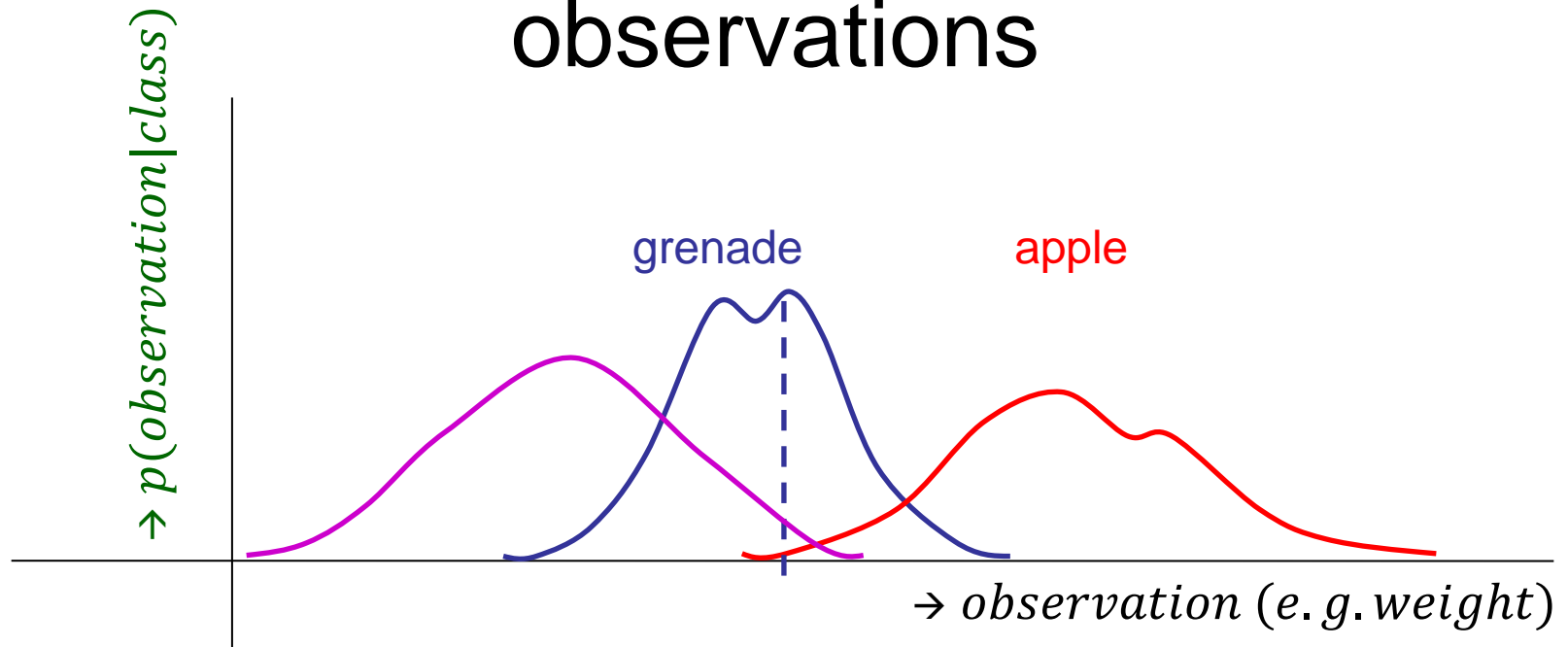


Sum rule:

