



ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE  
Fakulta jaderná a fyzikálně inženýrská



# **Generativní modely dat popsaných stromovou strukturou**

## **Generative models of tree structured data**

Bakalářská práce

Autor: **Jakub Bureš**  
Vedoucí práce: **Doc. Ing. Václav Šmídl, Ph.D.**  
Konzultant: **Doc. Ing. Tomáš Pevný, Ph.D.**  
Akademický rok: 2019/2020

1. Seznamte se s popisem dat pomocí stromové struktury. Zvláštní pozornost věnujte metodám více instančního učení (multiple instance learning). Seznamte se s konceptem vnořeného prostoru (embedded space) a jeho reprezentace pomocí neuronových sítí.
2. Seznamte se se základními generativními modely dat popsaných vektorem příznaků. Zvláštní pozornost věnujte metodám typu autoencoder a jejich variační formě. Demonstrujte vlastnosti modelů na jednoduchých příkladech. V maximální míře využijte dostupné knihovny pro generativní modely.
3. Navrhněte několik příkladů typů dat se stromovou strukturou a pro každý z nich navrhněte generativní model. Navrhněte algoritmus pro určení jeho parametrů z dat a diskutujte vhodnost jednotlivých architektur neuronových sítí.
4. Seznamte se s různými druhy apriorních rozložení používaných na latentní proměnné autoencoderu. Odvoďte algoritmy odhadu jejich parametrů a srovnajte jejich výsledky se základním modelem. Diskutujte výsledné odhady.
5. Vyvinutou metodu aplikujte na vhodně zvolená reálná data a diskutujte vliv zvoleného apriorního rozložení na výsledky.

- Zadání práce (zadní strana) -

*Poděkování:*

Chtěl bych zde poděkovat především svému školiteli panu Doc. Ing. Václavu Šmídlovi, Ph.D. za pečlivost, ochotu, vstřícnost a odborné i lidské zázemí při vedení mé bakalářské práce. Dále děkuji svému konzultantovi panu Doc. Ing. Tomáši Pevnému, Ph.D.

*Čestné prohlášení:*

Prohlašuji, že jsem tuto práci vypracoval samostatně a uvedl jsem všechnu použitou literaturu.

V Praze dne 7. července 2020

Jakub Bureš

## Generativní modely dat popsaných stromovou strukturou

*Obor: Matematické inženýrství*

*Druh práce:* Bakalářská práce

Konzultant: Doc. Ing. Tomáš Pevný, Ph.D.  
Katedra počítačů FEL ČVUT Praha Technická 1902/2 166 27 Praha 6 - Dejvice

**Klíčová slova:** klíčová slova (nebo výrazy) seřazená podle abecedy a oddělená čárkou

## Generative models of tree structured data

[illegible]

**Key words:** keywords in alphabetical order separated by commas

# Obsah

<b>1</b>	<b>Teorie</b>	<b>7</b>
1.1	Optimalizace . . . . .	7
1.1.1	Gradient Descent . . . . .	7
1.1.2	Metoda nejmenších čtverců . . . . .	8
1.2	Úvod do pravděpodobnosti a Bayesovská statistika . . . . .	9
1.2.1	Hustoty pravděpodobnosti . . . . .	10
1.2.2	Bayesovská metoda nejmenších čtverců . . . . .	12
1.2.3	Divergence . . . . .	13
1.2.4	ELBO . . . . .	14
1.3	Teorie grafů . . . . .	16
<b>2</b>	<b>Generativní modely</b>	<b>17</b>
2.1	Variační autoencoder . . . . .	18
2.1.1	Naivní přístup . . . . .	19
2.1.2	Variační Bayseova metoda . . . . .	19
<b>3</b>	<b>Stromové struktury</b>	<b>20</b>
	<b>Závěr</b>	<b>21</b>

# Kapitola 1

## Teorie

### 1.1 Optimalizace

Předpokládejme že máme trénovací data obsahující  $n$  pozorování  $x$ , nebo-li  $\mathbf{X} = (x_1, \dots, x_n)$ . Dále máme ke každému  $x$  právě jedno pozorování  $t$ , psáno  $\mathbf{t} = (t_1, \dots, t_n)$ , komplexně zapsáno zobrazením  $(x_1, \dots, x_n) \mapsto (t_1, \dots, t_n)$ .

Naším cílem je najít nejlepší proložení dat, čili fit, pomocí polynomické funkce řádu  $n$  ve tvaru

$$y(x, \theta) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_n x^n = \sum_{i=0}^n \theta_i x^i,$$

která je lineární v neznámých parametrech  $\theta = (\theta_0, \theta_1, \dots, \theta_n)$ . Takové modely nazýváme lineární a jejich vlastnosti budeme nadále využívat.

Abychom našli ten nejlepší možný fit, je nutno minimalizovat tzv. ztrátovou funkci (loss function)  $E(\Theta)$ . Tato funkce znázorňuje eukleidovskou vzdálenost od pozorovaných bodů  $\mathbf{t}$  k hledané funkci  $y(x, \theta)$ . Je tvaru

$$E(\theta) = \frac{1}{2} \sum_{n=1}^N [y(x_n, \theta) - t_n]^2 \quad (1.1)$$

Minimalizací ztrátové funkce získáme parametry

$$\theta = \arg \min_{\theta} E(\theta)$$

a jsme tudíž schopni sestavit předpis polynomu, který nejlépe daná data proloží. Tuto úlohu nazveme optimalizací.

V této bakalářské práci budeme výhradně používat gradientní metodu Gradient descent a její vylepšenou verzi ADAM.

