



دانشگاه صنعتی شریف

دانشکده مهندسی مکانیک

شبکه‌های عصبی بازگشتی

گزارش تمرین ۳- تشریحی

رباتیک اجتماعی و شناختی

مریم کریمی جعفری، ۹۹۱۰۶۶۱۷

استاد

دکتر طاهری

بهار ۱۴۰۳

سوال ۱:

(الف)

مشکل vanishing gradient یعنی مقدار گرادیان در مرحله backpropagation آنقدر کوچک شود که بعد از آپدیت پارامترها مقدار آن‌ها عملاً تفاوتی نکند. حال بررسی می‌کنیم که LSTM چگونه این مشکل را حل کرده است.

Backpropagation through time in LSTMs

گرادیان در شبکه‌ها عصبی بازگشتی به صورت زیر محاسبه می‌شود:

$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

اگر بخواهیم ناپدید شدن گرادیان رخ دهد باید مقدار بالا به ازای تمامی T ها صفر شود زیرا یک سری است که باید به صفر همگرا شود.

$$(S_1, S_2, S_3, \dots)$$

$$S_n = \sum_{t=1}^n \frac{\partial E_t}{\partial W}$$

مقدار گرادیان برای k time step به صورت زیر است:

$$\frac{\partial E_k}{\partial W} = \frac{\partial E_k}{\partial h_k} \frac{\partial h_k}{\partial c_k} \dots \frac{\partial c_2}{\partial c_1} \frac{\partial c_1}{\partial W}$$

$$= \frac{\partial E_k}{\partial h_k} \frac{\partial h_k}{\partial c_k} \left(\prod_{t=2}^k \frac{\partial c_t}{\partial c_{t-1}} \right) \frac{\partial c_1}{\partial W}$$

این مقدار عبارت زیر است که در RNN ها برای مقادیر بالای k به دلیل وجود تابع \tanh ، به صفر می‌گردد:

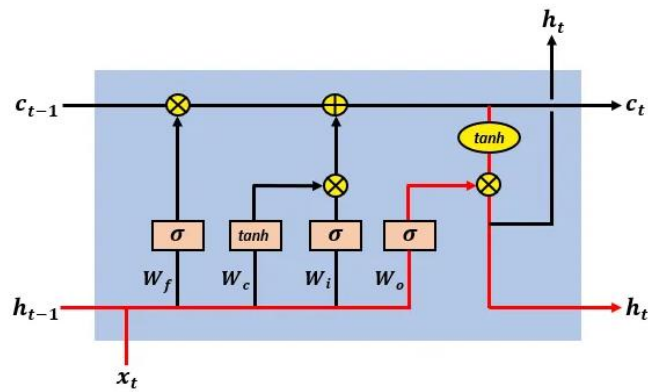
$$\prod_{t=2}^k \frac{\partial c_t}{\partial c_{t-1}}$$

مقدار c_t برابر است با:

$$c_t = c_{t-1} \otimes \sigma(W_f \cdot [h_{t-1}, x_t]) \oplus$$

$$\tanh(W_c \cdot [h_{t-1}, x_t]) \otimes \sigma(W_i \cdot [h_{t-1}, x_t])$$

$$c_t = c_{t-1} \otimes f_t \oplus \tilde{c}_t \otimes i_t$$



این c_t تابع مقادیری است که باید برای محاسبه مشتقش، از قانون مشتق زنجیره‌ای استفاده کرد:

$$\begin{aligned} \frac{\partial c_t}{\partial c_{t-1}} &= \frac{\partial}{\partial c_{t-1}} [c_{t-1} \otimes f_t \oplus \tilde{c}_t \otimes i_t] \\ &= \frac{\partial}{\partial c_{t-1}} [c_{t-1} \otimes f_t] + \frac{\partial}{\partial c_{t-1}} [\tilde{c}_t \otimes i_t] \\ &= \frac{\partial f_t}{\partial c_{t-1}} \cdot c_{t-1} + \frac{\partial c_{t-1}}{\partial c_{t-1}} \cdot f_t + \frac{\partial i_t}{\partial c_{t-1}} \cdot \tilde{c}_t + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} \cdot i_t \end{aligned}$$

در ادامه با محاسبه مشتق‌ها داریم:

$$\begin{aligned}\frac{\partial c_t}{\partial c_{t-1}} &= \sigma'(W_f \cdot [h_{t-1}, x_t]) \cdot W_f \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot c_{t-1} \\ &+ f_t \\ &+ \sigma'(W_i \cdot [h_{t-1}, x_t]) \cdot W_i \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot \tilde{c}_t \\ &+ \sigma'(W_c \cdot [h_{t-1}, x_t]) \cdot W_c \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot i_t\end{aligned}$$

اگر در نظر بگیریم:

$$\begin{aligned}A_t &= \sigma'(W_f \cdot [h_{t-1}, x_t]) \cdot W_f \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot c_{t-1} \\ B_t &= f_t \\ C_t &= \sigma'(W_i \cdot [h_{t-1}, x_t]) \cdot W_i \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot \tilde{c}_t \\ D_t &= \sigma'(W_c \cdot [h_{t-1}, x_t]) \cdot W_c \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot i_t\end{aligned}$$

با جایگذاری در فرمول مشتق c_t خواهیم داشت:

$$\frac{\partial c_t}{\partial c_{t-1}} = A_t + B_t + C_t + D_t$$

در نهایت فرمول گرادیان به صورت زیر خواهد شد:

$$\frac{\partial E_k}{\partial W} = \frac{\partial E_k}{\partial h_k} \frac{\partial h_k}{\partial c_k} \left(\prod_{t=2}^k [A_t + B_t + C_t + D_t] \right) \frac{\partial c_1}{\partial W}$$

مشکل ناپدید شدن گرادیان به دلیل مشتق تابع \tanh بود. در هر یک از عبارات A_t ، C_t و D_t عبارت مشتق \tanh وجود دارد. اما در عبارت B_t که همان forget gate است، این عبارت وجود ندارد. پس دروازه فراموشی در حل مشکل vanishing gradient نقش دارد.

