نرگس کافی    ۹۹۲۱۰۷۵۳

① الف) one-to-many: این شبکه ... یک ورودی دریافت می‌کند و ... و دنباله‌ای خروجی تولید می‌شود:



هر کدام از $f$ ها می‌تواند یک چند لایه بازگشتی محض باشد. طول ثابت ... دنباله‌ای برای این نیست ...

... وزن‌های ورودی لازم. پس این بایاس و ورودی از یک time step است یا اصلا مستقل از زمان است.

کاربرد: مثلا چون تولید موسیقی دنباله‌ای کار می‌کند مانند تولید موسیقی بادیس درست آهنگین یا آوا و موسیقی ... تولید

... با یک کلمه شروع شده تولید می‌شود.

many-to-one: این شبکه دنباله‌ای دارد و ورودی را طول بلند می‌گیرد و سپس پردازش تمام دنباله دنباله یک خروجی تولید می‌شود:



همانند $f$ قبلی هر کدام یک چند لایه بازگشتی می‌باشد.

کاربرد: مثلا "برای دسته‌بندی متن sequential ... استفاده می‌شود. مثلا Sentiment Analysis که طبق نشست احساسات و دسته‌بندی می‌کند تعداد مثبت منفی درست شده باشد.

... با میزان احتمال متعلق بودن ... اضافه آنان استفاده می‌شود.

many-to-many (synced): در این حالت به ازای هر ورودی یک خروجی داشته باشیم



خروجی همزمان با ورودی ها، با تاخیر کوچکی تولید می‌شود.

تولید می‌شود.

کاربرد: این مدل کاربرد گسترده ای دارد از این مهم ترین آن ها می‌توان به تشخیص فعالیت در ویدیو اشاره کرد
که در آن به ازای هر فریم یک فعالیت پیش بینی می‌شود، مثلا پیش بینی

یا آنالیز سری های زمانی که در آن به ازای هر نمونه، یک خروجی real-time است.

many-to-many (unsynced): در این حالت، ما یک دنباله ورودی را با طول مشخص به دنباله خروجی با طول دلخواه تبدیل می‌کنیم.
طول دنباله های ورودی و خروجی می‌تواند متفاوت باشد.
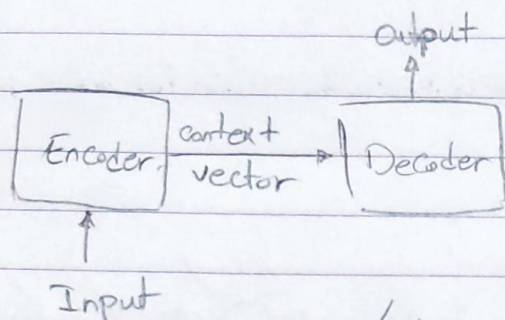


طول دنباله های ورودی و خروجی می‌تواند متفاوت باشد.

کاربرد: در ترجمه ماشینی استفاده می‌شود، به طور مثال ترجمه فارسی ~ انگلیسی، که در آن، کل جمله ورودی خوانده شده و سپس جمله انگلیسی تولید می‌شود.

به این مدل های Seq2Seq، مدل های Sequence Modelling می‌گویند که در آن‌ها طول input, output sequence
متفاوت است.

بنابراین این مدل‌ها از نوع many-to-many می‌باشند.

مدل encoder-decoder میشه گفت یه نوع معماریِ شبکس که از دوبخش encoder و decoder تشکیل شده که وظیفه هرکدوم

از اینجا ماشین، مثلا برای سرت خواهد بود.



Output

Encoder — context vector → Decoder

Input

ورودی به رشته ورودی به encoder میده این بخش نقش فشرده سازی اطلاعات رو داره که به طلاعات رو به context vector

کوچک میکنه. سپس این نمایش فشرده شده رو به عنوان ورودی به بخش decoder میده و نقش بخش decoder، استخراج اطلاعات

تولید کند.

از طرفی دو تو این معماری که برای بخش های encoder و decoder از رمزنگار و رمزگشا استفاده میکنیم از LSTM استفاده

میکنیم. دلیل استفاده از LSTM بخاطر قابلیت long-term dependency هستش که نسبت به سایر شبکه ها بهتر عمل میکنه.



ravi de vous rencontrer <END>

LSTM LSTM LSTM LSTM LSTM  →  h_t c_t  →  LSTM LSTM LSTM LSTM LSTM

nice to meet you          <start> ravi de vous rencontrer

تا اینجا متوجه میشیم این دو بخش برای encoder مراحل پیش پردازش رو شامل میشه برای cell state و hidden state رو باز

میسازه context vector است سپس بخش decoder رو میسازه و بخش decoder رو از زمان اولش با <start>

context vector شروع میشه و تا زمانی که به رمزگشایی LSTM میرسیم تا زمانی که به بخش نهایی <END>

تموم.

