

$$V_{aa} = U, W_{ab} = w$$

$$\frac{\partial L}{\partial h_{t+1}} = ?$$

$$h_{t+1} = \sigma(Ux_{t+1} + Wh_t + b)$$

المشتق

1
الف

حسب دالين شكل بايس صفياس $b=0$ و $U=1$ (نماذج)

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \times \frac{\partial}{\partial h_t} (\sigma(Ux_{t+1} + Wh_t)) =$$

$$\frac{\partial L}{\partial h_{t+1}} \times \sigma'(Ux_{t+1} + Wh_t)w$$

ب

$$\frac{\partial L}{\partial h_0} = \frac{\partial L}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_1}{\partial h_0} =$$

$$\frac{\partial L}{\partial h_t} (\sigma'(Ux_{t+1} + Wh_t)w) (\sigma'(Ux_{t-1} + Wh_{t-2})w) \dots$$

$$(\sigma'(Ux_t + Wh_0)w) \Rightarrow$$

$$\Rightarrow \frac{\partial L}{\partial h_0} = \frac{\partial L}{\partial h_t} w^t \prod_{i=0}^{t-1} \sigma'(Ux_{i+1} + Wh_i)$$

برای سنجش به دست آمدیم

الف

$$\frac{\partial L}{\partial h_0} = \frac{\partial L}{\partial h_t} w^t \prod_{i=0}^{t-1} \sigma'(Vx_i + wh_i)$$

Sigmoid = $\frac{1}{1+e^{-x}} \Rightarrow \sigma' = \sigma(x)(1-\sigma(x))$
 $= \sigma(x) - \sigma^2(x) = 0$

~~$\sigma(x) = 1 \Rightarrow x = 0$~~

الف

$$\frac{\partial L}{\partial h_0} = \frac{\partial L}{\partial h_t} w^t \prod_{i=0}^{t-1} \sigma(Vx_i + wh_i)$$

$$\sigma' = \sigma(1-\sigma) = \sigma - \sigma^2 \Rightarrow \sigma'' = \sigma(1-\sigma) - 2\sigma\sigma(1-\sigma) = 0$$

$$\sigma - \sigma^2 - 2\sigma^2 + 2\sigma^3 = 0 \Rightarrow 2\sigma^3 - 3\sigma^2 + \sigma = 0$$

$$\sigma(2\sigma^2 - 3\sigma + 1) = 0$$

جواب

$$\begin{cases} \sigma = 0 \rightarrow \text{not possible} \\ \sigma = \frac{1}{2} \rightarrow \boxed{x=0} \\ \sigma = 1 \rightarrow \text{not possible} \end{cases}$$

برای سنجش به دست آمدیم

$$\sigma(1-\sigma) = \frac{1}{2} \times (1 - \frac{1}{2}) = \frac{1}{4}$$

الف

$$\frac{\partial L}{\partial h_0} = \frac{\partial L}{\partial h_t} w^t x(\frac{1}{4})$$

$$\frac{\partial L}{\partial h_0} \rightarrow \infty$$

الف

الف

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} + \frac{\partial L}{\partial h_{t+2}} \frac{\partial h_{t+2}}{\partial h_t}$$

$$= \frac{\partial L}{\partial h_{t+1}} (\sigma(Ux_t + Wh_t + Mh_{t-1})w) +$$

$$\frac{\partial L}{\partial h_{t+2}} (\sigma(Ux_{t+1} + Wh_{t+1} + Mh_t) \cdot M)$$

حال h_t دارای ۲ جمله است h_t مشخص است که رابطه‌های داریم و تعداد جملات
رابطه‌های به این شکل از هر تیر به این است که به همین خاطر چون
تعداد زیادی جمله را داریم با هم جمع می‌کنیم پس اثر vanishing کمتری شود

درس (ج)
 gradient clipping
 در روش τ به ازای τ می‌کنیم
 در واقع یک τ در نظر می‌گیریم و حدی را تعیین می‌کنیم