#### 1.1.1 Gradient Descent

Jedná se iterativní optimalizační metodu pomocí které hledáme minimum dané funkce. My se snažíme zminimalizovat funkci  $L = (\mathbf{y} - \mathbb{X}\theta)^T (\mathbf{y} - \mathbb{X}\theta)$ , neboli funkci (1.1) pouze přepsanou pomocí maticového zápisu. Minimalizujeme  $L$ , tedy derivujeme dle vektoru  $\Theta$  a dostaneme

$$\nabla_{\theta} L = 2\mathbb{X}^T (\mathbb{X}\theta - \mathbf{y}),$$

kde  $\nabla_\theta$  značí gradient funkce  $L$  přes všechny hodnoty  $\theta$ . Použijeme bod  $\mathbf{a}$  funkce  $L(\theta)$  jako výchozí bod, ze kterého se pohybujeme ve směru záporného gradientu s krokem  $\gamma \in \mathbb{R}_+$ . Tento postup provádíme, dokud nejsme v minimu funkce a můžeme matematicky zapsat následujícím zápisem:

$$a_{n+1} = a_n - \gamma \nabla_\theta L(a_n)$$

## ADAM

Předchozí metoda není při větším množství dat tak rychlá, jak bychom pro výpočet minima funkce potřebovali. Používáme proto adaptivní iterační gradientní metodu ADAM (Adaptive Moment Estimation), která navíc používá druhý moment gradientu. Zatímco Gradient descent má krok stále stejný, u metody ADAM je krok  $\gamma$  adaptivní. Popřípadě můžeme ladit i zapomínací koeficienty, což už je ale mimo rámec této práce a my nebudeme při výpočtech využívat.

### 1.1.2 Metoda nejmenších čtverců

Uvažujme předeterminovaný systém  $n$  lineárních rovnic a  $p$  neznámých parametrech  $\theta = (\theta_0, \theta_1, \dots, \theta_p)$

$$\sum_{j=0}^p x_i^j \theta_j = y_i,$$

kde  $i \in (0, 1, \dots, n)$  a  $n > p$ . Přepíšeme pomocí maticového zápisu

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p \end{pmatrix} \cdot \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_0 \\ \varepsilon_1 \\ \vdots \\ \varepsilon_p \end{pmatrix}$$

Pro jednoduchost budeme tímto maticovým zápisem rozumět následující rovnici

$$\mathbf{y} = \mathbb{X} \cdot \theta + \epsilon \quad (1.2)$$

Naším cílem je získání parametrů  $\theta$ , proto obě strany rovnice vynásobíme zleva  $\mathbb{X}^T$ . Tím nám rovnice přejde do tvaru

$$\mathbb{X}^T \cdot \mathbf{y} = \mathbb{X}^T \cdot \mathbb{X} \cdot \hat{\theta}$$

Ted' už stačí rovnici zleva vynásobit inverzní maticí  $(\mathbb{X}^T \cdot \mathbb{X})^{-1}$ . Dostaneme tak konečné řešení

$$\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y} \quad (1.3)$$

Nicméně i zde lze parametry  $\theta$  odhadovat pomocí gradientní metody a to způsobem, který je popsán rovnicí (1.1).



## 1.2 Úvod do pravděpodobnosti a Bayesovská statistika

**Definice 1.2.1.** (Kolmogorova definice pravděpodobnosti). Mějme množinu  $\Omega$  vybavenou  $\sigma$ -algebrou  $\mathcal{A}$ , tedy souborem podmnožin obsahujícím  $\Omega$  a uzavřeným na doplňky a spočetná sjednocení. Pak libovolnou funkci  $P : \mathcal{A} \rightarrow \mathbb{R}$ , která splňuje :

1.  $(\forall A \in \mathcal{A})(P(A) \geq 0)$ .
2.  $P(\Omega) = 1$
3.  $\forall A_j$  disjunktní platí  $P(\sum_{j=1}^{\infty} A_j) = \sum_{j=1}^{\infty} P(A_j)$

**Věta 1.2.1.** (Vlastnosti  $P$ ). Mějme pravděpodobnostní prostor  $(\Omega, \mathcal{A}, P)$  a necht  $(\forall j \in \mathbb{N})(A_j \in \mathcal{A})$  a  $B \in \mathcal{A}$ . Pak platí:

1.  $P(\emptyset) = 0$ ,
2. Aditivita:  $P(\sum_{j=1}^n A_j) = \sum_{j=1}^n P(A_j)$ ,
3. Monotonie:  $A \subset B \Rightarrow P(A) \leq P(B)$ ,
4. Subtraktivita:  $A \subset B \Rightarrow P(B \setminus A) = P(B) - P(A)$ ,
5. Omezenost:  $(\forall A \in \mathcal{A})(P(A) \leq 1)$ ,
6. Komplementarita:  $A \in \mathcal{A} \Rightarrow P(A^C) = 1 - P(A)$

**Definice 1.2.2.** (Podmíněná pravděpodobnost). Necht'  $A, B \in \mathcal{A}$  a  $P(B) > 0$ . Pak definujeme podmíněnou pravděpodobnost:

$$P(A|B) = \frac{P(A, B)}{P(B)} \quad (1.4)$$

**Věta 1.2.2.** (Součinové pravidlo). Necht'  $A_1, \dots, A_n \in \mathcal{A}$  a dále necht' také  $P(A_1, \dots, A_n) > 0$ . Potom platí:

$$P(A_1, \dots, A_n) = P(A_1) \cdot P(A_2|A_1) \cdot P(A_3|A_2, A_1) \cdot \dots \cdot P(A_n|A_1, \dots, A_{n-1}) \quad (1.5)$$

**Věta 1.2.3.** (Bayseova věta). Necht'  $A \in \mathcal{A}$  a  $P(B) \neq 0$ . Potom platí:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1.6)$$

**Poznámka.**  $P(A)$  nazýváme prior a  $P(A|B)$  nazýváme posterior.

**Věta 1.2.4.** (Nezávislost jevů). Necht'  $A_j \in \mathcal{A} (\forall j \in \mathbb{N})$ . Potom jevy nazveme nezávislé právě tehdy když platí podmínka

$$P(A_1, \dots, A_k) = \prod_{i=1}^k P(A_i) \quad (1.7)$$

### 1.2.1 Hustoty pravděpodobnosti

Primárním cílem generativního modelování je hledání distribuce nebo-li hustoty pravděpodobnosti daných dat. Výhodou je, že pro hustotu pravděpodobnosti můžeme využívat stejné pravidlo podmíněnosti (1.4), součinné pravidlo (1.5) a Bayeseovo pravidlo (1.6). Toto se pro nás ukáže jako naprosto klíčové. Budeme uvažovat pouze spojitá rozdělení pravděpodobnosti náhodné veličiny  $X$ .