$$p(x) = \int p(x, y) dy$$



# Classification with continuous observations



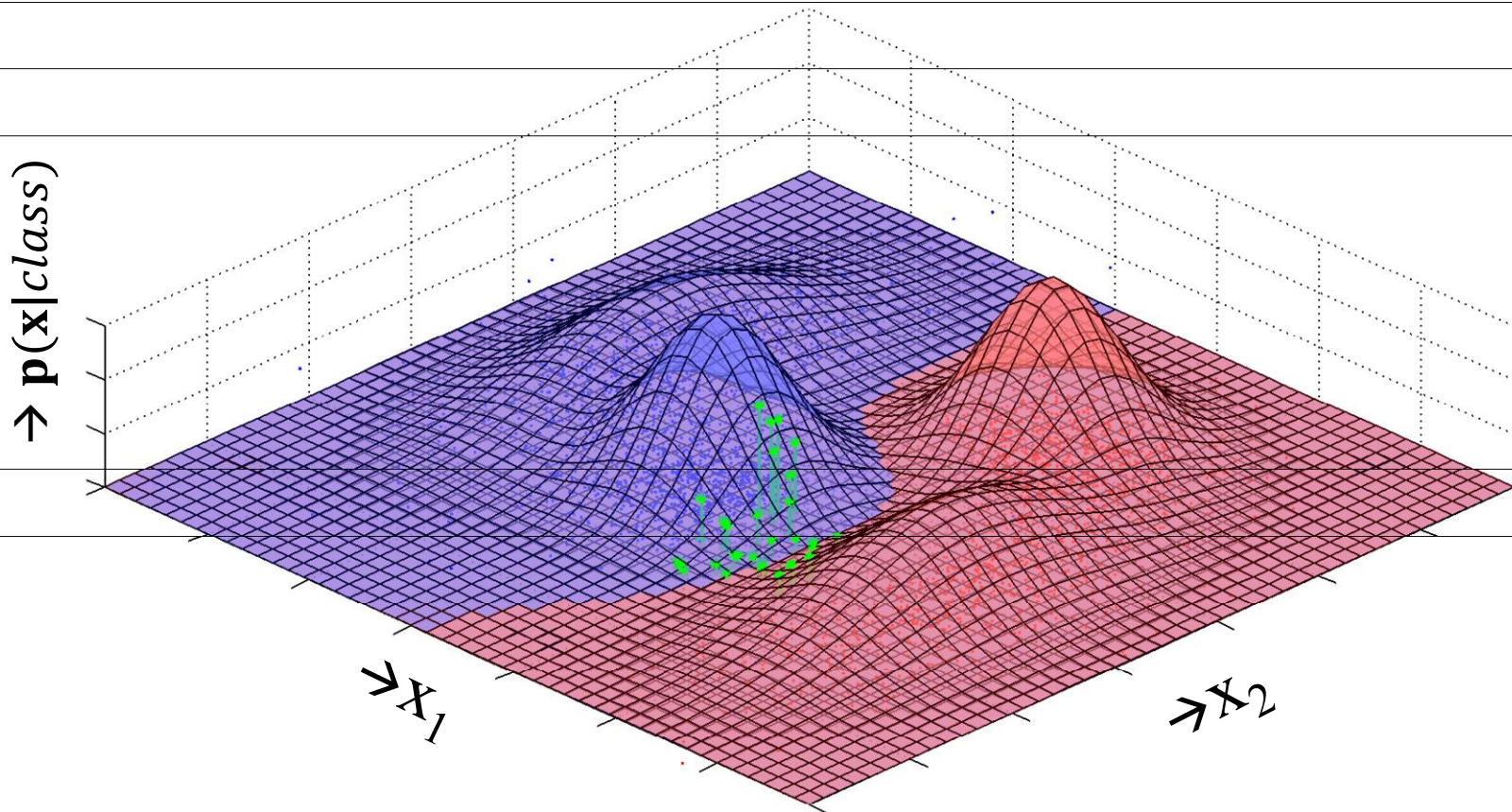
- Maximum a-posteriori classification rule says: select the more likely class

$$P(\text{class}|\text{observation}) = \frac{p(\text{observation}|\text{class})P(\text{class})}{p(\text{observation})}$$

$$P(\text{observation}) = \sum_{\text{class}} p(\text{observation}|\text{class})P(\text{class})$$