ج: برای بهبود عمکرد encoder کلمه مورد نظر را به decoder کمک میکنیم تا خروجی بهتری بدهد.

دقت در جهت حاصل شود و خروجی به کمک مناسب attention استفاده میکنیم attention و decoder

در این مرحله جمع وزن دار حالت و مجموع وزن ها و خروجی و مجموع وزن ها و خروجی

align میشود ...

بطور کلی این روش با اشتراک attention در cellها ...

این روش در بهبود مانند "then" آنجا می‌گیرد.



attention
distribution

attention
scores

les    pauvres    sont              ⟨start⟩

encoder                    decoder

② 

الف) برای Back P.T. Time ابتدا باید RNN را unfold و باز کنیم و هم برای هر لایه برای حساب کنیم.

حالت زمانی هر خروجی و تعداد وزن ها و بایاس های مشترک است. حال باید خروجی را با خروجی مورد نظر

مقایسه کنیم و تعداد loss را حساب کنیم، حال با برای هر loss نیز مشتق میگیریم. حال مشتق است.

مانند BP که خطا را محاسبه میکنیم، با استفاده از chain rule مشتق را بدست می‌آوریم و والی پالی

به ورودی و لایه ها می‌رسیم و تعداد وزن ها و بایاس ها آپدیت میشوند. تفاوت این با Back Propagation این در

برای حالت loss باید از بیشتر شناخت در حالت unfolding و BP کار می‌کنیم و غیره...

مشکلاتی که در شبکه های بازگشتی وجود دارد exploding gradient, vanishing gradient

است. کنترل این عنی سرعت صعود رشد دنباله کنترل نمی شود، و پدیده های اسپایک داریم، NaN می شود و بعد

از یایان اجرا می شود.

دلیل این اتفاق در شبکه های بازگشتی این است که در محاسبه رشد به صورت زیر عمل می کنیم. $\frac{\partial L}{\partial w} = \sum_{i=0}^{T} \frac{\partial h_i}{\partial w} \alpha \sum_{i=0}^{T} (\prod_{k \ge n} \frac{\partial h_i}{\partial h_{i-1}}) \frac{\partial h_k}{\partial w}$ به ازای اجرا داریم.

به عبارت دیگر بحث ① $\|\frac{\partial h_i}{\partial h_{i-1}}\|_2 < 1$ که باعث vanishing gradient می شود. و ② $\|\frac{\partial h_i}{\partial h_{i-1}}\| > 1$

که باعث exploding gradient می شود. در رشد دنباله تعداد زیادی جمله ضرب می شوند و در محاسبه رشد

long period dependency ها داشتیم. آنقدر ضرب می شود و درحالی که معذبانیم پست بیاست و NaN می شود و ایجاد مشکل

دلیل دیگر که این اتفاق در شبکه های بازگشتی می افتد این است که این تعداد زیاد ورودی می تواند اسپایک بزرگ یا شور باشد.

مسئله ایجاد دیگر در شبکه بازگشتی است. ایجاد دنباله بلند در روش حل truncated BPTT و

gradient clipping در شبکه آنها دو روش می باشد.

به منظور مدیریت بهتر در شبکه LSTM ما سه مکانیزم forget gate، input gate، و update gate استفاده می کنیم.

بیشتر ساختار cell state کنار hidden state مدیر بهتر است. در مدیریت long term dependency استفاده می شود.

حال می دانیم در درمدل هستند اندازه گرادیان راحت تأثیر گره ها می شود این گره ها ویژگی توابع حال سلول (یا طور دقیق تر

مشتق آن حال). حال اگر هرکدام از این ها ۱ کمتر شوند دچار vanishing g. و یا بزرگ تر از ۱ شوند و exploding اتفاق

می افتد. حال می دانیم در LSTM تابع حال سلول حاصل جمع دو مشتق است. بنابراین این گره ها back propagate

نتیجه شبکه عمیق نیست. وزن دقت بازگشتی هم از تابع فعل سازی forget gate است. بنابراین

این gate با تابع ثابت باشد ... activation این باشد، اگر این گرادیان دهد، می‌شود vanishing و بنابراین مقدار عقب

بنابراین این forget gate این بیشتر می‌شود دهد، می‌شود exploding. بنابراین این شبکه LSTM عمل نمی‌کند

exploding , vanishing می‌کند بنابراین بهتر.

