

## Lekce 2 [35]

Describe maximum likelihood estimation, as minimizing NLL, cross-entropy and KL divergence. [10]

### Self information

- $I(x) = -\log P(x)$
- Jak moc jsme překvapeni, když dostaneme  $x \sim P$
- Pro nezávislé jevy se počítá, pro jevy s pností 1 je rovna 0

### Entropie

- $H(P) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)]$
- Množství překvapení v distribuci  $P$

### Cross-entropy

- $H(P, Q) = -\mathbb{E}_{x \sim P}[\log Q] = \mathbb{E}_{x \sim P}[I_Q(x)]$
- V podmíněném případě pak  $H(Y|X) = H(X, Y) - H(X)$ , tedy jak moc jsem překvapen když se dozvím obojí oproti tomu, když se dozvím jen  $X$
- Měří, jak moc budu překvapený, když budu tahat  $x$  z distribuce  $P$ , ale své překvapení budu měřit na základě distribuce  $Q$

### Kullback-Leibler Divergence

- $D_{KL}(P||Q) = H(P, Q) - H(P) = -\mathbb{E}_{x \sim P}[\log P - \log Q]$
- Není symetrická
- V zásadě říká, jak "špatná" je moje distribuce  $Q$ . Konkrétně zjišťuje o jak moc více budu překvapen, když tahám  $x$  z  $P$ , ale překvapení měřím skrze  $Q$ 
  - Čím jsou si  $Q$  a  $P$  podobnější, tím menší tohle "překvapení navíc" bude

### MLE

- Samo o sobě je to takové hledání parametrů modelu, aby

$$\theta_{\text{ML}} = \arg \max_{\theta} p_{\text{model}}(\mathbb{X}; \theta) \quad (1)$$

Což se dá dále upravovat, až se dostaneme NLL, binární crossentropii, a KL divergenci

$$\begin{aligned} \theta_{\text{ML}} \quad p(y|x; \theta) &= \arg \max_{\theta} \prod_{i=1}^n p(y_i | x_i; \theta) \\ &= \arg \min_{\theta} -\sum_{i=1}^n \log p(y_i | x_i; \theta) \quad \leftarrow \text{NLL} \\ &= \arg \min_{\theta} \mathbb{E}_{x, y \sim p_{\text{data}}} [-\log p(y | x; \theta)] \\ &= \arg \min_{\theta} H(\hat{p}_{\text{data}}, p(y | x; \theta)) \\ &= \arg \min_{\theta} D_{KL}(\hat{p}_{\text{data}} || p(y | x; \theta)) + \dots \end{aligned}$$

$$\begin{aligned}
\theta_{\text{ML}} &= \arg \max_{\theta} p_{\text{model}}(\mathbb{Y}|\mathbb{X}; \theta) \\
&= \arg \min_{\theta} \sum_{i=1}^m -\log p_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}; \theta) \quad \leftarrow \text{NLL} \\
&= \arg \min_{\theta} \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} [-\log p_{\text{model}}(y|\mathbf{x}; \theta)] \\
&= \arg \min_{\theta} H(\hat{p}_{\text{data}}, p_{\text{model}}(y|\mathbf{x}; \theta)) \quad \leftarrow \text{binary crossentropy loss} \\
&= \arg \min_{\theta} D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}(y|\mathbf{x}; \theta)) + H(\hat{p}_{\text{data}}) \quad \leftarrow \text{KL-divergence}
\end{aligned}$$

Define mean squared error and show how it can be derived using MLE. [5]

$$MSE = \mathbb{E}[(\hat{y} - y)^2] = \frac{1}{m} \sum_{i=1}^m ((f(x_i; \theta) - y_i)^2) \quad (2)$$

Pokud se nám při regresi nechce odhadovat celá distribuce, můžeme si usnadnit práci, predikovat pouze její střední hodnotu a říct, že ta distribuce je normální s nějakým rozptylem (a s tou naší střední hodnotou).

Dává to smysl, protože normální rozdělení má mezi rozděleními se stejnou střední hodnotou a rozptylem maximální entropii, tedy nejméně navíc vnesené informace.