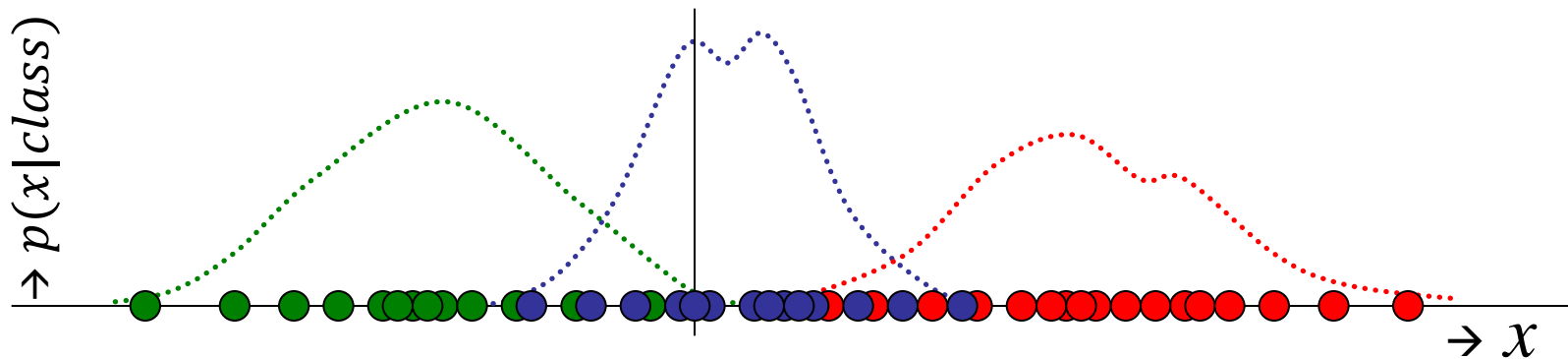
# Multivariate observations

From now, univariate observations will be denoted as  $x$  and multivariate as  $\mathbf{x} = [x_1, x_2, \dots, x_D] = [\textit{weight}, \textit{diameter}, \dots]$



# Estimation of parameters

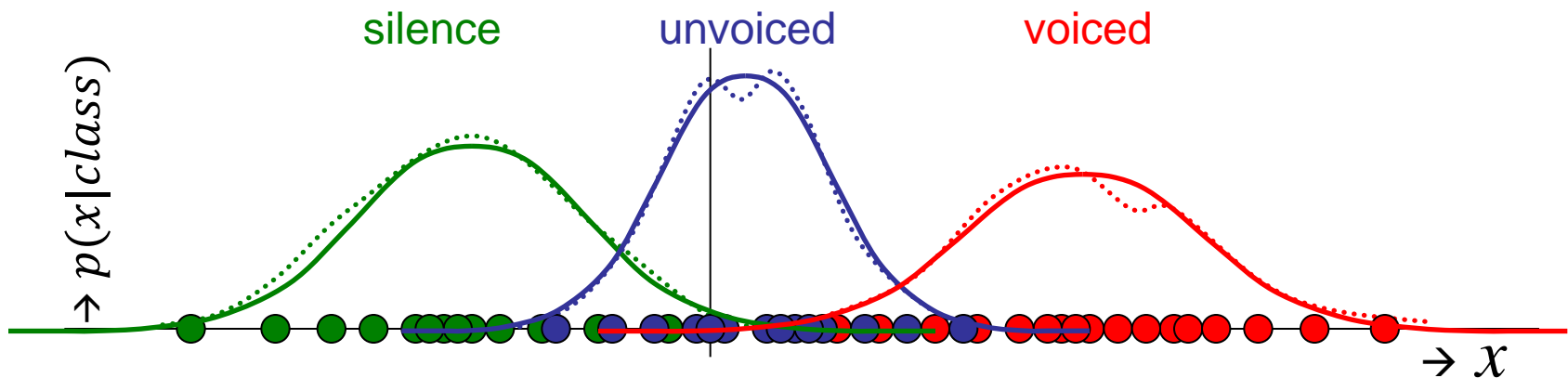
- Usually, we do not know the true distributions  $p(x|class)$



# Estimation of parameters

... we only see some training examples.