**Definice 1.2.3.** (*Hustota pravděpodobnosti*). Hustotou pravděpodobnosti náhodné veličiny  $X$  rozumíme spojitou funkci  $p(x)$ , která splňuje následující dvě podmínky:

1.  $p(x) \geq 0$
2.  $\int_{-\infty}^{\infty} p(x) dx = 1$

**Poznámka.** Hustotu pravděpodobnosti lze definovat také pro vícerozměrné funkce  $p(\mathbf{x})$ . Podmínky, které musí vícerozměrná hustota splňovat jsou analogické:

1.  $p(\mathbf{x}) \geq 0$
2.  $\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} p(\mathbf{x}) d\mathbf{x} = 1$

**Definice 1.2.4.** (*Střední hodnota náhodné veličiny*). Má-li náhodná veličina  $X$  spojitou hustotu pravděpodobnosti  $p(x)$ , definujeme její střední (očekávanou) hodnotu  $\mathbb{E}[X]$ , alternativně značeno  $\langle X \rangle$ , vztahem

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} xp(x)dx \quad (1.8)$$

**Definice 1.2.5.** (*Rozptyl náhodné veličiny*). Má-li náhodná veličina  $X$  spojitou hustotu pravděpodobnosti  $p(x)$ , definujeme rozptyl (varianci)  $\mathbb{D}[X]$ , alternativně značeno  $\text{var}(X)$ , vztahem

$$\mathbb{D}[X] = \int_{-\infty}^{\infty} x^2 p(x)dx = \int_{-\infty}^{\infty} (x - \mathbb{E}[X])^2 dx \quad (1.9)$$

**Definice 1.2.6.** (*Entropie*). Má-li náhodná veličina  $X$  spojitou hustotu pravděpodobnosti  $p(x)$ , definujeme entropii náhodné veličiny  $\mathbb{H}[X]$  vztahem

$$\mathbb{H}[X] = - \int_{-\infty}^{\infty} p(x) \ln p(x) dx \quad (1.10)$$

V následujícím textu uvedeme jednotlivá rozdělení a pro přehlednost jejich charakteristiky, které v této práci využíváme.

#### 1.2.1.1 Rovnoměrné rozdělení

Začneme jedním z nejjednodušších rozdělení. Rovnoměrné rozdělení, někdy také uniformní, přiřazuje všem hodnotám stejnou pravděpodobnost. Je definováno na intervalu  $(a, b)$  a můžeme ji vyjádřit následujícím způsobem.

$$U(a, b) = \begin{cases} \frac{1}{b-a}, & \text{pro } x \in (a, b) \\ 0, & \text{jinak} \end{cases} \quad (1.11)$$

- $\mathbb{E}[X] = \frac{1}{2}(a + b)$
- $\mathbb{D}[X] = \frac{1}{12}(b - a)^2$
- $\mathbb{H}[X] = \ln(b - a)$

### 1.2.1.2 Normální rozdělení

Nejdůležitější hustota pravděpodobnosti pro spojitě proměnné se nazývá normální nebo také Gaussovo rozdělení. Jeho hustota je definována  $\forall x \in \mathbb{R}$  pomocí dvou parametrů  $\mu \in \mathbb{R}$  a  $\sigma^2 > 0$  jako

$$\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\} \quad (1.12)$$

- $\mathbb{E}[X] = \mu$
- $\mathbb{D}[X] = \sigma^2$
- $\mathbb{H}[X] = \frac{1}{2} \ln 2\pi e \sigma^2$

Budeme využívat i d-rozměrnou variantu Gaussova rozdělení, které je definováno vztahem

$$\mathcal{N}(\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}, \quad (1.13)$$

kde  $\Sigma$  je  $d \times d$  matice, kterou nazveme kovarianční a  $\mu$  je vektor středních hodnot.

- $\mathbb{E}[X] = \mu$
- $\mathbb{D}[X] = \Sigma$
- $\mathbb{H}[X] = \frac{1}{2} \ln \det(2\pi e \Sigma)$

### 1.2.1.3 Gamma rozdělení

Gamma rozdělení je definováno stejně jako normální rozdělení pomocí dvou parametrů  $\alpha > 0$  a  $\beta > 0$ . Jeho hustota pravděpodobnosti má smysl pro  $\forall x > 0$  a můžeme ji najít v několika možných tvarech. My uvedeme tento:

$$\Gamma(\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp\{-\beta x\} \quad (1.14)$$

Stejně jako u Gaussova rozdělení uvedeme některé důležité charakteristiky.

- $\mathbb{E}[X] = \frac{\alpha}{\beta}$
- $\mathbb{D}[X] = \frac{\alpha}{\beta^2}$
- $\mathbb{H}[X] = \alpha - \ln \beta + \ln \Gamma(\alpha) + (1 - \alpha)\psi(\alpha)$

### 1.2.1.4 Inverzní gamma rozdělení

Inverzní gamma rozdělení je gamma rozdělení akorát pro převrácenou hodnotu  $x$ , je tedy opět popsáno dvěma parametry  $\alpha > 0$  a  $\beta > 0$  a definováno pro  $\forall x > 0$ . Jeho hustotu můžeme zapsat následovně:

$$i\Gamma(\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} \exp\left\{-\frac{\beta}{x}\right\} \quad (1.15)$$

Střední hodnota a rozptyl  $i\Gamma(\alpha, \beta)$  nejsou ale definována pro  $\alpha > 0$ , platí:

