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

# Lecture 5
## Generative Models for Discrete Data. Classical Probabilistic Models

Dmytro Progonov,

PhD, Associate Professor,

Department Of Physics and Information Security Systems

# Content

- Discrete distributions;

- Continuous distributions;

- Joint probability distributions;

- Characteristics of distributions.

Generative Models for Discrete Data.
Classical Probabilistic Models

2/22

# Common distributions

Discrete distributions
- Binomial and Bernoulli distributions
- Multinomial and multinoulli distributions
- Poisson distribution
- Empirical distribution

Continuous distribution
- Gaussian (normal) distribution
- Degenerate distribution
- Student's t distribution
- Laplace distribution
- Gamma distribution
- Beta distribution
- Pareto distribution

Joint probability distribution
- Multivariate Gaussian
- Multivariate Student t distribution
- Dirichlet distribution

Generative Models for Discrete Data.
Classical Probabilistic Models

# Discrete distributions (1/4).
# Binomial and Bernoulli distributions

Toss a coin $n$ times. Let $X \in \{0,1,2 \cdots n\}$ be the number of heads. If the probability of head is $\theta$ then $X$ has ***binomial distribution***:

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

where

$$\binom{n}{k} \triangleq \frac{n!}{(n-k)!\,k!}$$

is binomial coefficient (number of ways to choose $k$ items from $n$).

Let $X \in \{0,1\}$ be a binary random variable (success or fail). Then $X$ has ***Bernoulli distribution***:

$$X \sim Ber(x|\theta) = \theta^{\mathbb{I}(x=1)} (1-\theta)^{\mathbb{I}(x=0)},$$

where

$$\mathbb{I}(x) = \begin{cases} 1 & x = true; \\ 0 & x = false; \end{cases}$$

is indicator function.

Generative Models for Discrete Data.
Classical Probabilistic Models

4/22

# Discrete distributions (2/4).
# Multinomial and multinulli distributions

Outcomes of tossing K-sided die. Then random vector
$\mathbf{x} = (x_1, x_2 \cdots x_K)$ has ***multinomial distribution***:

$$Mu(\mathbf{x}|n, \theta) = \binom{n}{x_1 \cdots x_K} \prod_{j=1}^{K} \theta_j^{x_j}$$

where

$$\binom{n}{x_1 \cdots x_K} \triangleq \frac{n!}{x_1! \, x_2! \cdots x_K!}$$

is multinomial coefficient.

Let $\mathbf{x} = \big(\mathbb{I}(x = 1), \cdots, \mathbb{I}(x = K)\big)$ be a binary random vector. Then $\mathbf{x}$ has ***Multinoulli distribution***:

$$Mu(x|1, \theta) = \prod_{j=1}^{K} \theta_j^{\mathbb{I}(x_j = 1)}$$

Generative Models for Discrete Data.
Classical Probabilistic Models

5/22

# Discrete distributions (3/4).
# Poisson distribution

The $X \in \{0,1,2 \cdots\}$ has **_Poisson distribution_** with parameter $\lambda > 0$ if its pmf is:

$$X \sim Pos(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

The Poisson distribution is often used as a model for counts of rare events like radioactive decay and traffic accidents.

pmf – probability mass function.

Generative Models for Discrete Data.
Classical Probabilistic Models

6/22

# Discrete distributions (4/4).
# Empirical distribution

Given a set of data, $\mathcal{D} = \{x_1, \cdots, x_N\}$, the **_empirical distribution_** is:

$$p_{emp}(A) \triangleq \frac{1}{N}\sum_{i=1}^{N}\delta_{x_i}(A),$$

where $\delta_x(A) -$ Dirac measure, defined by

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

Generative Models for Discrete Data.
Classical Probabilistic Models

7/22

# Continuous distributions (1/8).
# Gaussian (normal) distribution

Probability distribution function for ***Gaussian distribution*** is given by

$$X \sim \mathcal{N}(x|\mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

where $\mu = \mathbb{E}[X]$ – the mean; $\sigma^2 = var[X]$ – variance; $\sqrt{2\pi\sigma^2}$ – the normalization constant.