در روش اندازه‌گیری همان برای است صرفاً در مورد راسی اندازه در نظر می‌گیریم
| چرا بدین اندازه جهت‌دار است؟ |

زیرا در بدین اندازه جهت‌دار عوض نمی‌شود و در بدین جهت
عوض می‌شود
 مقدار در روش اندازه‌گیری اندازه‌های عوض می‌شود
تأثیری در نتیجه ندارد

۲) الف) این روش را با مثال علم دانش آموز توضیح می دهیم که ~~یک~~ ~~از~~ ~~دانش~~ ~~آموز~~ ~~جواب~~ ~~شکل~~
ا را نداند b به a بستنی داشته باشد ~~و~~ ~~اگر~~ ~~استاد~~ ~~می~~ ~~تواند~~ ~~بگوید~~ ~~استاد~~ ~~باقی~~ ~~می~~ ~~ماند~~

~~شکل اصلی این است~~
این شکل این است که اگر دانش آموز صفا x_t را به میگوید صفا y_t حقیقت خوب
نیست و در سیدی از خود دارند $converge$ نمی شود ولی در روش
teacher forcing علاوه بر x_t خودی ~~اصلی~~ ~~در~~ ~~صفا~~ ~~مدل~~ ~~قبل~~ ~~را~~ ~~انداز~~ ~~می~~ ~~کنیم~~ ~~می~~ ~~دهیم~~ ~~که~~ ~~در~~ ~~واقع~~ ~~باید~~
روش ~~می~~ ~~کنیم~~ ~~را~~ ~~مجبوری~~ ~~کنیم~~ ~~که~~ ~~باید~~ ~~را~~ ~~بگوید~~

۲) ب) شکل اصلی این است که ~~شکل~~ ~~مندی~~ ~~که~~ ~~صفا~~ ~~خود~~ ~~مدل~~ ~~قبل~~ ~~را~~ ~~خود~~ ~~دارد~~
در حالی که در مثال سنت همین چیزی را دارد و مجبور است از خودی خودش استفاده
کند یعنی منابع این شکل را با $distribution$ $shift$ می شناسند

۲) ج) در این تکنیک یک حقیقت بر روی x با پارامتر θ در نظر می گیریم ~~و~~ ~~اگر~~ ~~استاد~~ ~~می~~ ~~تواند~~ ~~بگوید~~ ~~استاد~~ ~~باقی~~ ~~می~~ ~~ماند~~
دانش آموزی به ~~شکل~~ ~~که~~ ~~استاد~~ ~~می~~ ~~تواند~~ ~~بگوید~~ ~~استاد~~ ~~باقی~~ ~~می~~ ~~ماند~~ x یک باشد (مثلاً) بر حسب اصلی را به دانش آموز می دهیم
با ~~شکل~~ ~~که~~ ~~استاد~~ ~~می~~ ~~تواند~~ ~~بگوید~~ ~~استاد~~ ~~باقی~~ ~~می~~ ~~ماند~~ x یک باشد (مثلاً) بر حسب اصلی را به دانش آموز می دهیم
 $bias$ $exposure$

تأثیر کمی کم شود

۲
بسیار

۱) ۳ ۲

44

در روش greedy در هر مرحله max احتمال به عنوان انتخاب می شود اما در روش جستجوی

محلی این کار را انجام نمی دهیم. بلکه most likely

فرض کنیم $k=3$ است. سه مقدار بزرگتر را به عنوان انتخاب می کنیم پس به ازای هر کدام

جوابی رویم و خوبی های هر کدام را بررسی می کنیم (فرض کنیم که کاندید دانه $3 \times 4 = 20$)

پس از این 20 تا دوباره 3 تا را انتخاب می کنیم و ...