MLE potom vyjde

$$\begin{aligned}
\arg \max_{\theta} p(y|\mathbf{x}; \theta) &= \arg \min_{\theta} \sum_{i=1}^m -\log p(y^{(i)}|\mathbf{x}^{(i)}; \theta) \\
&= \arg \min_{\theta} -\sum_{i=1}^m \log \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \hat{y}(\mathbf{x}^{(i)}; \theta))^2}{2\sigma^2}\right) \\
&= \arg \min_{\theta} -m \log(2\pi\sigma^2)^{-1/2} - \sum_{i=1}^m -\frac{(y^{(i)} - \hat{y}(\mathbf{x}^{(i)}; \theta))^2}{2\sigma^2} \\
&= \arg \min_{\theta} \sum_{i=1}^m \frac{(y^{(i)} - \hat{y}(\mathbf{x}^{(i)}; \theta))^2}{2\sigma^2} = \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}(\mathbf{x}^{(i)}; \theta))^2
\end{aligned}$$

To  $1/m$  jsme si nakonec přimysleli, protože můžeme. MSE tedy dává smysl jako loss funkce, ale **pouze pokud má náš estimátor pevný rozptyl (tj. vlastně chybu)  $\sigma^2$** .

Describe gradient descent and compare it to stochastic (i.e., online) gradient descent and minibatch stochastic gradient descent. [5]

Pokud máme nějaký loss  $L$  a nějaká trénovací data, chceme při tréninku minimalizovat

$$J(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{p}_{\text{data}}} L(f(\mathbf{x}; \theta), y) \quad (3)$$

Což můžeme udělat v krocích pomocí tzv. **gradient descent** s learning rate  $\alpha$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta) \quad (4)$$

Druhy

- V běžném GD počítáme  $J$  a jeho gradient ze všech trénovacích dat

- V **online (stochastic) gradient descent** nasamplujeme pouze jedno dato
- V **minibatch SGD** vybereme  $m$  samplů, ze kterých poté odhadujeme střední hodnotu  $J$

Formulate conditions on the sequence of learning rates used in SGD to converge to optimum almost surely. [5]

SGD skoro jistě konverguje k optimu, pokud je naše **loss spojitá a konvexní** a zároveň pro learning raty platí

$$\forall i : \alpha_i > 0, \quad \sum_i \alpha_i = \infty, \quad \sum_i \alpha_i^2 < \infty \quad (5)$$

Tedy můžeme složením krůčků dojít kamkoli, ale zároveň musí platit  $\alpha \rightarrow 0$ . Konkrétně to *na druhou* se tam vyskytuje za MSE, říká v podstatě že "nabraná chyba bude konečná".

Write down the backpropagation algorithm. [5]

Chceme spočítat derivaci posledního vrcholu ( $u_n$ ) vzhledem ke všem předešlým vrcholům. Tím získáme derivaci loss vůči parametrům, což je to, co potřebujeme do SGD.

1. Spustíme forward propagation, kterým spočteme hodnoty všech vrcholů
2. Nastavíme  $g_n = 1$
3. Od konce počítáme  $g_i$  jako  $\sum_{j:i \in P(u^{(j)})} g^{(j)} \frac{\partial u^{(j)}}{\partial u^{(i)}}$ , využití chain rule

Write down the mini-batch SGD algorithm with momentum. Then, formulate SGD with Nesterov momentum and show the difference between them. [5]

$$\begin{aligned} g &\leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) \\ v &\leftarrow \beta v - \alpha g \\ \theta &\leftarrow \theta + v \end{aligned} \quad (6)$$

Navíc oproti SGD tam je to  $v$ , které zajišťuje roli "kudy jsme šli minule". V Nestor momentum je to pak jen trochu pozměněno, *momentum krok* se dělá *před* výpočtem gradientu, tak, aby ten samotný gradient byl přesnější.

$$\begin{aligned} \theta &\leftarrow \theta + \beta v \\ g &\leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) \\ v &\leftarrow \beta v - \alpha g \\ \theta &\leftarrow \theta - \alpha g \end{aligned} \quad (7)$$

Write down the AdaGrad algorithm and show that it tends to internally decay learning rate by a factor of  $1/t$  in step  $t$ . Then write down the RMSProp algorithm and explain how it solves the problem with the involuntary learning rate decay. [10]

$$\begin{aligned} g &\leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) \\ r &\leftarrow r + g^2 \\ \theta &\leftarrow \theta - \frac{\alpha}{\sqrt{r + \epsilon}} g \end{aligned} \quad (8)$$

Velikosti gradientu normalizujeme, aby se parametry s různými rozptyly měnily zhruba stejně. To, co máme uloženo v  $r$ , si můžeme představit zhruba jako  $\sigma^2$  jednotlivých složek, takže vydělení learning ratu  $\sqrt{r + \epsilon}$  provede normalizaci.