۳ مثلاً exposure bias این است به این ظاهر میکنیم درصورتیکه آمده است در مدل teacher forcing نشون میدیم

بلبع و توصیه است راه عنوان ورودی به لایه زمانی میدیم و این اتفاق باعث ناپایداری بین train-test میشود.

exposure bias آسیب میشود. این اتفاق در ترند منجر به overfitting میشود و بنابراین باعث مدل generalize !

شبیه: مثلاً میتونیم از قبل scheduled sampling برای تا که ! exposure bias تو سعی کرده شده اند . شبیه راه دیگه

آموزش و تست توزیع نزدیک بهم نزاشته باشند . علاوه شبکه بسیار ضعیف خواست شد

چنین میگرند اگر بخواهیم داده های صحیح راه فقط دادن بعم نتوانیم اصلاً مدل را آموزش دهیم ، و باین آموزش

زمان بسیار زیادی میبرد . زیرا که داده های صحیح بخش از دورک مرحله هستن .

الف) معرفی discriminator تابع شایستگی $V(G,D)$ است. میدانیم generator بهصورت زیر تعریف میشود:

$$V(G,D) = E_{x \sim P_{data}} [\log D(x)] + E_{x \sim P_g} [\log (1-D(x))]$$

$$= \int_{x} P_{data}(x) \log (D(x)) dx + \int_{z} P_z(z) \log (1-D(g(z))) dz$$

$$= \int_{x} P_{data}(x) \log (D(x)) + P_g(x) \log (1-D(x)) dx$$

حال برای تابع شایستگی $(a,b) \in R^2 \setminus \{0,0\}$ و $y = a \log y + b \log (1-y)$ در بازه $[0,1]$ برابر است با $\frac{a}{a+b}$ مقدار بیشینه. بنابراین مقادیر شایستگی متناسب با discriminator حل کرده و مقدار بهینه $D_G^*(x) = \dfrac{P_{data}(x)}{P_{data}(x) + P_g(x)}$

$$C(G) = \max_{D} V(G,D) = E_{x \sim P_{data}} [\log D_G^*(x)] + E_{z \sim P_z} [\log(1 - D_G^*(G(z)))] \quad (\text{ب}$$

$$= E_{x \sim P_{data}} [\log D_G^*(x)] + E_{x \sim P_g} [\log(1 - D_G^*(x))]$$

$$= E_{x \sim P_{data}} \left[ \log \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \right] + E_{x \sim P_g} \left[ \log \frac{P_g(x)}{P_{data}(x) + P_g(x)} \right]$$

مقدار مینیمم تابع $C(G)$ مقدار استخراج می‌شود. با فرض اینکه $P_g = P_{data}$، مقدار آن $-\log 4$ خواهد شد.

$$P_g = P_{data} \longrightarrow D_G^*(x) = \frac{1}{2} \implies C(G) = \log \frac{1}{2} + \log \frac{1}{2} = -\log 4$$

$$E_{x \sim P_{data}} [-\log 2] + E_{x \sim P_g} [-\log 2] = -\log 4$$

$$C(G) = -\log(4) + KL \left( P_{data} \| \frac{P_{data} + P_g}{2} \right) + KL \left( P_g \| \frac{P_{data} + P_g}{2} \right) \quad (KL = \text{Kullback-Leibler})$$

$$\implies C(G) = -\log 4 + 2 \cdot JSD(P_{data} \| P_g)$$

* تابع فاصله Jensen-Shannon بین دو توزیع غیرمنفی و متقارن است و مینیمم آن صفر است. مقدار

مینیمم global برای $C(G)$ برابر $-\log 4$ است، و در صورتی به دست می‌آید که JS بین دو توزیع

آموزش، توزیع مولد یکسان برابر صفر شود.

ج) از لحاظ تئوری در بالاترین سطح آموزش، مینیمایزر JS تضمین می‌کند که $P_{data} \simeq P_g$ ($P_{data} = P_g$) باشد. همچنین

به تعبیری تابع هدف جمع می‌کنیم.

