

Generativní modely dat popsané stromovou strukturou

Jakub Bureš

Katedra matematiky

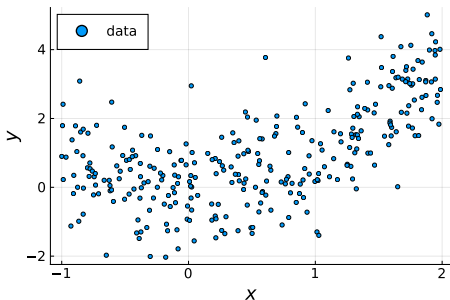
Vedoucí práce: doc. Ing. Václav Šmídl, Ph.D.

27. srpna 2020

- 1 Motivace
- 2 Generativní modely
 - Variační autoencdoer
- 3 Stromové struktury
- 4 Otázky oponenta

Motivace

Máme k dispozici následující data a chceme predikovat nová, hledáme tedy hustotu $p(y|x)$.



Obrázek: Zadaná data, ze kterých chceme predikovat nová.

Nic těžkého \Rightarrow problém vede na úlohu nejmenších čtverců.

Model

$$\mathbf{y} = \mathbb{X} \cdot \theta + \epsilon. \quad (1)$$

Předpokládáme $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ a položíme $X^T = (1 \ x \ x^2 \ \dots \ x^s)$, kde s je stupeň polynomu, jakým data prokládáme. Obdržíme tvar hustoty

$$p(y|x) = \mathcal{N}(X^T \cdot \theta, \sigma^2) \quad (2)$$

Problém nastane v okamžiku, kdy budeme chtít znát například $p(y|x = 20)$, čili extrapolace mimo interval daných dat. Odpověď nemusí být přesná.

Řešení?

- $\int p(y|x, \theta) p(\theta) d\theta = p(y|x)$
- **Generativní model**

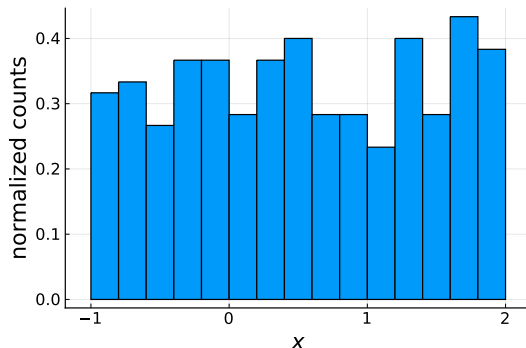
Generativní model

Mějme nějakou množinu datových záznamů $\mathbf{x} = \{x_1, \dots, x_n\}$, představující nezávislé proměnné a nějakou množinu $\mathbf{y} = \{y_1, \dots, y_n\}$, jakožto závislé proměnné. Generativní model je potom takový model, který se učí sdruženou hustotu pravděpodobnosti $p(x, y)$.

- Odhad hustoty pravděpodobnosti $p(x, y)$
- Součinné pravidlo $p(y, x) = p(y|x) \cdot p(x)$
- Pokusíme se tedy v dalším kroku najít $p(x)$.

Generativní model

Hustotu $p(x)$ můžeme určit například pomocí histogramu x -ových složek



Obrázek: Histogram x -ových složek zadaných dat.

Histogram odpovídá uniformnímu rozdělení

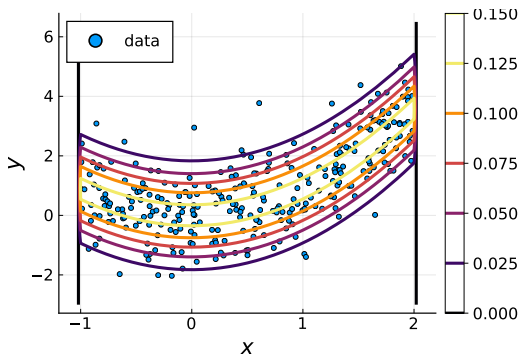
$$p(x) = U(-1, 2)$$

(3)

Generativní model

Pro sdruženou hustotu $p(y, x)$ pak dostaneme vztah

$$p(y, x) = U(-1, 2) \cdot \mathcal{N}(X^T \cdot \theta, \sigma^2) \quad (4)$$



Obrázek: Contour plot sdružené hustoty $p(y, x)$.