Pokud zůstávají gradienty dlouho stejné, tj  $g \approx g_0$ , tak po  $t$  krocích algoritmu je  $r \approx t \cdot g_0^2$ , a proto

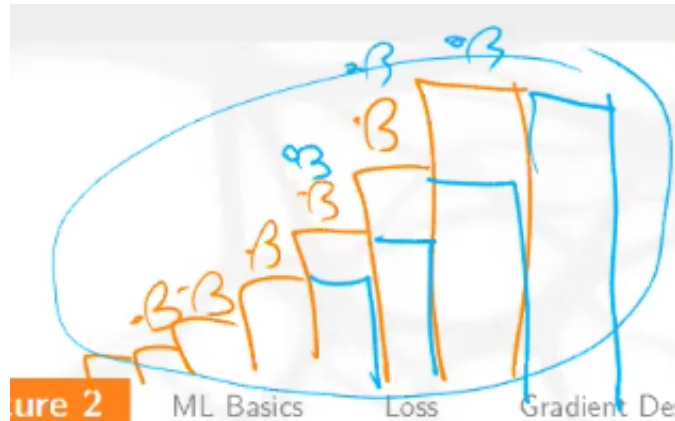
$$\frac{\alpha}{\sqrt{r + \epsilon}} \approx \frac{\alpha/\sqrt{t}}{\sqrt{g_0^2 + \epsilon/t}}, \quad (9)$$

jinými slovy, jako kdybychom learning rate škálovali  $1/\sqrt{t}$ , což zpravidla nechceme, protože to může být moc rychlé.

RMSProp funguje podobně, ale  $r$  počítáme tak, aby zhruba odpovídalo střední hodnotě poslechných  $g^2$  — počítáme exponenciální průměr posledních několika hodnot.

$$\begin{aligned} \mathbf{g} &\leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) \\ \mathbf{r} &\leftarrow \beta \mathbf{r} + (1 - \beta) \mathbf{g}^2 \\ \theta &\leftarrow \theta - \frac{\alpha}{\sqrt{\mathbf{r} + \epsilon}} \mathbf{g} \end{aligned} \quad (10)$$

Jak tento exponenciální průměr funguje je ukázáno na následujícím obrázku.



Write down the Adam algorithm. Then show why the bias-correction terms  $(1 - \beta^t)$  make the estimation of the first and second moment unbiased. [10]

Adam je spojením momentum a RMSProp.

$$\begin{aligned} \mathbf{g} &\leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) \\ t &\leftarrow t + 1 \\ \mathbf{s} &\leftarrow \beta_1 \mathbf{s} + (1 - \beta_1) \mathbf{g} \text{ (biased first moment estimate)} \\ \mathbf{r} &\leftarrow \beta_2 \mathbf{r} + (1 - \beta_2) \mathbf{g}^2 \text{ (biased second moment estimate)} \\ \hat{\mathbf{s}} &\leftarrow \mathbf{s} / (1 - \beta_1^t), \hat{\mathbf{r}} \leftarrow \mathbf{r} / (1 - \beta_2^t) \\ \theta &\leftarrow \theta - \frac{\alpha}{\sqrt{\hat{\mathbf{r}} + \epsilon}} \hat{\mathbf{s}} \end{aligned} \quad (11)$$

První moment odpovídá momentum, druhý používáme kvůli normalizaci LR, stříškové verze složí jako korekce biasů. Po  $t$  krocích totiž  $r$  vypadá jako

$$\mathbf{r}_t = (1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} \mathbf{g}_i^2 \quad (12)$$

Tedy jako bych dělal vážený průměr nějakých prvků s celkovou vahou

$$(1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} = (1 - \beta_2) \frac{1 - \beta_2^t}{1 - \beta_2} = 1 - \beta_2^t. \quad (13)$$

Jinými slovy,

$$\mathbb{E}[\mathbf{r}_t] \approx \mathbb{E}[\mathbf{g}^2] \cdot (1 - \beta_2^t) \quad (14)$$

A biasu se tedy zbavím vydělením.