- $\mathbb{E}[X] = \frac{\beta}{\alpha-1}$ , pro  $\alpha > 1$
- $\mathbb{D}[X] = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)^2}$ , pro  $\alpha > 2$
- $\mathbb{H}[X] = \alpha + \ln \beta + \ln \Gamma(\alpha) - (1 + \alpha)\psi(\alpha)$

### 1.2.2 Bayesovská metoda nejmenších čtverců

Uvažujme standardní problém na lineární regresi (1.2), avšak více specifikujme šum  $\epsilon_i$ , kterými je zatížen každý bod  $y_i$  pro  $\forall i \in \hat{n}$  a to následovně

$$\mathbf{y} = \mathbb{X} \cdot \theta + \epsilon, \quad (1.16)$$

Platí  $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)$  pro jehož složky platí že jsou iid (independent and identically distributed) s rozdělením  $\epsilon_i \sim \mathcal{N}(0, 1)$  a tudíž  $p(\epsilon_i) \propto \exp\left\{-\frac{1}{2}\epsilon_i^2\right\}$ .

**Poznámka.** Zanedbáváme normalizační konstantu hustot, proto využíváme znak úměrnosti  $\propto$ .

Z rovnice (1.16) jednoduchou úpravou dostaneme

$$\epsilon = \mathbf{y} - \mathbb{X} \cdot \theta \quad (1.17)$$

Této rovnici odpovídá následující přepis pomocí hustot, přesněji vícerozměrného Gaussova rozdělení

$$p(\epsilon) \propto p(\mathbf{y}|\mathbb{X}, \theta) \propto \exp\left\{-\frac{1}{2}(\mathbf{y} - \mathbb{X}\theta)^\top (\mathbf{y} - \mathbb{X}\theta)\right\} \quad (1.18)$$

Snažíme se získat hustotu  $p(\theta|\mathbf{y}, \mathbb{X})$ , kterou získáme pomocí Bayesovy věty (1.6).

$$p(\theta|\mathbf{y}, \mathbb{X}) = \frac{p(\mathbf{y}|\mathbb{X}, \theta)p(\theta|\mathbb{X})}{p(\mathbf{y}|\mathbb{X})} \propto p(\mathbf{y}|\mathbb{X}, \theta)p(\theta|\mathbb{X}). \quad (1.19)$$

K tomu abychom mohli pokračovat ve výpočtu  $p(\theta|\mathbf{y}, \mathbb{X})$ , potřebujeme určit  $p(\theta|\mathbb{X})$ . Jelikož je  $\theta$  nezávislé na  $\mathbb{X}$ , můžeme psát pouze  $p(\theta)$ .

Pro hustotu  $p(\theta)$  předpokládáme následující vztah:

$$p(\theta) = \mathcal{N}(0, \alpha^{-1}\mathbb{I}) \propto \exp\left\{-\frac{1}{2}\theta^\top \theta \alpha\right\} \quad (1.20)$$

Nyní můžeme pokračovat dosazením do (2.8):

$$\begin{aligned} p(\mathbf{y}|\mathbb{X}, \theta)p(\theta|\mathbb{X}) &\propto \exp\left\{-\frac{1}{2}(\mathbf{y} - \mathbb{X}\theta)^\top (\mathbf{y} - \mathbb{X}\theta)\right\} \exp\left\{-\frac{1}{2}\theta^\top \theta \alpha\right\} \\ &\propto \exp\left\{-\frac{1}{2}(\mathbf{y}^\top \mathbb{Y} - \theta^\top \mathbb{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbb{X} \theta + \theta^\top \mathbb{X}^\top \mathbb{X} \theta + \theta^\top \theta \alpha)\right\} \\ &\propto \exp\left\{-\frac{1}{2}[\mathbf{y}^\top \mathbf{y} - \theta^\top \mathbb{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbb{X} \theta + \theta^\top (\mathbb{X}^\top \mathbb{X} + \alpha \mathbb{I}) \theta]\right\} \end{aligned} \quad (1.21)$$

Pro dokončení je důležitý předpoklad tvaru řešení a to:

$$p(\theta|\mathbf{y}, \mathbb{X}) \propto \exp\left\{-\frac{1}{2}(\theta - \hat{\theta})^\top \Sigma^{-1} (\theta - \hat{\theta}) + z\right\} \propto \exp\left\{-\frac{1}{2}(\theta - \hat{\theta})^\top \Sigma^{-1} (\theta - \hat{\theta})\right\} \exp\{z\}$$

který dále pomocí prvního tvaru upravíme tak, abychom dokázali určit  $\hat{\theta}$ ,  $\Sigma$  a  $z$ . Roznásobením dostaneme

$$p(\theta|\mathbf{y}, \mathbb{X}) \propto \exp\left\{-\frac{1}{2}(\theta^\top \Sigma^{-1} \theta - \hat{\theta}^\top \Sigma \theta - \theta^\top \Sigma^{-1} \hat{\theta} + \hat{\theta}^\top \Sigma^{-1} \hat{\theta}) + z\right\} \quad (1.22)$$

z čehož už při porovnání výrazu  $\theta^T (\mathbb{X}^T \mathbb{X} + \alpha \mathbb{I}) \theta$  v konečném tvaru rovnice (1.21) s výrazem  $\theta^T \Sigma^{-1} \theta$  v předchozí rovnici (1.22), plyne předpis pro

$$\Sigma^{-1} = \mathbb{X}^T \mathbb{X} + \alpha \mathbb{I} \quad (1.23)$$

Tento je výsledek je pro nás velmi důležitý a budeme jej i nadále využívat. Přímo porovnávejme další dva výrazy z těchto rovnic

$$-\mathbf{y}^T \mathbb{X} \theta = -\hat{\theta}^T \Sigma^{-1} \theta$$

Nyní z této rovnice jednoduchou úpravou a dosazením za  $\Sigma$  dostaneme další velmi důležitý předpis pro  $\hat{\theta}$ , a to

$$\hat{\theta} = \Sigma \mathbb{X}^T \mathbf{y} = (\mathbb{X}^T \mathbb{X} + \alpha \mathbb{I})^{-1} \mathbb{X}^T \mathbf{y} \quad (1.24)$$

