

(7)

روش های Normalizing Flow یک نوعی از مدل های generative هستند

که بدلی حل مشکلات مدل های Autoregressive و VAE بشمارد² است:

در اینجا بدلی بدست آشنی نوعی حاشیه ای (مجموعه دایم):

$$P_z(x; \theta) = P_z \left(f_{\theta}^1(x) \right) \left| \det \left[\frac{\partial f_{\theta}^1(x)}{\partial x} \right] \right|$$

فعلی که در اینجا به کار رفته است، تبدیل دو تابع از این سی آر که تغییر معین یک یالتی مثال به از تبدیل invertible
ی دصو Flow بدلی مفی است که تبدیلات invertible می تواند با یک دیگر ترکیب شوند
که تبدیلات invertible (معکوس پذیر) به شدی شکل خود حال این روش ها (برسای

transformations های یایی روی یک تابع توزیع کاری کند) بر اساس قاعده $\frac{d}{dx}$ و الگوریتم توزیع را در لایه

یانی حسابی کند و بدلی ترتیب با استفاده از یک تابع ساده و تابع تغییر بدلیت معینده ای را
می سازند تبدلی توزیع داده این شبکه به روش بالا حسابی شود

$$\log P_i(z_i) = \log P_{i-1}(z_{i-1}) - \log \left(\det \left[\frac{\partial f_i}{\partial z_{i-1}} \right] \right)$$

حال اگر θ تابع اولیه خوب انتخاب شود، اصطلاحاً خوش معین باشد می توانیم حسابت
را به سادگی انجام دهیم. تابع دد ها نیزه اینجا به قدرت زیر تغییر می شود

$$Loss = -\frac{1}{|S|} \sum_{x \in S} \log P(x)$$

بنابراین با توجه به این روابط و تابع دد های توانیم شبکه را آموزش دهیم، backprop انجام
دهیم و پارامترهای مدل را بدست آوریم. همچنین بدلی می توانیم از این روش استفاده کنیم

همچنین برای تولید نمودار توزیع اولیه که معمولاً بدست می آید نموداری انجام می دهیم

پس با جدولی که تبدیل های P_1, P_2, \dots, P_k یک نمودار توزیع مقعده بدست می آید.

۸

~~این روش برای این کار مناسب نیست~~

۸

این روش این کار را انجام می دهد که به جای این که Q -learning را روی حالت $state$ / $action$ اجرا کنیم

انجام دهیم سیستم کل وضعیت $state$ ها را $[state, action, reward, next_state]$ در یک جدول جدید نگه می دارد. در واقع فاز یادگیری از یک تجربه جدا است و بدینای این است که از جدول نمونه برداری انجام دهیم.

~~در واقع این روش در الگوریتم Deep Reinforcement Learning استفاده می شود که بتواند ارزش کل~~

بدست آورد در تمامی حالات استفاده کنیم و مدل را به خوبی آموزش دهیم. البته این روش

سعی می کند $trade-off$ بین تجربه کردن و استفاده از تجربه ($exploration$ and $exploitation$)

را با الگوریتم ϵ -greedy رعایت کند و نهایتاً به ضریبی مطلوب برسد نسبت به تجربه کردن

دست پیدا کند (این الگوریتم در DQN هم استفاده می شود)

فایده استفاده از این روش ۱) استفاده بیشتر از تجربیات قبلی یادگیری جدید را

۲) Converge کردن بهتر وقتی که نیاز به یک بخش یادگیری داریم

مغایب استفاده از روش های چند مرحله ای یادگیری مانند $Q(\lambda)$ سرعت یادگیری می شود.