منابع:

- Arbel, Nir (2018). How LSTM networks solve the problem of vanishing gradients. <https://medium.datadriveninvestor.com/how-do-lstm-networks-solve-the-problem-of-vanishing-gradients-a6784971a577>

(ب)

گیت فراموشی: با توجه به ورودی جدید و h لایه قبل و با کمک تابع فعالسازی سیگموید، مشخص می‌کند که چه اطلاعاتی اط قبل که داخل cell state بود مفید هستند و می‌توانند باقی بمانند. در این مرحله اطلاعاتی که دیگر در این زمان نیازی به آن‌ها نیست حذف می‌شوند. این مرحله از اشباع شدن حافظه با اطلاعات غیر مرتبط جلوگیری می‌کند.

گیت ورودی: با توجه به ورودی جدید و h لایه قبل و با کمک تابع فعالسازی سیگموید، مشخص می‌کند که چه مقدار از این اطلاعات جدید برای ذخیره شدن در cell state مفید هستند. اطلاعات جدید از تابع فعالسازی tanh عبور کرده‌اند. این اطلاعات به اطلاعات گذشته از گیت فراموشی اضافه می‌شوند.

گیت خروجی: با توجه به ورودی و h لایه قبل مشخص می‌کند که چه مقدار از اطلاعات داخل cell state بعد از عبور از دو دروازه قبلی، برای خروجی hidden state یا همان h همین لایه مناسب هستند. به طور خلاصه موارد زیر توسط این گیت‌ها اتفاق می‌افتد:

فیلتر کردن اطلاعات: گیت‌های فراموشی و ورودی به LSTM اجازه می‌دهند تا اطلاعات نامربوط را فیلتر کرده و بر روی اطلاعات مرتبط تمرکز کنند.

جمع‌آوری اطلاعات: گیت ورودی اطلاعات جدید را به طور گزینشی به حافظه اضافه می‌کند و به LSTM کمک می‌کند تا الگوهای long-range را در طول زمان ایجاد کند.

تمرکز بر الگوهای long-range: گیت خروجی بخش‌های مرتبط از اطلاعات ذخیره‌شده در حافظه را فیلتر می‌کند و به LSTM اجازه می‌دهد تا بر ارتباطات long-range تمرکز کند.

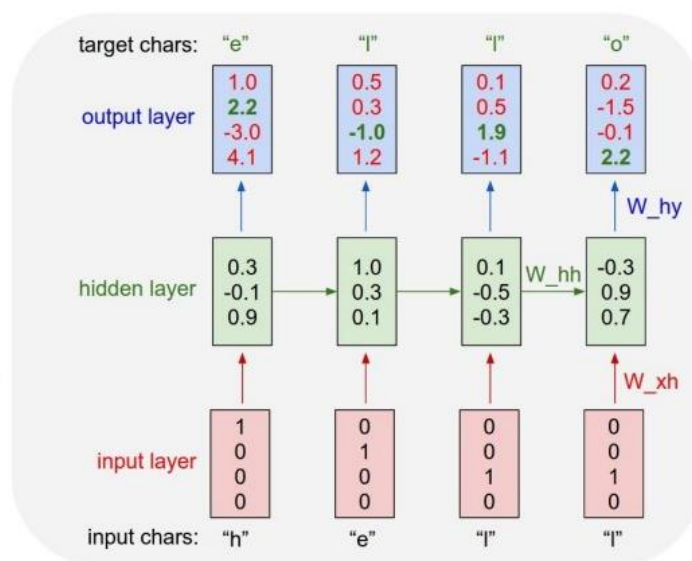
(پ)

اولین راهی که به ذهن می‌رسد این است که داده‌های ناقص به کلی حذف شوند که البته این کار ممکن است باعث کاهش حجم اطلاعات و یا از بین رفتن اطلاعات شود. یا اینکه باید داده‌های حذف شده و یا ناقص را با مقادیری جایگزین کرد (imputation). یکی از این روش‌ها این است که با میان‌یابی (interpolation) جاهای خالی را پر کرد. خود میان‌یابی می‌تواند خطی و یا غیرخطی باشد. در روشی دیگر داده‌های ناقص را با میانگین feature مربوطه پر می‌کنند. یا اینکه می‌توان از شبکه‌های عصبی برای تکمیل کردن داده‌ها استفاده کرد. یعنی

شبکه عصبی مثل حتی یک MLP ساده train کنیم که بتواند به بهترین نحو داده‌های ناقص را کامل کند. autoencoderها و یا شبکه‌های generative می‌توانند در این مواقع به کار روند.

(ت)

Teacher forcing مبحثی نیست که از قبل با آن آشنا نباشیم. در کی از اسلایدهای درس مثال زیر را داریم:



این مثال برای پیش بینی کلمه 'hello' به کار برده شده است. ورودی را به صورت برداری شامل ۴ عدد در نظر بگیرید که هر عدد به ترتیب متعلق به حروف h, e, l, و o است. اولین ورودی h است و خروجی توالی (sequence) اول باید e باشد. اما با توجه به اعداد خروجی، مدل o را پیش بینی کرده است زیرا مقدار آن یعنی ۴.۱، بیشترین مقدار است. teacher forcing بیان می‌کند که در مرحله آموزش، به جای اینکه خروجی هر مرحله را (که در اینجا در مرحله اول o است) به عنوان ورودی مرحله بعد بدهیم، مقدار صحیح آن یعنی همان e را (در این مرحله از مثال) بدهیم که در تصویر بالا مشخص است.

مزایا:

Faster Convergence

با این روش فرایند آموزش مدل سریع‌تر اتفاق می‌افتد. یعنی مدل زودتر به دقت مطلوب می‌رسد.

Reduced Error Propagation

این روش انتشار خطا در طول فرایند را کاهش می‌دهد. یعنی از بزرگ و انباشته شدن خطایی که در مراحل اولیه آموزش وجود دارد، جلوگیری می‌کند. زیرا دیگر خطای زیادی در مراحل اولیه وجود ندارد. برای درک مفهوم انتشار خطا گلوله برفی را در نظر بگیرید که از بالای یک کوه برفی به پایین سر بخورد. این گلوله برفی در راه خود به دلیل تماس با برف‌های بیشتر، بزرگ و بزرگ‌تر می‌شود.