۲
۲
۲

اگر k بیش از حد بزرگ شود؛ اداسی به حساب می آید زیرا

تأثیر ابتکار محسوس کم می شود یعنی صواب آن همان انتخاب می شود که بیشترین احتمال را دارند

و با این که خوبی در نهایت درست است ولی ممکن است خوبی خیلی جالبی باشد

اگر k بیش از حد کوچک باشد؛ فضای جستجو خیلی کم می شود که باعث می شود

احتمال درست یافتن به یک دنباله مطلوب به سرت کاهش یابد

مسئله ۲: مشخص است که بافتی شود «انتظار» بیشتری داشته باشد

بعضاً خدش بدیع در بوجود میآورد.

presampling (۲) (۳) (۴) (۵)
 این هم خدش می‌دهد softmax نمونه‌گیری نمی‌دهد و top-k-
 صفای ک خدش می‌دهد با بیشترین احتمال sampling انجام می‌دهیم.

pure sampling ← $k=n$
 max (greedy) $k=1$

هرچه k کمتر باشد تنوع جلات کمتر است و جلات بی‌ارزش
 کمتری هم داریم

هرچه k بیشتر باشد تنوع بیشتری شود «انتظار» بیشتری شود و جلات بی‌ارزش
 بیشتری هم داریم.

3) حب می خواهم وقتی خدایی صفر شد صفر باشد یا نه؟ پس $h_{t-1} = 0$ هست می گویم
 $h_t = 1$ شدی یک جایی صفر شد و صفر خدای! مطابق این بودیم

x_t	h_{t-1}	h_t
0	0	1
0	1	1
1	0	0
1	1	1

→ واقعی
ایستادگی
نشان ندهم

حال خدایی هم می شود نصف h_t

h_t	y_t
0	1
1	0

به عبارت دیگر

$$h_t = \text{II}(w_2 h_{t-1} + b_2 + w_1 x_t)$$

$$y_t = \text{II}(w_3 h_t + b_3)$$

حد حالتی که $x_t = 1$ باشد خدایی h_t است $h_{t-1} = 0$

$$h_t = \text{II}(-x_t + h_{t-1} + 0.5) \geq 0$$

$$y_t = \text{II}(-h_t + 0.5) \geq 0$$

ماندگاری

$$w_1 = -1$$

$$b_2 = 0.5$$

$$w_2 = 1$$

$$b_3 = 0.5$$

$$w_3 = -1$$

$$h^{(t)} = \begin{bmatrix} h_1^{(t)} \\ h_2^{(t)} \end{bmatrix} \quad x^{(t)} = \begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

2

$r, c, c_0 \Rightarrow \text{scaler}$

$$h_1^{(t)} = 1 \Leftrightarrow x_1^{(t)} = x_2^{(t)} = 0 \Leftrightarrow \Pi(-x_1^{(t)} - x_2^{(t)} + 0.5 \geq 0) \Rightarrow w = \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix}$$

$$h_2^{(t)} = 1 \Leftrightarrow x_1^{(t)} = x_2^{(t)} = 1 \Leftrightarrow \Pi(x_1^{(t)} + x_2^{(t)} + 1.5 \geq 0) \quad b = \begin{pmatrix} 0.5 \\ -1.5 \end{pmatrix}$$

حال سبز دارم این دو $h_1^{(t)}$, $h_2^{(t)}$ را به کسب و کسب از خودی این عملیات با خروجی های صحت and شود!

$$\Pi(h_1^{(t)} + h_2^{(t)} + y^{(t-1)} - 1.5) \geq 0 \Rightarrow \begin{cases} r = 1 \\ c = -1.5 \\ c_0 = -0.5 \end{cases} \quad v = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

اگر بخواهیم این حرف بزنیم باید به جواب برسیم اگر

$$w = \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix} \quad b = \begin{pmatrix} 0.5 \\ -1.5 \end{pmatrix} \quad v = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \begin{matrix} r = 1 \\ c = -1.5 \\ c_0 = -0.5 \end{matrix}$$