If $p(X = x) = \mathcal{N}(x|\mu, \sigma^2) = \mathcal{N}(0,1)$ then $X$ has ***standard normal distribution***.

The ***central limit theorem*** tells us that sums of independent random variables have approximately Gaussian distribution.

Generative Models for Discrete Data.
Classical Probabilistic Models

8/22

# Continuous distributions (2/8).
# Degenerate distribution

In the limiting case where $\sigma^2 \to 0$, the Gaussian becomes an infinitely tall and infinitely thin "spike" centered at $\mu$:

$$\lim_{\sigma^2 \to 0} \mathcal{N}(x|\mu, \sigma^2) = \delta(x - \mu),$$

where $\delta(\cdot) -$ Dirac delta function, defined as

$$\delta(x) = \begin{cases} +\infty, & x = 0; \\ 0, & x \neq 0; \end{cases}$$

such that

$$\int\limits_{-\infty}^{+\infty} \delta(x) dx = 1.$$

Generative Models for Discrete Data.
Classical Probabilistic Models

9/22

# Continuous distributions (3/8). Student's $t$ distribution

Distribution robust to outliers is the **_Student distribution_** with probability distribution function :

$$\mathcal{T}(x|\mu, \sigma^2, v) \propto \left[1 + \frac{1}{v}\left(\frac{x - \mu}{\sigma}\right)^2\right]^{-\frac{v+1}{2}},$$

where $\mu = \mathbb{E}[X] -$ the mean; $\sigma^2 > 0 -$ scale parameter; $v > 0 -$ degrees of freedom.

To ensure finite variance $v > 2$; for $v \gg 5$ Student distribution rapidly approaches Gaussian distribution and loses its robustness properties.

If $v = 1$ distribution is known as **_Cauchy_** or **_Lorentz distribution_**.

Generative Models for Discrete Data.
Classical Probabilistic Models

10/22

# Continuous distributions (4/8).
# Laplace distribution

***Laplace distribution*** or ***double sided exponential distribution*** has probability distribution function:

$$Lap(x|\mu, b) \triangleq \frac{1}{2b} e^{-\frac{|x-\mu|}{b}},$$

where $\mu -$ location parameter; $b$ $(b > 0) -$ scale parameter.

Laplace distribution is more robust to outlier than Gaussian distribution and puts more probability density at 0, which is useful way to encourage sparsity in a model.

Generative Models for Discrete Data.
Classical Probabilistic Models

11/22

# Continuous distributions (5/8).
# Gamma distribution

***Gamma distribution*** is defined in terms of the shape $a > 0$ and the rate $b > 0$:

$$Ga(x|a, b) \triangleq \frac{b^a}{\Gamma(a)} x^{a-1} e^{-xb}$$

where $\Gamma(a) -$ gamma function:

$$\Gamma(a) \triangleq \int_0^{+\infty} u^{x-1} e^{-u} du.$$

***Inverse gamma distribution*** defined by

$$IG(x|a, b) \triangleq \frac{b^a}{\Gamma(a)} x^{-(a+1)} e^{-b/x}.$$

Generative Models for Discrete Data.
Classical Probabilistic Models

12/22

# Continuous distributions (6/8).
# Special cases of gamma distribution

***Exponential distribution*** is defined by $Expon(x|\lambda) \triangleq Ga(x|1,\lambda)$. Distribution describes the times between Poisson process, i.e. a process in which events occur continuously and independently at a constant average rate $\lambda$.

***Erlang distribution*** is the same as the Gamma distribution where $a$ is integer $- Erlang(x|\lambda) = Ga(x|a,\lambda), a \in \mathbb{Z}$.

***Chi-squared distribution*** is the sum of squared Gaussian random variables and is defined as $\chi^2(x|v) = Ga\left(x|\frac{v}{2},\frac{1}{2}\right)$.