Pro z nám zbývá

$$z = \mathbf{y}^T \mathbf{y} - \hat{\theta}^T \Sigma^{-1} \hat{\theta}.$$

Používáme ale znak úměrnosti a  $\exp \{z\}$  je pouze konstanta, můžeme ji tedy vynechat a dostaneme

$$p(\theta | \mathbf{y}, \mathbb{X}) \propto \exp \left\{ -\frac{1}{2} (\theta - \hat{\theta})^T \Sigma^{-1} (\theta - \hat{\theta}) \right\}$$

### 1.2.3 Divergence

Divergence je funkce  $D(\cdot \| \cdot) : S \times S \rightarrow \mathcal{R}$ , kde je  $S$  je prostor pravděpodobnostních rozdělení a které splňuje následující dvě podmínky:

1.  $D(q \| p) \geq 0$
2.  $D(q \| p) = 0$  pro  $p = q$

Divergence do jisté popisuje vzdálenost nebo rozdíl mezi dvěma distribucemi. Jelikož divergence nemusí splňovat podmínku symetrie a trojúhelníkové nerovnosti, nejedná se tedy o metriku, nýbrž o semimetriku.

#### f-divergence

Nejdůležitější skupinou divergencí jsou takzvané f-divergence. Jsou definovány pomocí konvexní funkce  $f(x)$ , kde  $x > 0$  a takové že  $f(1) = 0$ . Jsou tvaru

$$D_f(q \| p) = \int_{\text{supp}(q)} q(x) f\left(\frac{q(x)}{p(x)}\right) dx \quad (1.25)$$

kde  $\text{supp}(q)$  značí nosič funkce  $q(x)$ .

#### Kullback-Leiblerova divergence

Pro nás bude užitečná tzv. Kullback-Leiblerova divergence, kde za funkci  $f$  bereme přirozený logaritmus, značeno  $\log$ . To je rozhodně konvexní funkce pro kterou platí  $\log 1 = 0$ . Tvar KL-divergence je následující:

$$D_{KL}(q \| p) = \int_{\text{supp}(q)} q(x) \log \frac{q(x)}{p(x)} dx \quad (1.26)$$

### 1.2.4 ELBO

Předpokládejme že máme pozorování  $X$  a  $Z$  jsou skryté (latentní) proměnné. Posteriorní distribuci latentní proměnné  $Z$  můžeme napsat pomocí Bayesova pravidla (1.6), jehož jmenovatel se někdy také nazývá evidence, takto:

$$p(Z|X) = \frac{p(X|Z)p(Z)}{p(X)} = \frac{p(X|Z)p(Z)}{\int p(X,Z)dZ} \quad (1.27)$$

Dále zadefinujeme nový objekt

$$\log p(X) = \log \int p(X,Z)dZ \quad (1.28)$$

a abychom mohli pokračovat, využijeme pomocnou funkci  $q(Z|\theta)$

$$\log p(X) = \log \int p(X,Z)dZ = \log \int q(Z|\theta) \frac{p(X,Z)}{q(Z|\theta)}dZ = \log \mathbb{E}_q \left[ \frac{p(X,Z)}{q(Z|\theta)} \right] \quad (1.29)$$

Dále využijeme Jensenovu nerovnost, díky které získáme spodní hranici (lower bound), odtud Evidence Lower Bound, čili ELBO.

$$\log \mathbb{E}_q \left[ \frac{p(X,Z)}{q(Z|\theta)} \right] \geq \mathbb{E}_q \left[ \log \frac{p(X,Z)}{q(Z|\theta)} \right] \quad (1.30)$$

ELBO můžeme rozepsat pomocí součinnového pravidla (1.5), využít vlastností logaritmu a dle definice KL-divergence (1.26), přepsat do tvaru

$$\mathbb{E}_q \left[ \log \frac{p(X,Z)}{q(Z|\theta)} \right] = \mathbb{E}_q \left[ \log \frac{p(X|Z)p(Z)}{q(Z|\theta)} \right] = \mathbb{E}_q [\log p(X|Z)] - \mathbb{E}_q \left[ \log \frac{p(Z)}{q(Z|\theta)} \right] \quad (1.31)$$

$$= \mathbb{E}_q [\log p(X|Z)] - D_{KL}(q(Z|\theta) \| p(Z)) = \mathcal{L}(\theta) \quad (1.32)$$

Budeme-li maximalizovat ELBO přes všechny variační parametry  $\theta$ , získáme nejbližší možnou hodnotu k  $\log p(X)$ . Navíc je maximalizace ELBO ekvivalentní k minimalizaci KL-divergence mezi  $q(Z|\theta)$  a  $p(Z|\theta)$ , jelikož platí

$$\begin{aligned} D_{KL}(q(Z|\theta) \| p(Z|X)) &= \mathbb{E}_q \left[ \log \frac{q(Z|\theta)}{p(Z|X)} \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(Z|\theta)p(X)}{p(X|Z)p(Z)} \right] \\ &= -\mathbb{E}_q [\log p(X|Z)] + \mathbb{E}_q \left[ \log \frac{q(Z|\theta)}{p(Z|X)} \right] + \mathbb{E}_q [\log p(X)] \\ &= -\mathbb{E}_q [\log p(X|Z)] + D_{KL}(q(Z|\theta) \| p(Z)) + \log p(X) \end{aligned} \quad (1.33)$$

Z toho jednoduchou úpravou dostaneme konečný vztah

$$D_{KL}(q(Z|\theta) \| p(Z|X)) = -\mathcal{L}(\theta) + \log p(X) \quad (1.34)$$

### Příklad

Předvedeme příklad, jak ELBO využít v praxi.