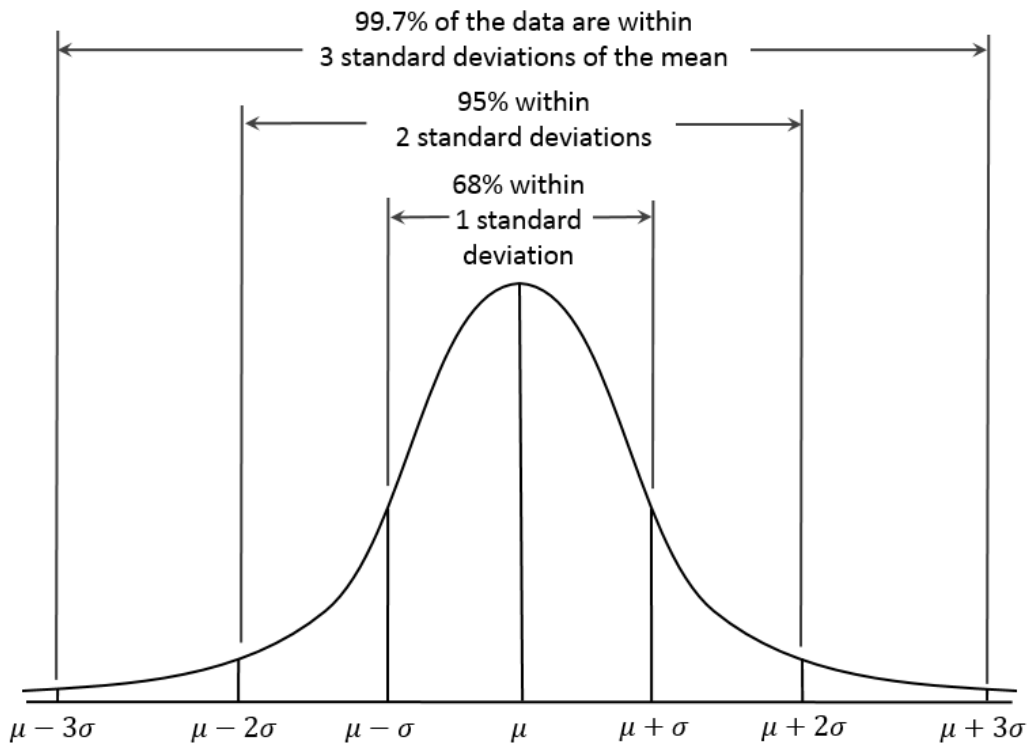
- Let's decide for some parametric model for  $p(x|class)$  (e.g. Gaussian distribution) and estimate its parameters from the data.



- Here, we are using the **frequentist approach**: Estimated distributions tell us that observation  $x$  will be more likely as we see more similar observations in the training data.
- From now, let's forget about classes. We will concentrate just on estimating probability density functions (e.g. one for each class).

# Gaussian distribution (univariate)

$$p(x) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



## ML estimates of parameters

$$\mu = \frac{1}{N} \sum_n x_n$$

$$\sigma^2 = \frac{1}{N} \sum_n (x_n - \mu)^2$$

# Why Gaussian distribution?

- Simple and easy to deal with
  - Just a quadratic function in log domain

$$\log \mathcal{N}(x; \mu, \sigma^2) = -\frac{\log(2\pi\sigma^2)}{2} - \frac{1}{2\sigma^2}(x - \mu)^2 = -\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2} + K$$

- Log likelihood of observed sequence  $\mathbf{x} = [x_1, x_2, x_3, \dots, x_N]$  is

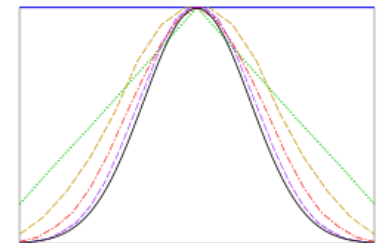
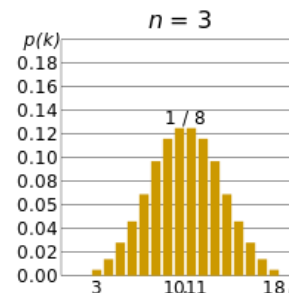
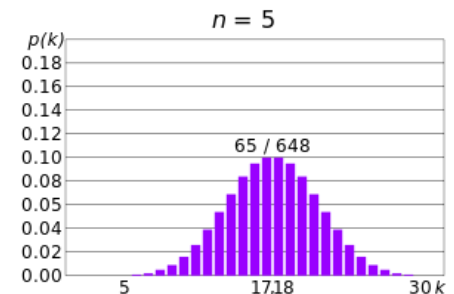
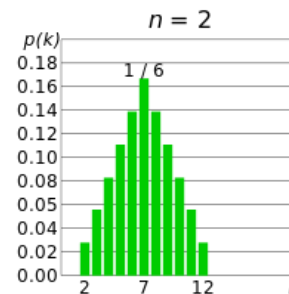
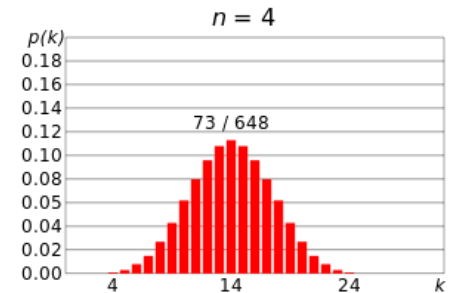
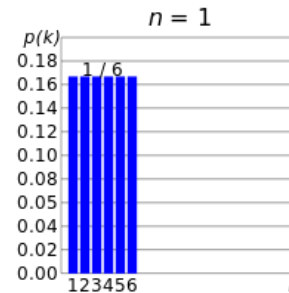
$$\log p(\mathbf{x}|\mu, \sigma^2) = \log \prod_n \mathcal{N}(x_n; \mu, \sigma^2) = \sum_n \log \mathcal{N}(x_n; \mu, \sigma^2)$$