Stable Training

این روش با حفظ یک رابطه‌ی ورودی-خروجی شفاف و تعریف‌شده، فرآیند آموزشی باثبات و ثابتی را فراهم می‌کند.

Explicit Supervision

با این روش وقتی مدل به جواب صحیح دسترسی داشته باشد، بهتر می‌تواند در مواقعی که خروجی از الگوهای خاص پیروی می‌کند، عمل کند. برای مثال در تسک‌هایی مثل ترجمه ممکن است مدل با دانستن خروجی واقعی بتواند مواردی مثل گرامر یا لغات زبان مقصد را بیاموزد.

Controlled Exploration

مدل باید بتواند برای یافتن خروجی مطلوب، exploration‌های درست انجام دهد. این روش کاوش‌های مدل را کنترل و هدایت شده می‌کند؛ به طوری که قدرت generalization مدل برای مثال‌های unseen بیشتر شده و خروجی‌های معنادارتری تولید می‌کند.

Easier Evaluation

در نهایت این روش می‌تواند باعث بهتر و آسان‌تر انجام شدن ارزیابی‌ها شود زیرا که خروجی که با کمک آن می‌توان error را محاسبه کرد مشخص است.

معایب:

Exposure Bias

در این پدیده، مدل ممکن است بیش از حد به سیگنال‌های ارائه شده توسط معلم وابسته شود و در طول استنتاج (یعنی بدون اجبار معلم) با مشکل مواجه شود. استنتاج (inference)، فرآیندی است که در آن مدل از دانش آموخته‌ی خود برای تولید خروجی‌های جدید بر اساس ورودی‌های جدید استفاده می‌کند.

Mismatch Between Training and Inference

اگر مدل در طول آموزش به شدت به Teacher Forcing تکیه کند، ممکن است در طول فرایند inference عملکرد خوبی نداشته باشد. در استنتاج، مدل باید به طور مرحله به مرحله توالی را تولید کند، بدون اینکه به خروجی واقعی دسترسی داشته باشد. به عبارت دیگر، ناسازگاری بین نحوه آموزش مدل و نحوه عملکرد آن در دنیای واقعی وجود دارد. در طول آموزش، مدل با استفاده از teacher forcing، ورودی "ایده‌آل" را در هر مرحله دریافت می‌کند. اما در استنتاج، مدل باید بر اساس ورودی که دریافت می‌کند، خروجی بعدی را به طور مستقل پیش‌بینی کند. این عدم تطابق می‌تواند منجر به کاهش عملکرد مدل در هنگام استفاده واقعی شود.

Lack of Real-World Noise

در واقعیت داده‌ها ممکن است نویزی باشند. مدلی که با teacher forcing آموزش دیده است ممکن است نتواند به خوبی داده‌های نویزی را handle کند.

Limited Exploration

از جمله مزایای این بود که explorationهای مدل جهت‌دار می‌شود. اما باید توجه کرد که در عین حال این کاوش‌ها محدود می‌شوند و مدل نمی‌تواند حتی کمی آزادی و خلاقیت فراتر از چیزی که به آن آموزش داده شده است، داشته باشد.

Incomplete Training Data

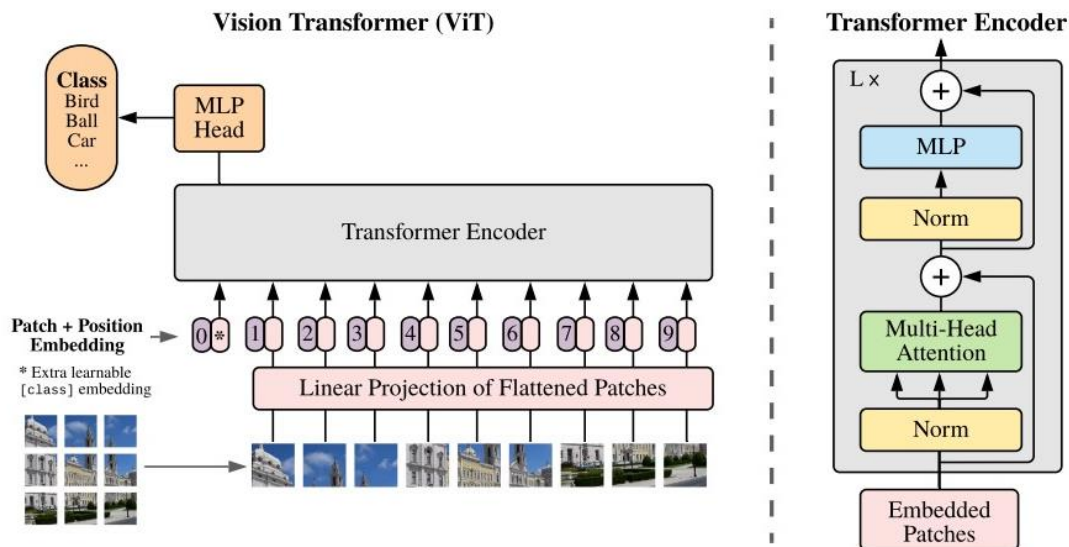
وقتی که قرار باشد مدل خیلی به ورودی درست تکیه کند، ممکن است در مواقعی که دیتا ناقص باشد نتواند خروجی مطلوبی تولید کند. این مشکل در تسک‌های ترجمه اتفاق می‌افتد.

Resource-Intensive

در نهایت استفاده از teacher forcing نیازمند به داشتن دیتای غنی است که جمع‌آوری این میزان داده نیز زمان زیادی می‌طلبد.

منابع:

- Wong, Wanshun (2019). What is Teacher Forcing?. <https://towardsdatascience.com/what-is-teacher-forcing-3da6217fed1c?gi=21e413ecd815>
- Van Otten, Neri (2023). Teacher Forcing In Recurrent Neural Networks (RNNs): An Advanced Concept Made Simple. <https://spotintelligence.com/2023/10/12/teacher-forcing-in-recurrent-neural-networks-rnns-an-advanced-concept-made-simple/>



فرض می‌کنیم ابعاد تصویر $H \times W$ با C کانال است. vision transformer تقریباً از همان ترنسفورمرهای عادی تشکیل شده است اما تصاویر را باید به فرم یک بردار برای ورودی آن‌ها تبدیل کرد. هر تصویر به دسته‌هایی تبدیل می‌شوند. این دسته‌ها (patches) هر کدام یک تصویر با ابعاد پایین‌تری هستند. برای مثال تصویر $H \times W \times C$ ما به N تصویر با ابعاد $P \times P \times C$ تقسیم می‌شود. در تصور بالا این مرحله مشخص است.