Uvažujme pouze sadu dvou pozorování  $y_1$  a  $y_2$  s normálním rozdělením  $\mathcal{N}_i(\theta, 1)$  pro  $i \in 1, 2$ . Dále uvažujme jeden parametr  $\theta \sim \mathcal{N}(0, \alpha)$  a nechť  $\alpha$  má inverzní gamma rozdělení, tedy  $\alpha \sim i\Gamma(0, 0)$ . Snažíme se

získat sdruženou distribuci parametrů  $\theta$  a  $\alpha$ , tedy  $p(\theta, \alpha|y_1, y_2)$ . Tuto distribuci můžeme přepsat pomocí definice podmíněné pravděpodobnosti (1.4) a řetězového pravidla (1.5) jako

$$p(\theta, \alpha|y_1, y_2) = \frac{p(\theta, \alpha, y_1, y_2)}{p(y_1, y_2)} = \frac{p(y_1|\theta)p(y_2|\theta)p(\theta)p(\alpha)}{p(y_1, y_2)} \quad (1.35)$$

Dosažením předpokladů do čitatele dostaneme:

$$p(y_1|\theta)p(y_2|\theta)p(\theta)p(\alpha) \propto \exp\left\{-\frac{1}{2}(y_1 - \theta)^2\right\} \cdot \exp\left\{-\frac{1}{2}(y_2 - \theta)^2\right\} \cdot \frac{1}{\sqrt{\alpha}} \exp\left\{-\frac{\theta^2}{2\alpha}\right\} \cdot \frac{1}{\alpha} \quad (1.36)$$

Zdánlivě se nám může zdát určení jmenovatele jako jednoduché, protože pravděpodobnost  $p(y_1, y_2)$  lze získat tzv. marginalizací, nebo-li vyintegrováním přes  $\theta$  a  $\alpha$ .

$$\begin{aligned} p(y_1, y_2) &= \int p(\theta, \alpha, y_1, y_2) d\theta d\alpha \\ &= \int p(y_1|\theta)p(y_2|\theta)p(\theta)p(\alpha) d\theta d\alpha \\ &= \int \exp\left\{-\frac{1}{2}(y_1 - \theta)^2\right\} \cdot \exp\left\{-\frac{1}{2}(y_2 - \theta)^2\right\} \cdot \frac{1}{\sqrt{\alpha}} \exp\left\{-\frac{\theta^2}{2\alpha}\right\} \cdot \frac{1}{\alpha} d\theta d\alpha \end{aligned} \quad (1.37)$$

Po bližším přezkoumání (1.37) zjistíme, že nelze přes  $\alpha$  vyintegrovat. Proto použijeme ELBO. Dle definice KL-divergence a za předpokladu  $q(\theta, \alpha) = q(\theta)q(\alpha)$  můžeme psát:

$$D_{KL}(q(\theta, \alpha|\mu, \sigma, \gamma, \delta) \| p(\theta, \alpha|y_1, y_2)) = \int_G q(\alpha)q(\theta) \ln \left\{ \frac{q(\alpha)q(\theta)}{p(\theta, \alpha|y_1, y_2)} \right\} d\theta d\alpha = \diamond \quad (1.38)$$

kde  $G = \text{supp}(q(\theta)) \times \text{supp}(q(\alpha))$ . Nezapomínejme, že  $q(\theta)$  a  $q(\alpha)$  jsou distribuce, pro které si apriori zvolíme

$$\begin{aligned} q(\theta) &= \mathcal{N}(\mu, \sigma) \\ q(\alpha) &= i\Gamma(\gamma, \delta) \end{aligned}$$

Dle (1.2.3) navíc víme, že platí  $\int_G q(\alpha)q(\theta) d\theta d\alpha = 1$ . Výraz budeme rozepisovat pomocí pravidel pro logaritmy a postupně upravovat. Výraz  $p(y_1, y_2)$  v integrálu je konstanta, kterou můžeme pro jednoduchost zanedbat. Výsledek budeme na konci maximalizovat a konstanta polohu maxima nemění.

$$\begin{aligned} \diamond &\propto \int_G q(\alpha)q(\theta) \ln \frac{q(\alpha)q(\theta)}{p(y_1|\theta)p(y_2|\theta)p(\theta)p(\alpha)} d\theta d\alpha \\ &= \int_G q(\theta)q(\alpha) (-\ln p(y_1|\theta) - \ln p(y_2|\theta) - \ln p(\theta) - \ln p(\alpha) + \ln q(\theta) + \ln q(\alpha)) d\alpha d\theta \end{aligned} \quad (1.39)$$

Poslední dva výrazy jsou tzv. entropie pro Gaussovo rozdělení, resp. inverzní gamma rozdělení. Můžeme využít již známých výsledků:

$$\begin{aligned} \int q(\theta) \ln q(\theta) d\theta &\propto -\frac{1}{2} \ln \sigma \\ \int q(\alpha) \ln q(\alpha) d\alpha &= -\gamma - \ln \delta \Gamma(\gamma) + (1 + \gamma) \psi(\gamma) \end{aligned}$$

Vypočítejme zbývající výrazy, kde pro jednoduchost budeme pro střední hodnoty využívat značení pomocí špičatých závorek:

$$\begin{aligned}\int_G q(\theta)q(\alpha) \ln p(y_1|\theta) d\alpha d\theta &= \left\langle -\frac{1}{2}(y_1 - \theta)^2 \right\rangle = -\frac{1}{2}(y_1^2 - 2y_1\mu + \mu^2 + \sigma) \\ \int_G q(\theta)q(\alpha) \ln p(y_2|\theta) d\alpha d\theta &= \left\langle -\frac{1}{2}(y_2 - \theta)^2 \right\rangle = -\frac{1}{2}(y_2^2 - 2y_2\mu + \mu^2 + \sigma) \\ \int_G q(\theta)q(\alpha) \ln p(\theta) d\alpha d\theta &= \left\langle -\frac{\theta^2}{2\alpha} - \frac{1}{2} \ln \alpha \right\rangle = -\frac{1}{2} \left( (\mu^2 + \sigma) \frac{\gamma}{\delta} + \ln \delta - \psi(\gamma) \right) \\ \int_G q(\theta)q(\alpha) \ln p(\alpha) d\alpha d\theta &= \langle -\ln \alpha \rangle = \psi(\gamma) - \ln \delta\end{aligned}$$

Nyní máme všechny výrazy pro výpočet distribuce  $q(\theta, \alpha)$  numericky a to pomocí maximalize KL-divergence přes parametry  $\mu, \sigma, \gamma, \delta$ .