③ مسأله در موضوع یادگیری GAN ها مجزد گشتن (oscillate) در خردی تولیدکنند که مرتبط کاری موردی تقاضی یا

موضوع دیگری. اما در generator تغییر مسیر میدهد تا discriminator را گول بزند. حتی است به حالت اول

generator برگردد. در واقع generator سعی میکند بیشترین خطا را برای discriminator تولید کنند. از آنجا که generator

شروع کند، توسعه کند خردی سیاسی (را بیدیر چیرگی موردی) را بیشترین خطای برای discriminator تولید کنند به خطا آرایه موردی

هر iteration در generator نیاز... بیا را در... بیشترین طبع طبع این بعد که ته شل در discriminator که اثر در شناس یا

بسیار ساده خواهد بود که بیشترین خردی را برای discriminator نگاه میشود. در هر چرخش generator

توسعه کنند discriminator میشود آن چرخش over-optimize! نشر و discriminator میشود فقط به دیدگاه خارج ته

خارج شود. بنابراین generator حول یک تعداد خردی موردی میچرخد. به این شل میتوانند نگاه کنند شود.

④ ١) وقتی که تابع مدیریت شبکه در جهت $E_{x \sim p_g}[\log(1-D(x))]$ باشد دچار پایداری شبکه generator میشود

چون میشود میشود در جهت شد. آنگاه که discriminator بیشتر شل آنگاه حال اتفاق میشود میشود نگاه میشود میکند آرایه

برای generator سیلم و تقریبا صفر است. باعث میشود پایداری سیلی میشود شود که آن saturation بیشتر

ولی با استفاده از $E_{x \sim p_g}[\log D(x)]$ میشود discriminator را ابتدای طبع شناس کند که مدل گول میشود بیشتر شل ود

پایدار است و سریعا جهت خردی موردی

(الف) در این شرایط که discriminator به نقطه بهینه خود نزدیک می‌شود:

$$\lim_{\|D-D^*\| \to 0} \nabla_\theta E_{z \sim P_z} \left[ \log(1-D(g_\theta(z))) \right] = 0$$

$$\Rightarrow \nabla_{\theta_g} \log(1-D(G(z^i))) = 0$$

(ب) در این صورت که نزدیک به نقطه بهینه مقدار با مسئله vanishing gradient ایجاد می‌شود، میل می‌کند نخواهیم داشت.

(ج) در این صورت نزدیک ...
$$E_{z \sim P(z)} \left[ -\nabla_\theta \log D^*(g_\theta(z)) \Big|_{\theta=\theta_0} \right] = \nabla_\theta \left[ KL(P_{g_\theta} \| P_r) - 2JSD(P_g \| P_r) \right] \Big|_{\theta=\theta_0}$$

و بنابراین گرادیان جدید تابع ... شامل KL-divergence ، JS-divergense می‌باشد. بنابراین می‌شود ... شبه مدل قابل ... واقعی تر تولید ... اما ممکن است ... نوسان تولید نشان بدهد.

⑦ اگر $P_r$ و $P_g$ در توزیع اصلی ... ... باشند و در دو منیفلد (manifold) $M$، $P$ متفاوت قرار بگیرند ...

... اگر این دو منیفلد full dimension باشند و کامل بر روی هم align نشده باشند، خواهیم داشت:
$$JSD(P_r \| P_g) = \log 2$$
$$KL(P_r \| P_g) = +\infty$$
$$KL(P_g \| P_r) = +\infty$$

به همین صورت این divergenceها در صورت ... منیفلد و طوری که بتوانند ... متفاوت باشند، نشان می‌دهند. ممکن است خروجی شبه مدل ... KL divergence بهینه می‌شود و همچنین شامل

نکته: این شبکه با استفاده از gradient descent بهینه‌سازی می‌شود و تمام وزن‌ها آموزش داده می‌شوند.

اسنیزم

(۸) تمام بلوک‌های A شبیه هم هستند و قابلیت ... ... ... دارند.



## Generator Network



## Discriminator Network



$\rightarrow \quad \begin{matrix} HR \\ SR \end{matrix} ?$

generator loss function:

$$l^{SR} = \underbrace{l_x^{SR}}_{\text{Content loss}} + 10^{-3} \underbrace{l_{Gen}^{SR}}_{\text{adverserial loss}}$$

$$( l_{Gen}^{SR} = \sum_{n=1}^{N} -\lg D_{\theta_D} ( G_{\theta_G} (I^{LR}))$$

$$( l_x^{SR} = MSE(HR, SR))$$

discriminator loss function:

$$\min_{\theta_G} \max_{\theta_D} E_{I^{HR} \sim P_{train}(I^{HR})} \left[ \log D_{\theta_D}(I^{HR}) \right] + E_{I^{LR} \sim P_G(I^{LR})} \left[ \log (1 - D_{\theta_D}(G_{\theta_G}(I^{LR}))) \right]$$