$$= -\frac{1}{2\sigma^2} \sum_n x_n^2 + \frac{\mu}{\sigma^2} \sum_n x_n - N \left( \frac{\mu^2}{2\sigma^2} + K \right)$$

Sufficient statistics  
(second, first and zero order)

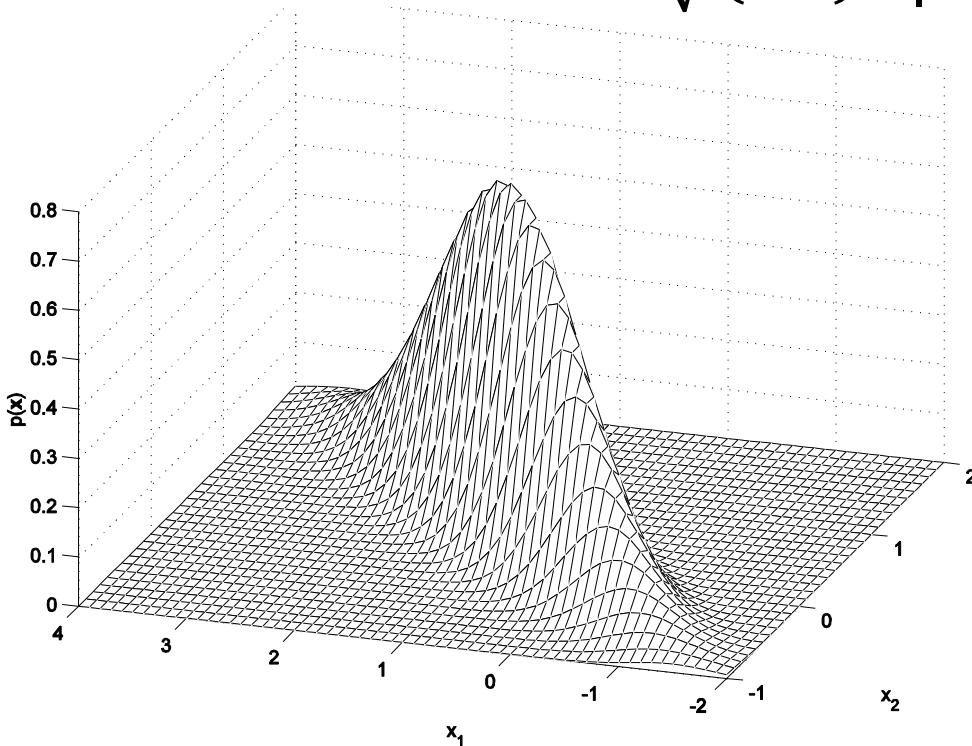
# Why Gaussian distribution?

- Naturally occurring
- Central limit theorem: Summing values of many independently generated random variables gives Gaussian distributed observations
- Examples:
  - Summing outcome of N dices
  - Galton's board  
<https://www.youtube.com/watch?v=03tx4v0i7MA>



# Gaussian distribution (multivariate)

$$p(x_1, \dots, x_D) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$



**ML estimates of parameters**

$$\boldsymbol{\mu} = \frac{1}{N} \sum_n \mathbf{x}_n$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_n (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T$$



# Maximum likelihood estimation of parameters

- Let's choose a parametric distribution  $p(\mathbf{x}|\boldsymbol{\eta})$  with parameters  $\boldsymbol{\eta}$ 
  - Gaussian distribution with parameters  $\mu, \sigma^2$
- ... and let's have some observed training data  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ , which we assume to be i.i.d. generated from this distribution.
- We might obtain maximum likelihood estimates of the parameters  $\hat{\boldsymbol{\eta}}^{ML}$  by maximizing the likelihood of the observed data

$$\hat{\boldsymbol{\eta}}^{ML} = \arg \max_{\boldsymbol{\eta}} p(\mathbf{X}|\boldsymbol{\eta}) = \arg \max_{\boldsymbol{\eta}} \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\eta})$$

- Later, we will see that, under some assumptions, this estimates gives us the most likely parameters.