### 1.3 Teorie grafů

Pro popis složitějších datových struktur, zejména těch stromových, můžeme využít teorie grafů. Strom je totiž speciální případ grafu. Pokusme se to řádně zdefinovat.

**Definice 1.3.1.** (Graf) Grafem  $G$  se rozumí dvojice  $(V, H)$ , kde  $V$  je množina vrcholů grafu  $G$ ,  $H$  je množina hran tohoto grafu a tyto množiny jsou vzájemně disjunktní.

**Definice 1.3.2.** (Cesta v grafu) Cestou v grafu rozumíme posloupnost vrcholů a hran  $(v_0, h_1, v_1, \dots, h_t, v_t)$ , kde vrcholy  $v_0, \dots, v_t$  jsou navzájem různé vrcholy grafu  $G$  a pro každé  $i = 1, 2, \dots, t$  je  $e_i = \{v_{i-1}, v_i\} \in H$

**Definice 1.3.3.** (Souvislost grafu) Řekneme že graf  $G$  je souvislý, jestliže pro každé dva vrcholy  $v_0$  a  $v_1$  existuje v  $G$  cesta z  $v_0$  do  $v_1$ .

**Definice 1.3.4.** (Cyklus v grafu) Cyklem v grafu  $G$  rozumíme posloupnost vrcholů a hran  $(v_0, h_1, v_1, \dots, h_t, v_t = v_0)$ , kde vrcholy  $v_0, \dots, v_{t-1}$  jsou navzájem různé vrcholy grafu  $G$  a pro každé  $i = 1, 2, \dots, t$  je  $e_i = \{v_{i-1}, v_i\} \in H$

**Definice 1.3.5.** (Strom) Strom je souvislý graf neobsahující cyklus.

V podstatě si to můžeme představit opravdu jako strom - má jeden kořen, v první úrovni se dělí na  $K_1$  větví, každá další větev se v druhé úrovni dělí na  $K_{2_i}$  a tak dále. My se budeme v této práci zabývat pouze kořenem a první úrovní větví.



## Kapitola 2

# Generativní modely

Ve strojovém učení se setkáváme s dvěma hlavními typy modelů a to jsou generativní modely a diskriminativní modely. Jak už napovídá název této práce, budeme se zde zabývat výhradně generativními modely.

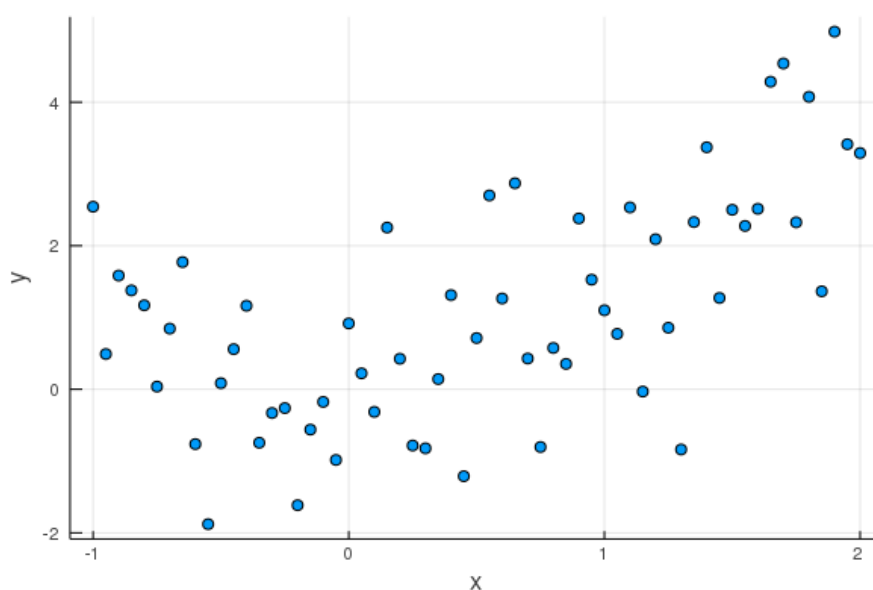
**Definice 2.0.1.** (*Generativní model*) Mějme nějakou množinu datových záznamů  $X$  a nějakou množinu  $Y$ . Cílem je tuto množinu  $X$  klasifikovat pomocí množiny  $Y$ . Generativní model je potom takový model, který se učí sdruženou distribuci  $p(X, Y)$ .