در ادامه هر تصویر flatten می‌شود تا به یک بردار یک‌بعدی تبدیل شود. اما از آنجایی که همچنان سائز تصویر بزرگ است، با فرض اینکه ورودی transformer برداری با ابعاد $1 \times D$ است، تصویر توسط linear projection مپ می‌شود که به آن patch embedding می‌گویند. حال این بردار می‌تواند به عنوان ورودی به ترنسفورمر داده شود. اما در کنار بردار embed شده تصویر، برداری تحت عنوان position embedding با همان ابعاد نیز با آن جمع شده و سپس به ترنسفورمر وارد می‌شوند. این بردار نمادی از محل قرارگیری تصویر patch از تصویر اصلی است.

حال ما N ورودی داریم که در تصویر بالا این تعداد است. باید به این نکته توجه کرد که در ویژن ترنسفرها در کنار تصاویر، یک توکن اضافه‌تر در ابتدا به عنوان توکن صفر نیز به شبکه داده می‌شود. این توکن همان [class] مشابح توکن [cls] در BERT است. این توکن در پایان شبکه نمادی از دسته‌ای است که تصویر به آن تعلق دارد؛ یعنی از این توکن در آخر برای classification تصویر استفاده می‌شود. مقدار این توکن یکی از learnable parameterها است که باید train شود. در لایه اول با نام z_0^{class} و در لایه آخر با نام z_L^0 شناخته می‌شود. مقدار در آخر به یک شبکه MLP با تابع فعالسازای GELU داده می‌شود تا فرایند

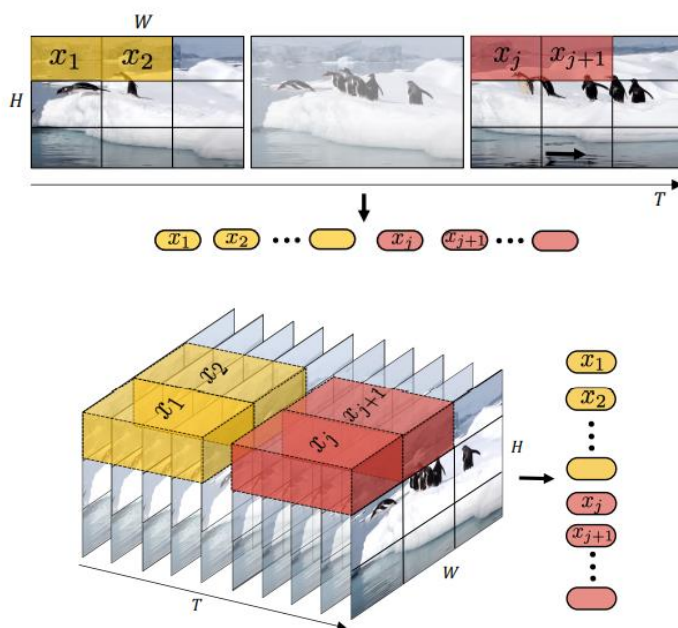
classification انجام گیرد. یعنی فرایند classification تنها به توکن ابتدایی که جزو تصویر نبود نیاز دارد. اما باید توجه کرد که این توکن خود به سایر توکن‌ها وابسته است.

در نهایت یکی از تفاوت‌های ساختار خود ترنسفورمر به ترنسفورمر عادی مقاله *attention is all you need* وجود فرایندهای نرمالیزیشن در شبکه است (LayerNorm). نرمالیزیشن قبل از هر بلوکی انجام می‌شود که در تصویر صفحه قبل به خوبی مشخص است.

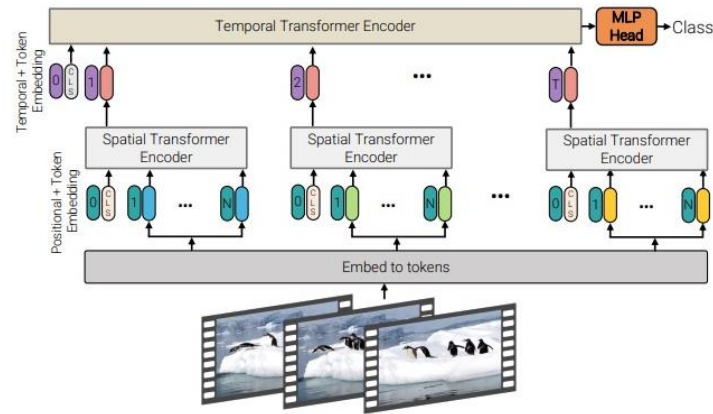
$$\begin{aligned} \mathbf{z}_0 &= [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, & \mathbf{E} &\in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \\ \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, & \ell &= 1 \dots L \\ \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, & \ell &= 1 \dots L \\ \mathbf{y} &= \text{LN}(\mathbf{z}_L^0) \end{aligned}$$

چگونه از *vision transformer* برای ویدیو استفاده کنیم؟

در قدم اول باید توجه کرد که در ویدیوها علاوه بر مکان، زمان نیز باید مدنظر قرار گیرد. ما باید به دنبال راه حلی باشیم که بتوان هر فیلم را به توکن‌هایی برای ورودی یک *transformer* از هر نوعی، تبدیل کرد. یک از راه‌ها این است که هر فریم به عنوان یک تصویر به *patch*‌هایی تقسیم شود. برای مثال اگر n_t فریم داشته باشیم که هر کدام $n_h * n_w$ تعداد *patch* تولید کند، تعداد کل *patch*‌ها $n_t * n_h * n_w$ خواهد بود. حال هر *patch* تحت *linear projection* به توکن‌های ورودی *vision transformer* تبدیل شود. راه حل دوم این است که *patch*‌ها را در عمق ایجاد کنیم.