Generative Models for Discrete Data.
Classical Probabilistic Models

13/22

# Continuous distributions (7/8).
# Beta distribution

**_Beta distribution_** has support over the interval $[0; 1]$ and is defined as follow:

$$Beta(x|a,b) \triangleq \frac{1}{B(a,b)} x^{a-1}(1-x)^{b-1}.$$

where $B(p,q) -$ beta function:

$$B(p,q) \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$

Generative Models for Discrete Data.
Classical Probabilistic Models

14/22

# Continuous distributions (8/8).
# Pareto distribution

The **_Pareto distribution_** is defined as follow:

$$Pareto(x|k,m) \triangleq km^k x^{-(k+1)} \mathbb{I}(x \geq m)$$

As $k \to +\infty$, the distribution approaches $\delta(x - m)$. Distribution has the long (heavy) tails and it is widely used for modeling the power-low dependencies.

Generative Models for Discrete Data.
Classical Probabilistic Models

15/22

# Joint probability distributions (1/3). Multivariate Gaussian

The ***multivariate Gaussian (normal) distribution*** in $D$ dimensions defined as:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right],$$

where $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] \in \mathbb{R}^D$ − the mean vector; $\boldsymbol{\Sigma} = cov[\mathbf{x}]$ − the $D \times D$ covariance matrix.

*Precision (concentration) matrix* is just the inverse covariance matrix:
$$\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}.$$

Generative Models for Discrete Data.
Classical Probabilistic Models

16/22

# Joint probability distributions (2/3). Multivariate Student distribution

The ***multivariate Student's t distribution*** in $\underline{D}$ dimensions defined as:

$$\mathcal{T}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, v) \triangleq \frac{\Gamma(v/2 + D/2)}{\Gamma(v/2)} \cdot \frac{|\boldsymbol{\Sigma}|^{-1/2}}{v^{D/2}\pi^{D/2}} \cdot \left[1 + \frac{1}{v}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]^{-\frac{v+D}{2}}$$

where $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] \in \mathbb{R}^D$ — the mean vector; $\boldsymbol{\Sigma}$ — the scale matrix.

Generative Models for Discrete Data.
Classical Probabilistic Models

17/22

# Joint probability distributions (3/3). Dirichlet distribution

The **_Dirichlet distribution_** is generalization of beta distribution, defined by:

$$Dir(\mathbf{x}|\boldsymbol{\alpha}) \triangleq \frac{1}{B(\boldsymbol{\alpha})} \prod_{d=1}^{D} x_d{}^{\alpha_d-1} \mathbb{I}(\mathbf{x} \in S_K)$$

where

$$S_D = \left\{ x: 0 \le x_d \le 1, \sum_{d=1}^{D} x_d = 1 \right\}$$

is probability simplex;

$$B(\boldsymbol{\alpha}) \triangleq \frac{\prod_{d=1}^{D} \Gamma(\alpha_d)}{\Gamma(\alpha_0)}, \alpha_0 \triangleq \sum_{d=1}^{D} \alpha_d,$$

is natural generalization of the beta function to $D$ variable.

Generative Models for Discrete Data.
Classical Probabilistic Models

18/22

# Characteristics of
# discrete and continuous distributions

| Distribution | Mean | Variance | Skewness | Kurtosis |
|---|---|---|---|---|
| Discrete distribution | | | | |
| Binomial $Bin(k\|n,\theta)$ | $np$ | $np(1-p)$ | $\dfrac{1-2p}{\sqrt{np(1-p)}}$ | $\dfrac{1-6p(1-p)}{np(1-p)}$ |
| Bernoulli $Ber(x\|\theta)$ | $p$ | $p(1-p)$ | $\dfrac{1-2p}{\sqrt{p(1-p)}}$ | $\dfrac{1-6p(1-p)}{p(1-p)}$ |
| Multinomial $Mu(\mathbf{x}\|n,\theta)$ | $np_i$ | $np_i(1-p_i)$ | $-$ | $-$ |
| Multinoulli $Mu(\mathbf{x}\|1,\theta)$ | $\mathbf{p}$ | $\mathbf{\Sigma}_{ij} = \begin{cases} p_i(1-p_i), i = j \\ -p_ip_j, i \neq j \end{cases}$ | $-$ | $-$ |
| Poisson $Pos(x\|\lambda)$ | $\lambda$ | $\lambda$ | $\lambda^{-1/2}$ | $\lambda^{-1}$ |