Variační autoencoder

Cílem je najít hustotu $p(\mathbf{x})$ vzorků $\{x_i\}_{i=1}^n$ za pomoci latentní proměnné \mathbf{z} , u které předpokládáme

$$p(\mathbf{z}) = \mathcal{N}(0, \mathbb{I}) \quad (5)$$

Jelikož lze každou náhodnou veličinu transformovat pomocí věty o transformaci náhodné veličiny, budeme hledat můžeme f_θ takovou, že

$$\mathbf{x} = f_\theta(\mathbf{z}). \quad (6)$$

Takovou funkci je obecně velmi těžké najít, proto hledáme f_θ následujícího modelu

$$\mathbf{x} = f_\theta(\mathbf{z}) + \epsilon. \quad (7)$$

kde $\epsilon \sim \mathcal{N}(0, \sigma^2 \cdot \mathbb{I})$ a tudíž volíme $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(f_\theta(\mathbf{z}), \sigma^2 \cdot \mathbb{I})$.
Využijeme následující formu aproximace

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}. \quad (8)$$

Víme-li, že

$$D_{KL}(q\|p) = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}, \quad (9)$$

můžeme použít ELBO.

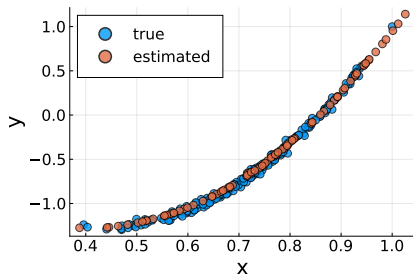
$$\begin{aligned} D_{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}|\mathbf{x})) &= \mathbb{E}_q[\log q(\mathbf{z}|\mathbf{x}) - \log p(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_q[\log q(\mathbf{z}|\mathbf{x}) - \log p(\mathbf{x}|\mathbf{z}) - \log p(\mathbf{z}) + \log p(\mathbf{x})]. \end{aligned} \quad (10)$$

Tuto rovnici můžeme přepsat pomocí další KL-divergence

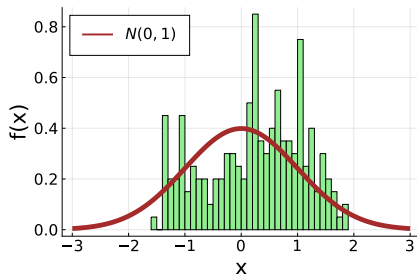
$$\log p(\mathbf{x}) - D_{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}|\mathbf{x})) = \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})), \quad (11)$$

kde pravá strana této rovnice je lower bound objektu $\log p(\mathbf{x})$, tedy ELBO.

Variační autoencoder



(a) Skutečné vzorky $\{\mathbf{x}, \mathbf{y}\}$ (modře) a jejich odhad pomocí VAE (červeně).



(b) Transformace odhadu vzorků $\{\mathbf{x}, \mathbf{y}\}$ zpět na Histogram vzorků \mathbf{z}

- Červené vzorky byly vygenerovány pomocí $f_{\theta}(z) : \mathbb{R}^1 \rightarrow \mathbb{R}^2$
- Histogram byl vygenerován pomocí $g(x, y) : \mathbb{R}^2 \rightarrow \mathbb{R}^1$

Příklad inspirovaný finanční aplikací.

- Máme k dispozici i -tého klienta, u které víme, zda-li splácel půjčku $y_i \in \{0, 1\}$, 0 značí ano, 1 značí ne.
- Dále máme k dispozici seznam jeho transakcí $X = \{x_{i,j}\}_{j=1}^{N_x}$, kde N_x počet těchto transakcí a u jednotlivých klientů se liší.
- Hodnoty a počty těchto transakcí se liší navíc, jestli daný klient splácel nebo nesplácel půjčku.
- Uvažujeme dvě gaussové směsi (GM), $m \in \hat{a}, n \in \hat{b}$