یک دیگر از راه‌هایی که وجود دارد این است که در ابتدا هر فریم از ویدیو را به ViT بدهیم سپس خروجی آن‌ها را به یک شبکه مثل RNN ، $LSTM$ و یا حتی یک $transformer$ دیگر بدهیم تا بتواند توالی زمانی را بررسی کند.



منابع:

- Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv preprint arXiv:2010.11929* (2020).
- Arnab, Anurag, et al. "Vivit: A video vision transformer." *Proceedings of the IEEE/CVF international conference on computer vision*. 2021.

سوال ۲:

سوال ۲ :

(الف) به طریقی برای هر t داریم :

$$h_t = (u_t + h_{t-1})(w)$$

(ب) روابط انتشار به جلوی این شبکه به صورت زیر است :

$$\begin{aligned} p &= u_0 + h_0 & r &= u_1 + q & t &= u_3 + s \\ q &= wp & s &= wr & \hat{y} &= wt \end{aligned}$$

$$\Rightarrow \hat{y} = w \left[u_3 + w(u_1 + w(u_0 + h_0)) \right]$$

ب. د. ت

$$\frac{\partial \text{Loss}}{\partial w} = \frac{\partial \text{Loss}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w}$$

$$\frac{\partial \text{Loss}}{\partial \hat{y}} = -2(y - \hat{y}) = 2(\hat{y} - y)$$

حلقه به کاسه y می بیند داریم :

$$\frac{\partial \hat{y}}{\partial w} = \frac{\partial w}{\partial w} t + w \frac{\partial t}{\partial w} = t + w \frac{\partial t}{\partial w}$$

$$\frac{\partial t}{\partial w} = \frac{\partial s}{\partial w} \Rightarrow \frac{\partial \hat{y}}{\partial w} = t + w \frac{\partial s}{\partial w} = t + w \left[r + w \frac{\partial r}{\partial w} \right]$$

$$\frac{\partial r}{\partial w} = \frac{\partial q}{\partial w} = p \Rightarrow \frac{\partial \hat{y}}{\partial w} = t + w \left[r + wp \right]$$

$$\Rightarrow \frac{\partial \hat{y}}{\partial w} = u_3 + w(u_1 + w(u_0 + h_0)) + w \left[(u_1 + w(u_0 + h_0)) + w(u_0 + h_0) \right] =$$

$$u_3 + 2wu_1 + 3w^2(u_0 + h_0)$$

$$\Rightarrow \frac{\partial \text{Loss}}{\partial w} = 2(\hat{y} - y) \left[u_3 + 2wu_1 + 3w^2(u_0 + h_0) \right] \textcircled{A}$$

مشتق جزئی روی x_3 (اول):

$$\left(\frac{\partial \text{Loss}}{\partial w}\right)_1 = 2(\hat{y} - y)(t) = 2(\hat{y} - y) \left[x_3 + w(x_1 + w(x_0 + h_0)) \right]$$

مشتق جزئی روی w (دوم):

$$\left(\frac{\partial \text{Loss}}{\partial w}\right)_2 = 2(\hat{y} - y)(wr) = 2(\hat{y} - y)(w)(x_1 + w(x_0 + h_0))$$

مشتق جزئی روی w^2 (سوم):

$$\left(\frac{\partial \text{Loss}}{\partial w}\right)_3 = 2(\hat{y} - y)(w^2 p) = 2(\hat{y} - y)(w^2)(x_0 + h_0)$$

صحنه‌ها که در رابطه \otimes دیده می‌شوند، وزن‌ها بتوان 2 نیز ظاهر شده‌اند. اگر این وزن‌ها خیلی کوچک باشند، مقدار $\frac{\partial \text{Loss}}{\partial w}$ خیلی ناچیز تغییر خواهد کرد و حتی ممکن است به صفر میل کرده و مشکل ناپدید شدن گرایان $\text{vanishing gradient}$ را به وجود آورد. از طرفی اگر این وزن‌ها خیلی بزرگ باشند، از آنجایی که بتوان 2 نیز می‌رسد، مقدار گرایان ممکن است به بی‌نهایت میل کند و مشکل انفجار گرایان $\text{exploding gradient}$ را به وجود آورد.

سوال ۳:

6 / May 2024

دوشنبه | اردیبهشت

۲۷ شوال ۱۴۴۵

سوال ۳: تمام حالت هایی که در بیت می توانند داشته باشند به همراه حالت هایی که رقم carry می توانند داشته باشند، در جدول زیر آورده اند؟ رقم carry نشان دهنده این است که جمع درست ۲ شده است یا نه:

x_1	x_2	carry	y
0	0	0	0
0	1	0	1
1	0	0	0
1	1	0	0 + 1 carry
0	0	1	1
0	1	1	0 + 1 carry
1	0	1	0 + 1 carry
1	1	1	1 + 1 carry

لیست همان حالت قبل

الحال جدیدی که اضافه می شود

باز به جدول بالا ما در کل ۸ حالت برای carry داریم :
 $\begin{cases} C: \text{carry} \Rightarrow 1 \\ NC: \text{no carry} \Rightarrow 0 \end{cases}$
 اما چهار حالت برای x ما داریم : $[0], [1], [0], [1]$

همین ضربی و هم کادو حالت ۰، ۱، نشان می دهد. اما ضرب از این روی تواند NC، C باشد؟

پس در کل ۸ حالت داریم : $(NC, 0), (C, 0), (NC, 1), (C, 1)$

جدول قبل $(NC, 0)$ یعنی هیچ رقم carry از state قبل نداشته باشیم و ۰ می باشد
 در ضربی کادو، هر بوک RNN می تواند یکی از ۴ حالت بالا باشد و همین می تواند از یکی از ۴ حالت بالا به بوک RNN فعلی برسد و یعنی بوک های RNN قبل هم می توانند حالت های