### Příklad

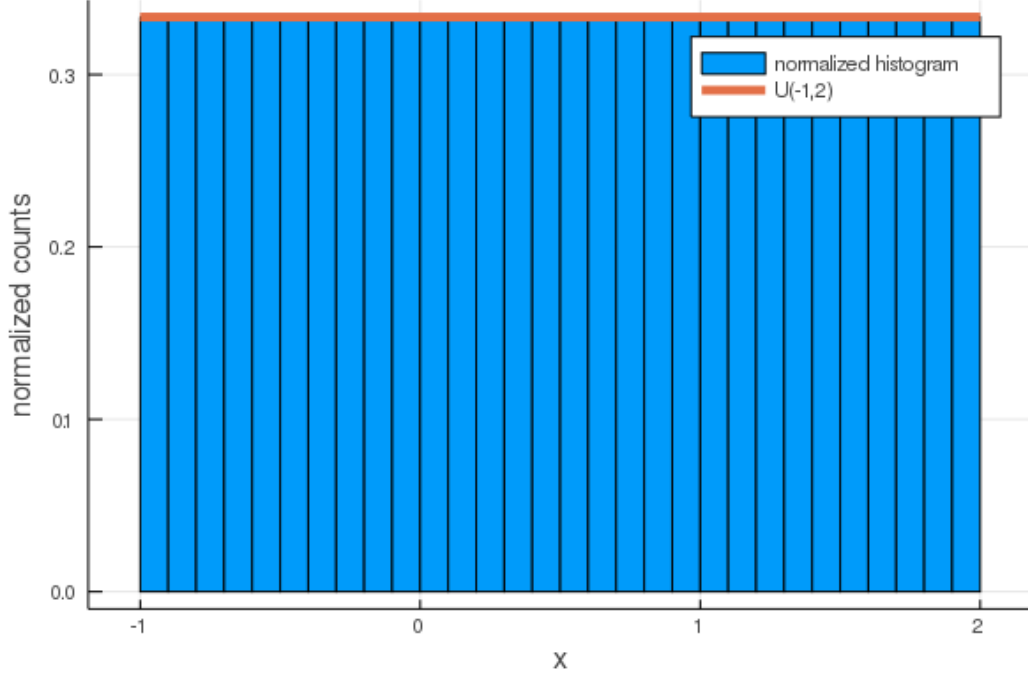
Jeden ze způsobů jak odhadnout distribuci  $p(Y, X)$  je využití součinového pravidla (1.5), pomocí kterého získáme

$$p(Y, X) = p(Y|X)p(X) \quad (2.1)$$

Problém je tedy převeden na hledání distribucí  $p(Y|X)$  a  $p(X)$ . Pro ilustraci uvažujme následující množinu datových záznamů.



Určit distribuci  $p(X)$  není nic těžkého, jelikož jsou tato data na ose  $x$  rozdělena rovnoměrně. To můžeme určit například z histogramu  $x$ -ových souřadnic jednotlivých bodů. Ten vypadá následovně:



Obrázek 2.1: Normalizovaný histogram  $x$ -ových souřadnic a jeho distribuce  $p(X)$ .

Histogram je znormalizován a oranžovou čarou je zde znázorněno rovnoměrné rozdělení  $p(X) = U(-1, 2)$ . Nyní přejdeme k hledání distribuce  $p(Y|X)$ . Tu můžeme určit pomocí metody nejmenších čtverců (1.3), protože víme že pro takovou distribuci platí  $p(Y|X) = \mathcal{N}(X\theta, \sigma^2 I)$ , kde  $\sigma^2$  je rozptyl šumu  $\varepsilon_i$ .

## 2.1 Variační autoencoder

Cílem je najít hustotu  $p(x)$  vzorků  $\{x^{(i)}\}_{i=1}^N$ , jehož empirická hustota se dá zapsat pomocí delta funkce

$$p_{\text{emp}}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x^{(i)}) \quad (2.2)$$

Předpokládáme následující vztahy  $x = f_{\theta}(z) + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$  a vzájemnou nezávislost  $x^{(i)}$ . Z toho můžeme určit distribuce:

$$\begin{aligned} p(x|z) &= \mathcal{N}(f_{\theta}(z), \sigma^2 I), \\ p(z) &= \mathcal{N}(0, I), \end{aligned} \quad (2.3)$$

a proto má smysl využít následující formu apromaximace

$$p(x) = \int p(x|z)p(z)dz \quad (2.4)$$

### 2.1.1 Naivní přístup

K nalezení  $p(x)$  je třeba najít parametry  $\theta$  transformace  $f(z)$ , proto zkusme využít KL-divergence a hledat tak  $\theta$  minimalizací  $D_{KL}(p_{\text{emp}}(x) \| p_{\theta}(x))$

$$\begin{aligned}\hat{\theta} &= \arg \min \sum_{i=1}^N \log p(x^{(i)}) \\ &= \arg \min \sum_{i=1}^N \log \int \mathcal{N}(f_{\theta}(z), \sigma^2 I) \mathcal{N}(0, I) dz \\ &= \arg \min \sum_{i=1}^N \log \sum_{j=1}^N \exp \left\{ -\frac{1}{2\sigma^2} (x - f_{\theta}(z)) \right\}\end{aligned}\tag{2.5}$$

### 2.1.2 Variační Bayseova metoda

Lepší metodou se ukazuje vzorkovat z podmíněné distribuce  $q(z|x)$  a využít ELBO:

$$\begin{aligned}D_{KL}(q(z|x) \| p(z|x)) &= \mathbb{E}_q [\log q(z|x) - \log p(z|x)] \\ &= \mathbb{E}_q [\log q(z|x) - \log p(x|z) - \log p(z) + \log p(x)]\end{aligned}\tag{2.6}$$

Tuto rovnici můžeme přepsat pomocí KL-divergence

$$\log p(x) - D_{KL}(q(z|x) \| p(z|x)) = \mathbb{E}_q [\log p(x|z)] - D_{KL}(q(z|x) \| p(z))\tag{2.7}$$

kde pravá strana této rovnice je lower bound  $\log p(x)$ . Jestliže vybereme parametrickou formu distribuce

$$q(z|x) = \mathcal{N}(\mu_{\theta}(x), \sigma_{\phi}^2(x))\tag{2.8}$$

můžeme parametry  $\theta$  a  $\phi$  minimalizovat zároveň a to následovně:

### Příklad

## **Kapitola 3**

# **Stromové struktury**

# **Závěr**

Text závěru....

# Literatura

- [1] S. Allen, J. W. Cahn: *A microscopic theory for antiphase boundary motion and its application to antiphase domain coarsening*. Acta Metall., 27:1084-1095, 1979.
- [2] G. Ballabio et al.: *High Performance Systems User Guide*. High Performance Systems Department, CINECA, Bologna, 2005. [www.cineca.it](http://www.cineca.it)
- [3] J. Becker, T. Preusser, M. Rumpf: *PDE methods in flow simulation post processing*. Computing and Visualization in Science, 3(3):159-167, 2000.