Generative Models for Discrete Data.
Classical Probabilistic Models

19/22

# Characteristics of discrete and continuous distributions

| Distribution | Mean | Variance | Skewness | Kurtosis |
|---|---|---|---|---|
| Continuous distribution | | | | |
| Gaussian $\mathcal{N}(x\|\mu,\sigma^2)$ | $\mu$ | $\sigma^2$ | $0$ | $0$ |
| Student's $\mathcal{T}(x\|\mu,\sigma^2,v)$ | $0\ (v>1)$ | $\begin{cases}\dfrac{v}{v-2}, v>2 \\ +\infty, v\in(1;2]\end{cases}$ | $0\ (v>3)$ | $\begin{cases}\dfrac{6}{v-4}, v>4 \\ +\infty, v\in(2;4]\end{cases}$ |
| Laplace $Lap(x\|\mu,b)$ | $\mu$ | $2b^2$ | $0$ | $3$ |
| Gamma $Ga(x\|a,b)$ | $a/b$ | $a/b^2$ | $2/\sqrt{a}$ | $6/a$ |
| Beta $Beta(x\|a,b)$ | $\dfrac{a}{a+b}$ | $\dfrac{ab}{(a+b)^2(a+b+1)}$ | $\dfrac{2(b-a)\sqrt{a+b+1}}{(a+b+2)\sqrt{ab}}$ | $\dfrac{6[(a-b)^2(a+b+1)-ab(a+b+2)]}{ab(a+b+2)(a+b+3)}$ |
| Pareto $Pareto(x\|k,m)$ | $\begin{cases}+\infty, k\le 1 \\ \dfrac{km}{k-1}, k>1\end{cases}$ | $\begin{cases}+\infty, k\in(0;2] \\ \dfrac{km^2}{(k-1)^2(k-2)}, k>2\end{cases}$ | $\dfrac{2(k+1)}{k-3}\sqrt{\dfrac{k-2}{k}}, k>3$ | $\dfrac{6(k^3+k^2-6k-2)}{k(k-3)(k-4)}, k>4$ |

Generative Models for Discrete Data. Classical Probabilistic Models

20/22

# Characteristics of
# discrete and continuous distributions

| Distribution | Mean | Variance | Skewness | Kurtosis |
|---|---|---|---|---|
| Joint distribution | | | | |
| Multivariate Gaussian $\mathcal{N}(\mathbf{x}\|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ | $\boldsymbol{\mu}$ | $\boldsymbol{\Sigma}$ | – | – |
| Multivariate Student's $\mathcal{T}(\mathbf{x}\|\boldsymbol{\mu}, \boldsymbol{\Sigma}, v)$ | $\boldsymbol{\mu}, v > 1$ | $\dfrac{v}{v-2}\boldsymbol{\Sigma}, v > 2$ | – | – |
| Dirichlet $Dir(\mathbf{x}\|\boldsymbol{\alpha})$ | $\dfrac{\alpha_i}{\sum_k \alpha_k}$ | $\dfrac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0{}^2(\alpha_0 + 1)}, \alpha_0 = \sum_k \alpha_k$ | – | – |

Generative Models for Discrete Data.
Classical Probabilistic Models

21/22

# Conclusion

- Common discrete, continuous and joint probability distributions were considered;

- Key parameters of these distributions were presented.

Generative Models for Discrete Data.
Classical Probabilistic Models

22/22