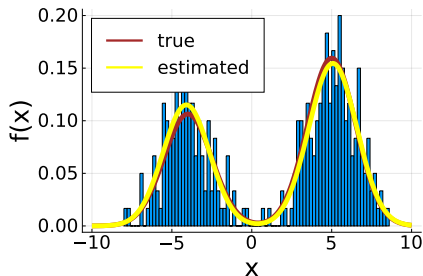
$$\begin{aligned} p(x_m|y=1) &= w_1 \cdot \mathcal{N}(\mu_1, \sigma_1^2) + (1 - w_1) \cdot \mathcal{N}(\mu_2, \sigma_2^2), \\ p(x_n|y=0) &= w_2 \cdot \mathcal{N}(\mu_3, \sigma_3^2) + (1 - w_2) \cdot \mathcal{N}(\mu_4, \sigma_4^2), \end{aligned} \quad (12)$$

- Dále uvažujeme

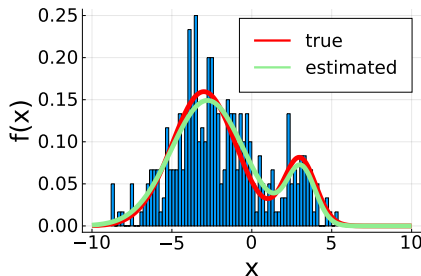
$$\begin{aligned} p(N_{kx}^{(1)} | y = 1) &= \text{Po}(\lambda_1) & k \in \hat{c}, \\ p(N_{lx}^{(0)} | y = 0) &= \text{Po}(\lambda_2) & l \in \hat{d}. \end{aligned} \tag{13}$$

- Jak nyní rozhodnout, do které třídy (schopný či neschopný splácet) patří nový klient?
- Nejprve odhadneme všechny parametry všech rozdělení pomocí MLE.

$$\begin{aligned}\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1^2, \hat{\sigma}_2^2, \hat{w}_1 &= \arg \max_{\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, w_1} \log \left(\prod_{m=1}^a p(x_m | y = 1) \right) \\ \hat{\mu}_3, \hat{\mu}_4, \hat{\sigma}_3^2, \hat{\sigma}_4^2, \hat{w}_2 &= \arg \max_{\mu_3, \mu_4, \sigma_3^2, \sigma_4^2, w_2} \log \left(\prod_{n=1}^b p(x_n | y = 0) \right), \\ \hat{\lambda}_1 &= \frac{1}{c} \sum_{k=1}^c N_{\text{kx}}^{(1)}, \\ \hat{\lambda}_0 &= \frac{1}{d} \sum_{l=1}^d N_{\text{lx}}^{(0)}\end{aligned}\tag{14}$$



(c) $p(x|y = 1)$



(d) $p(x|y = 0)$

Obrázek: Dvě GM, kde červenou a hnědou barvou jsou nakresleny skutečné distribuce, zeleně a žlutě jsou jejich MLE.

- MLE zde funguje i přesto, že MLE gaussovske směsi nemá analytické řešení.

- V dalším kroku sestavíme

$$\begin{aligned} p(\mathbf{x}, N_x^{(1)} | y = 1) &= \left(\prod_{i=1}^{N_x^{(1)}} p(x_i | y = 1) \right) \cdot p(N_x^{(1)} | y = 1), \\ p(\mathbf{x}, N_x^{(0)} | y = 0) &= \left(\prod_{j=1}^{N_x^{(0)}} p(x_j | y = 0) \right) \cdot p(N_x^{(0)} | y = 0) \end{aligned} \quad (15)$$

- S jejich pomocí provedeme věrohodnostní poměry

$$\begin{aligned} \Lambda_0(\mathbf{x}) &= \frac{p(\mathbf{x}, N_x^{(0)} | y = 0)}{p(\mathbf{x}, N_x^{(1)} | y = 1) + p(\mathbf{x}, N_x^{(0)} | y = 0)} \in \langle 0, 1 \rangle \\ \Lambda_1(\mathbf{x}) &= \frac{p(\mathbf{x}, N_x^{(1)} | y = 1)}{p(\mathbf{x}, N_x^{(1)} | y = 1) + p(\mathbf{x}, N_x^{(0)} | y = 0)} \in \langle 0, 1 \rangle. \end{aligned} \quad (16)$$

Děkuji za pozornost.

Co je to regulární zobrazení ve větě o transformaci náhodné veličiny?

Jaký je rozdíl mezi metrikou a semimetrikou?

Proč je jasné, že

$$\bar{x}_i = \frac{1}{N_x^{(i)}} \sum_{l=1}^{N_x^{(i)}} x_l^{(i)} \quad (17)$$

bude mít normální rozdělení? Platí to obecně nebo za nějakých předpokladů?