# ML estimate for Gaussian

$$\begin{aligned}\arg \max_{\mu, \sigma^2} p(\mathbf{x}|\mu, \sigma^2) &= \arg \max_{\mu, \sigma^2} \log p(\mathbf{x}|\mu, \sigma^2) = \arg \max_{\mu, \sigma^2} \sum_n \log \mathcal{N}(x_n; \mu, \sigma^2) \\ &= \arg \max_{\mu, \sigma^2} \left( -\frac{1}{2\sigma^2} \sum_n x_n^2 + \frac{\mu}{\sigma^2} \sum_n x_n - N \frac{\mu^2}{2\sigma^2} - N \frac{\log(2\pi\sigma^2)}{2} \right)\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \mu} \log p(\mathbf{x}|\mu, \sigma^2) &= \frac{\partial}{\partial \mu} \left( -\frac{1}{2\sigma^2} \sum_n x_n^2 + \frac{\mu}{\sigma^2} \sum_n x_n - N \frac{\mu^2}{2\sigma^2} - N \frac{\log(2\pi\sigma^2)}{2} \right) \\ &= \frac{1}{\sigma^2} \left( \sum_n x_n - N\mu \right) = 0 \quad \Rightarrow \quad \hat{\mu}^{ML} = \frac{1}{N} \sum_n x_n\end{aligned}$$

and similarly:  $\hat{\sigma}^{2, ML} = \frac{1}{N} \sum_n (x_n - \mu)^2$

# Categorical distribution



4	22	50	14	6	3	1	100
<i>lightest</i>	<i>lighter</i>	<i>light</i>	<i>middle</i>	<i>heavy</i>	<i>heavier</i>	<i>heaviest</i>	
0.0 - 0.1	0.1 - 0.2	0.2 - 0.3	0.3 - 0.4	0.4 - 0.5	0.5 - 0.6	0.6 - 0.7	[kg]

$$p(x|\boldsymbol{\pi}) = \text{Cat}(x|\boldsymbol{\pi}) = \pi_x$$

- Also referred to as **Discrete distribution**
- Special binary case is **Bernoulli distribution**
- $x \in \{\textit{lightest}, \textit{lighter}, \textit{light}, \textit{middle}, \textit{heavy}, \textit{heavier}, \textit{heaviest}\}$   
or  $x$  can be simply the index of a category  $x \in \{1, 2, \dots, C\}$
- $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_C]$  - probabilities of the categories are the parameters
- Likelihood of an observed training set  $\mathbf{x} = [x_1, x_2, \dots, x_N]$

$$P(\mathbf{x}|\boldsymbol{\pi}) = \prod_n \text{Cat}(\mathbf{x}_n|\boldsymbol{\pi}) = \prod_n \pi_{x_n} = \prod_c \pi_c^{m_c}$$

where  $m_c$  is number of observations from category  $c$ .

- (e.g. the numbers from the table)

# ML estimate for Categorical

$$\begin{aligned}\arg \max_{\boldsymbol{\pi}} p(\mathbf{x}|\boldsymbol{\pi}) &= \arg \max_{\boldsymbol{\pi}} \log p(\mathbf{x}|\boldsymbol{\pi}) = \arg \max_{\boldsymbol{\pi}} \log \prod_{n=1}^N \text{Cat}(x_n|\boldsymbol{\pi}) \\ &= \arg \max_{\boldsymbol{\pi}} \log \prod_c \pi_c^{m_c} = \arg \max_{\boldsymbol{\pi}} \sum_c m_c \log \pi_c\end{aligned}$$

We need to use Lagrange multiplier  $\lambda$  to enforce the constraint  $\sum_k \pi_k = 1$

$$\frac{\partial}{\partial \pi_c} \left( \sum_k m_k \log \pi_k - \lambda \left( \sum_k \pi_k - 1 \right) \right) = \frac{m_c}{\pi_c} - \lambda = 0$$

$$\Rightarrow \pi_c = \frac{m_c}{\lambda} = \frac{m_c}{\sum_k m_k} = \frac{m_c}{N}$$