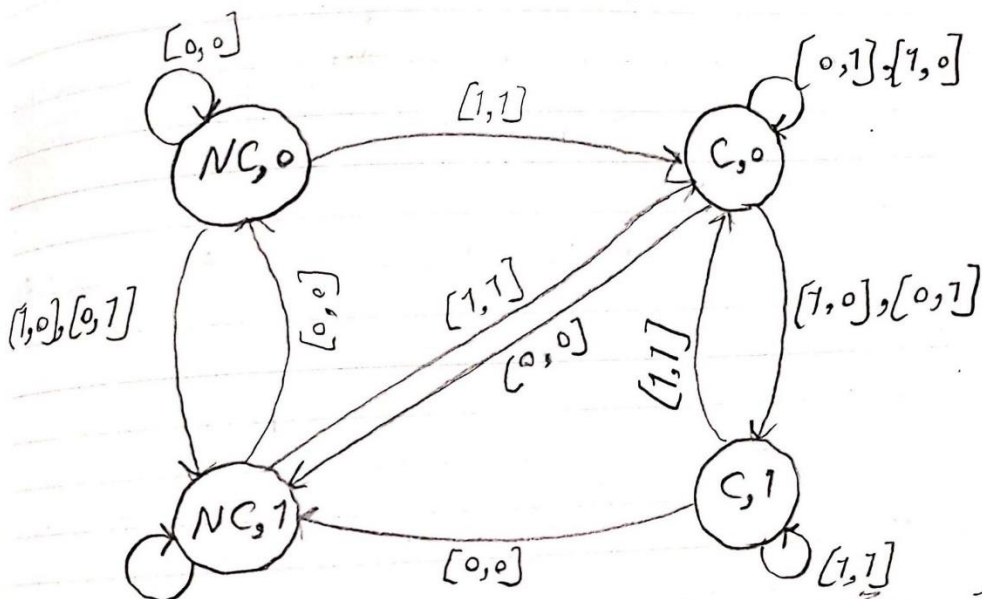
سه شنبه | اردیبهشت
دسته باستان که به نتیجه نهایی ختم شود. حالت های زیر در کل به بودی است:

From \ to	(NC, 0)	(NC, 1)	(C, 0)	(C, 1)
(NC, 0)	[0, 0]	[1, 0] و [0, 1]	[1, 1]	—
(NC, 1)	[0, 0]	[1, 0] و [0, 1]	[1, 1]	—
(C, 0)	—	[0, 0]	[1, 0] و [0, 1]	[1, 1]
(C, 1)	—	[0, 0]	[1, 0] و [0, 1]	[1, 1]

این جدول بیان می کند که برای مثال در state قبلی در حالت (NC, 0) بودیم، یعنی بوک
RNN قبل با $z_t = 0$ را به دست آورده بود و رقم $carry_t = 0$ را یاد کرده و به بوک فعلی می فرستد.
حال برای اینکه بوک فعلی نیز $z_t = 0$ را به دست آورد، باید $carry_t = 0$ را یاد کند برای بوک بعدی ما به
بیت های ورودی یا همان $[x_1, x_2]$ یا $[0, 0]$ یا $[0, 1]$ یا $[1, 0]$ یا $[1, 1]$ بستگی دارد. بستگی از حالتی که
رقم $carry$ نداریم، یعنی توان به حالتی که رقم $carry$ داشته باشیم، همچنین $z_t = 1$ و $z_t = 0$
داریم.

برای مثال اگر بخواهیم روش ستاره را توضیح دهیم. بوک RNN قبلی $z_t = 0$ را در ضربی به دست
آورده و رقم $carry_t = 0$ را یاد کرده تا به بوک فعلی بماند. برای اینکه برای بوک بعدی $z_t = 1$
به دست شود، همان یک رقم $carry_t = 1$ داشته باشیم، باید بیت های ورودی عدد 1 داشته باشیم.
هم عدد 1، خود یک $z_t = 0$ و $carry_t = 1$ ایجاد می کند، از طرفی ما $carry_t = 1$ از بوک قبلی داریم

چهارشنبه | اردیبهشت
جدول صفات قبل را به شکل ریاضی زیر بنویسید و می‌توان نشان داد:



تصویر قبلی از این گراف را افزوده است.

گفتیم تعداد نمره‌های لایه مخفی ۳ است. فردی لایه مخفی ۱ همان بردار h_1 است و می‌تواند در تمام یک از حالت‌های $(C,1), \dots, (NC,0)$ هستند. پس هر یک از این ۴ حالت را با یک بردار مشخص می‌کنیم که در این آن صفات یک هستند. چنین حالت مختلف داریم که این حالت‌های زیر را امتحان کرده‌ایم:

$$(NC,0) \rightarrow [0,0,0,0]$$

$$(NC,1) \rightarrow [0,1,0,0]$$

$$(C,0) \rightarrow [1,1,1,1]$$

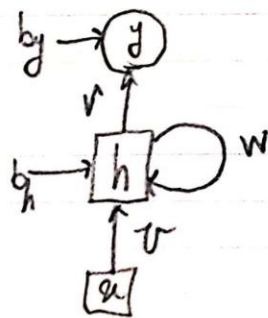
$$(C,1) \rightarrow [0,1,1,1]$$

$$\begin{matrix} h_1 & & h_2 & & h_3 \\ & \downarrow & & \downarrow & \\ & h_2 & & h_3 & \end{matrix}$$

حال با توجه به notation ها اینجا، جدولی برای input و output های این شبکه می‌نویسیم تا وزن‌ها را بیابیم. جدول در صفحه بعد آمده است.

	$h_1(t-1)$	$h_2(t-1)$	$h_3(t-1)$	$x_1(t)$	$x_2(t)$	$g(t)$	$h_1(t)$	$h_2(t)$	$h_3(t)$
1	0	1	0	1	0	1	0	1	0
2	0	1	0	1	1	0	1	1	1
3	1	1	1	1	0	0	1	1	1
4	1	1	1	0	0	1	0	1	0
5	0	1	0	0	1	1	0	1	0
6	0	1	0	1	1	0	1	1	1
7	1	1	1	0	0	1	0	1	0

$$\begin{aligned}
 t_{21} &\Rightarrow [NC, 1] \rightsquigarrow [NC, 1] \\
 t_{22} &\Rightarrow [NC, 1] \rightsquigarrow [C, 0] \\
 t_{23} &\Rightarrow [C, 0] \rightsquigarrow [C, 0] \\
 t_{24} &\Rightarrow [C, 0] \rightsquigarrow [NC, 1] \\
 t_{25} &\Rightarrow [NC, 1] \rightsquigarrow [NC, 1] \\
 t_{26} &\Rightarrow [NC, 1] \rightsquigarrow [C, 0] \\
 t_{27} &\Rightarrow [C, 0] \rightsquigarrow [NC, 1]
 \end{aligned}$$



$$\begin{cases}
 h_t \rightarrow 3 \times 1 \\
 x_t \rightarrow 2 \times 1
 \end{cases}$$

$$\begin{cases}
 W \rightarrow 3 \times 3 \\
 v \rightarrow 3 \times 2 \\
 b_h \rightarrow 3 \times 1 \\
 u \rightarrow 1 \times 3 \\
 b_y \rightarrow 1 \times 1
 \end{cases}$$

$$h_t = f(W h_{t-1} + v x_t + b_h)$$

$$g_t = f(u h_t + b_y)$$

نمایان از روشی که به صورت زیر است به این شکل در نظر گرفته می شود:

حال سعی می‌کنیم رابطه بین g_t و h_t را بیابیم؛ اگر به صورت g_t و h_t نگاه کنیم، متوجه می‌شویم که $h_2(t)$ همواره ۱ است اما همگامه h_1 و h_3 همگرا باشند، $g_t = 1$ شده است، برعکس پس می‌توانیم رابطه‌ای به شکل زیر بنویسیم که با این شرط نیز باید همراهش کنیم:

$$a_t = (1 - h_1(t)) + h_2(t) + (1 - h_3(t)) + b \quad \text{و} \quad f(a_t) = g_t$$

حالا می‌خواهیم t ها داریم:

t	a_t	g_t
1	$3+b$	1
2	$1+b$	0
3	$1+b$	0
4	$3+b$	1
5	$3+b$	1
6	$1+b$	0
7	$3+b$	1

$$g_t = \begin{cases} 1 & \text{if } a_t \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

در تمامی حالات که $g_t = 1$ شده است $a_t = 3+b$ باشد:

$$3+b > 0 \Rightarrow b > -3 \Rightarrow \boxed{b = -2,8}$$

$$\Rightarrow a_t = 1 - h_1 + h_2 + 1 - h_3 - 2,8$$

$$\Rightarrow a_t = -h_1 + h_2 - h_3 - 0,2 \Rightarrow g_t = f(-h_1 + h_2 - h_3 - 0,2)$$

$$\Rightarrow g_t = f\left(\begin{bmatrix} -1 & 1 & -1 \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} - 0,2\right) \Rightarrow \boxed{V = \begin{bmatrix} -1 & 1 & -1 \end{bmatrix}} \\ \boxed{b_g = -0,2}$$

در ادامه به بررسی فضای W ، U می‌رویم، مشابه حالت قبل، تمامی h_2 ها ۱ هستند، همچنین $h_1 = h_3$ است. پس برای راحتی در بیان بردار h و g را به صورت h_3 و h_2 و h_1 بنویسیم.

اگر $x_1(t)$ و $x_2(t)$ و $h_3(t-1)$ را با هم جمع کنیم، مسئله می شود که جمع این سه رابطه را
 y و $h_1(t)$ دارد به طوریکه جمع این سه به متغیری می شود که می چینی print شود و به رقم
 carry برای state بعد خواهیم داشت. یا به عبارتی دیگر جمع آنها بشود 1 بشود 1، y و $h_1(t)=0$
 خواهیم شد و اگر جمع آن ها 2 بشود، $y=0$ و $h_1(t)=1$ خواهد شد. پس خواهیم داشت:

$$h_1(t) = f(d_t), \quad d_t = h_3(t-1) + x_1(t) + x_2(t) + b$$

برای t داریم:

t	d_t	$h_1(t)$
1	$1+b$	0
2	$2+b$	1
3	$2+b$	1
4	$1+b$	0
5	$1+b$	0
6	$2+b$	1
7	$1+b$	0

$h_1(t) = \begin{cases} 1 & \text{if } d_t > 0 \\ 0 & \text{otherwise} \end{cases}$
 به طوریکه d_t برای t می شود:
 $2+b > 0 \Rightarrow b > -2 \Rightarrow \boxed{b \geq -1.8}$

$$\Rightarrow h_1(t) = f(x_1(t) + x_2(t) + h_3(t-1) - 1.8)$$

$$\Rightarrow h_3(t) = h_1(t)$$

$$\Rightarrow h_1(t) = f\left(\begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} h_1(t-1) \\ h_2(t-1) \\ h_3(t-1) \end{bmatrix} + \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} - 1.8\right) = h_3(t)$$

معادله $h_1(t)$ برای $h_2(t)$ خواهیم داشت:

$$h_2(t) = f(c_t), \quad c_t = h_2(t-1) + x_1(t) + x_2(t) + b$$

12 / May 2024

۲۳

۳ ذی القعدة ۱۴۴۵

یکشنبه | اردیبهشت

t	c_t	$h_2(t)$
1	$2+b$	1
2	$3+b$	1
3	$2+b$	1
4	$1+b$	1
5	$2+b$	1
6	$3+b$	1
7	$1+b$	1

$$3+b > 0$$

$$2+b > 0$$

$$1+b > 0 \Rightarrow b > -1 \Rightarrow b = -0.8$$

برای هر t داریم:

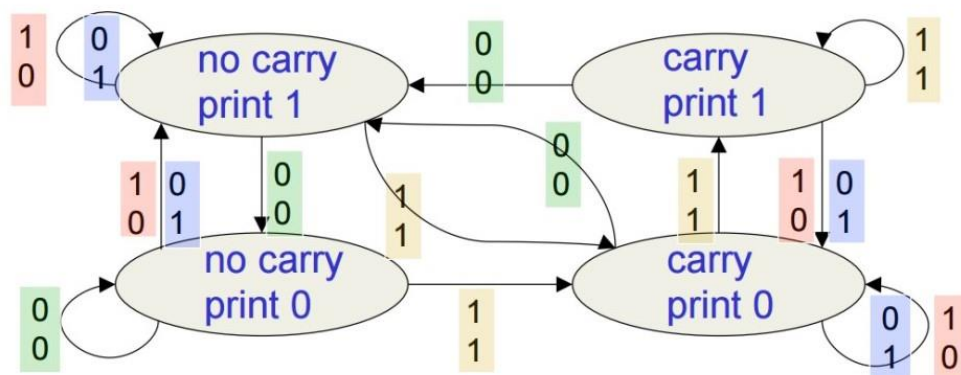
$$\Rightarrow c_t = h_2(t-1) + u_1(t) + u_2(t) - 0.5$$

$$\Rightarrow h_2(t) = f\left(\begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} h_1(t-1) \\ h_2(t-1) \\ h_3(t-1) \end{bmatrix} + \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} u_1(t) \\ u_2(t) \end{bmatrix} - 0.5\right)$$

در نهایت داریم:

$$h(t) = f\left(\begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} h_1(t-1) \\ h_2(t-1) \\ h_3(t-1) \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} u_1(t) \\ u_2(t) \end{bmatrix} + \begin{bmatrix} -1.8 \\ -0.8 \\ -1.8 \end{bmatrix}\right)$$

$$\Rightarrow \begin{cases} W = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ U = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \\ b_h = \begin{bmatrix} -1.8 \\ -0.8 \\ -1.8 \end{bmatrix} \end{cases}$$



سوال ۴:

(الف)

بگذارید با یک مثال از شبکه $GPT3$ جلو برویم. ورودی‌ها یا همان x ‌های این مدل بردارهایی با اندازه‌ای تقریباً برابر ۵۰ هزار هستند که از یک $word\ embedding$ عبور کرده و به فضای ۱۲ هزار بعدی می‌روند. D در این مدل همان ۱۰ هزار است. یعنی d اندازه هر x ورودی است. از طرفی w بیان می‌کند که ما چندمین ورودی هستیم. یعنی اگر برای مثال مدل GPT ما جمله در ورودی بگیرد، به طوری که هر x نشانگر یک کلمه باشد، w بیان می‌کند که x چندمین کلمه در جمله است. k هم می‌گوید که المان چندم P هستیم. برای هر x یک P داریم.

(ب)

سوال 4 :

(ب)

$$M \cdot \begin{bmatrix} \sin \omega_k t \\ \cos \omega_k t \end{bmatrix} = \begin{bmatrix} \sin(\omega_k(t+\phi)) \\ \cos(\omega_k(t+\phi)) \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} \sin \omega_k t \\ \cos \omega_k t \end{bmatrix} = \begin{bmatrix} \sin(\omega_k(t+\phi)) \\ \cos(\omega_k(t+\phi)) \end{bmatrix}$$

$$\Rightarrow \begin{cases} a \sin(\omega_k t) + b \cos(\omega_k t) = \sin(\omega_k(t+\phi)) & \times \cos(\omega_k t) \\ c \sin(\omega_k t) + d \cos(\omega_k t) = \cos(\omega_k(t+\phi)) & \times \sin(\omega_k t) \end{cases}$$

$$\Rightarrow \begin{cases} (a \sin \omega_k t)(\cos \omega_k t) + b \cos^2(\omega_k t) = \sin(\omega_k(t+\phi)) \cos(\omega_k t) \\ c \sin^2(\omega_k t) + d \cos \omega_k t (\sin \omega_k t) = \cos(\omega_k(t+\phi)) \sin(\omega_k t) \end{cases}$$

$$\Rightarrow \begin{cases} \frac{1}{2} a \sin(2\omega_k t) + \frac{1}{2} b (1 + \cos(2\omega_k t)) = \frac{1}{2} [\sin(2\omega_k t + \phi) + \sin(\omega_k \phi)] \\ \frac{1}{2} c (1 - \cos(2\omega_k t)) + \frac{1}{2} d \sin(2\omega_k t) = \frac{1}{2} [\sin(2\omega_k t + \phi) - \sin(\omega_k \phi)] \end{cases}$$

$$\begin{cases} a \sin(2\omega_k t) + b \cos(2\omega_k t) + b = \sin(2\omega_k t + \phi) + \sin(\omega_k \phi) \\ -c \cos(2\omega_k t) + d \sin(2\omega_k t) + c = \sin(2\omega_k t + \phi) - \sin(\omega_k \phi) \end{cases}$$

$$\Rightarrow \begin{cases} b = \sin \omega_k \phi \\ -c = -\sin \omega_k \phi \end{cases} \Rightarrow b = c$$

$$\Rightarrow \begin{cases} a \sin(2\omega_k t) + \sin(\omega_k \phi) \cos(2\omega_k t) = \sin(2\omega_k t + \phi) \\ \sin(\omega_k \phi) \cos(2\omega_k t) + d \sin(2\omega_k t) = \sin(2\omega_k t + \phi) \end{cases} \Rightarrow a = d = \cos(\omega_k \phi)$$

$$\Rightarrow M = \begin{bmatrix} \cos(\omega_k \phi) & \sin(\omega_k \phi) \\ -\sin(\omega_k \phi) & \cos(\omega_k \phi) \end{bmatrix}$$

این ماتریس متعام است :

$$|M| = \cos^2 + \sin^2 = 1$$

$$MM^T = \begin{bmatrix} \cos & \sin \\ -\sin & \cos \end{bmatrix} \begin{bmatrix} \cos & -\sin \\ \sin & \cos \end{bmatrix} = I$$

$$\begin{bmatrix} \cos & \sin \end{bmatrix} \cdot \begin{bmatrix} -\sin & \cos \end{bmatrix} = -\sin \cos + \sin \cos = 0$$

↓
dot

$$P_{t+k} = T^{(k)} P_t$$

این ماتریس برهمی خواهد بود که $d \times d$ است :

$$\begin{bmatrix} \cos(\omega_1 \phi) & \sin(\omega_1 \phi) & 0 & 0 & \dots & 0 & 0 \\ -\sin(\omega_1 \phi) & \cos(\omega_1 \phi) & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos(\omega_2 \phi) & \sin(\omega_2 \phi) & \dots & 0 & 0 \\ 0 & 0 & -\sin(\omega_2 \phi) & \cos(\omega_2 \phi) & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & \cos(\omega_d \phi) & \sin(\omega_d \phi) \\ 0 & 0 & 0 & 0 & \dots & -\sin(\omega_d \phi) & \cos(\omega_d \phi) \end{bmatrix}$$