

یکشنبه

Sunday

9 Apr 2023

۱۸ رمضان ۱۴۴۴

فروردین

$$h(x) = w^T x + b \quad \text{loss}(h(x), y) = \log(1 + e^{-y h(x)}) \quad (1) \text{ (الف)}$$

$$x' = x + \delta \quad h(x') - h(x) = w^T \delta$$

$$\text{argmax}_x w^T \delta \quad \text{جواب} \quad \delta^* = \epsilon \cdot \text{sign}(w)$$

$$s.b. \quad 11.5.11.5.5.5$$

y = -1



بنابراین میزان تغییر $h(x)$ حداکثر برای $h(x') - h(x) = w^T \delta^* = \epsilon \|w\|$ می باشد.

اگر $y = 1$ باشد آنگاه تغییر x به میزان $\epsilon \cdot \text{sign}(w)$ یعنی $x' = x + \epsilon \cdot \text{sign}(w)$ به میزان

$\epsilon \|w\|$ در فرجه $h(x)$ تغییر ایجاد می کند که اگر این میزان باعث شود مرز تقسیم گیری را به یک سمت

یک دنده تغییر حساب می شود.

اگر $y = -1$ باشد آنگاه تغییر x به میزان $\epsilon \cdot \text{sign}(w)$ یعنی $x' = x + \epsilon \cdot \text{sign}(w)$ به میزان

$\epsilon \|w\|$ در فرجه $h(x)$ تغییر ایجاد می کند که اگر این میزان باعث شود مرز تقسیم گیری را به یک سمت

یک دنده تغییر حساب می شود.

چون در مدل FGDMM مقدار نرزی بهینه برای $\eta^* = \epsilon \cdot \text{sign}(\nabla_{\eta} f(x))$ است از آنجایی که در

این مدل داریم $f(x) = w^T x + b$ بنابراین داریم $\eta = \epsilon \cdot \text{sign}(w)$ از آنجایی که اگر $y = 1$ باشد

با به درجهت $-\text{sign}(w)$ حرکت کنیم تا بهینه برعکس شدن تریک شود بنابراین تابع فرجه در حالت

آینه شدن تغییرات با به صورت نرمی در نظر می گیریم.

$$\min_{w, b} \frac{1}{2} E_{(x, y)} [\text{loss}(h(x), y)] + \frac{1}{2} E_{(x - \epsilon \cdot \text{sign}(w), y)} [\text{loss}(h(x), y)]$$

$$= \frac{1}{2} \log(1 + e^{-y(w^T x + b)}) + \frac{1}{2} \log(1 + e^{-y(w^T (x - \epsilon \cdot \text{sign}(w)) + b)})$$

$$= \frac{1}{2} \log(1 + e^{-y(w^T x + b)}) + \frac{1}{2} \log(1 + e^{-y(w^T x + b - \epsilon \|w\|)})$$

تفاوت
الگوریتم این فرایند را می توان به دو روش بیان کرد: در روش اول تغییرات را به صورت گام به گام و در روش دوم به صورت یکباره.

ج) اگر از حالت یادگیری استفاده با استفاده از منظم سازیم ۱ استفاده کنیم تابع

$$\min_{w, b} E_{(x, y)} [\log(h(x, y))] + \lambda \|w\|_1$$

$$= \log(1 + e^{-y(\omega^T x + b)}) + \lambda \|w\|_1$$

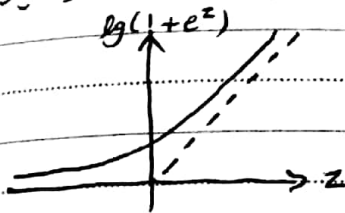
با شبیه تابع هزینه آلودگی خطا و آلودگی استفاده با منظم سازیم ۱ می بینیم مشاهده کردیم که در حالت
پسورد استفاده و هزینه افشان $\lambda \|w\|_1$ حدوده به تابع هزینه افشانی محدود و اکتفا مشاهده کردیم

که اکتفا مدل در بین بقیه دارد جفت است

اما در حالت آلودگی خطا و هزینه افشان $\lambda \|w\|_1$ و در بین تابع softmax قرار دارد و نسبت می کند

که اگر مدل با اکتفا خطی دارد و این را پیش بینی می کند، هزینه افشانی نامحدود به تابع هزینه افشانی محدود می شود

اگر مدل در بین بقیه کمتر و کمتر و به سمت پیش بینی خطای می رود، تأثیر این هزینه افشانی بیشتر و بیشتر



ی دارد این رفتار با مشاهده تابع softmax واضح است

مادامی که مدل استفاده از آلودگی استفاده به علت تحمل

یک هزینه مازاد به مدل، عملکرد و فضای زنی آن را محدود

می کند و باعث این عملکرد مدل می شود اما استفاده از آلودگی خطا به علت تحمل هزینه افشانی به جا

و مناسب به مدل می تواند عملکرد بهتری نسبت به حالت آلودگی استفاده با منظم سازیم ۱ داشته باشد

(۲) طبق بینار مقاله Towards evaluating the Robustness of NN

موردی که زیر را می توان برای انتخاب کلاس هدف در بین گرفتن:

۱. حالت میانگین (Average): کلاسی را از بین کلاسی های که کلاس مضیع دارد، مورد نظر

نقشه به صورت تعدادی با انتخاب مادی انتخاب کنیم

۲. حالت بهترین (Best): محله را برای تمام دسته های غیر مضیع انجام داده و یکی دسته ای را انتخاب

کنیم که سازه دین با کم ترین میزان perturbation را لازم داشت

۳. حالت بدترین (Worst): محله را برای تمام دسته های غیر مضیع انجام داده و یکی دسته ای را انتخاب

کنیم که سخت ترین با بیشترین میزان perturbation را لازم داشت

ب) برای اینکه یک محله هدفمند است یا نه باید سعی در کم کردن هزینه دانه بنوی داد با در نظر گرفتن

دسته درست و برای دسته هدف کنیم عبارت بهینه سازی به صورت زیر می شود:

$$\min_{x'} \text{loss}(h_{\theta}(x'), t)$$

↓
دسته هدف

(۳) رابطه مربوط به نحوه تغییر ساخته شده با محله PGD به صورت زیر است:

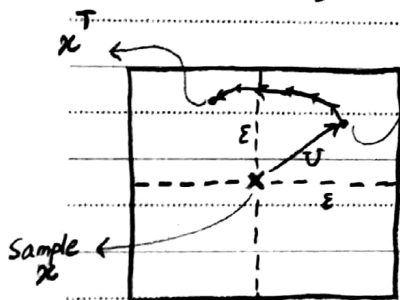
for $i=0$ to T :

$$\| \delta \|_{\infty} \leq \epsilon$$

$$x' = \text{clip}_{[-\epsilon, \epsilon]} \{ x' + \nabla(-\epsilon, \epsilon) \}$$

$$\delta' = \alpha \cdot \text{sign}(\nabla_x \text{loss}(\theta, x', y))$$

$$x^{i+1} = \text{clip}_{[\max(0, x - \epsilon), \min(1, x + \epsilon)]} \{ x' + \delta' \}$$



PGD و فلان FGM روشی iterative است. در روش FGM با یک تکرار فقط به تان به یکی از گوشه های مجاور $\| \delta \|_{\infty} = \epsilon$ می رود. در روش PGD ما در هر تکرار یک تکرار (با در نظر گرفتن ϵ به عنوان حد) اما در روش PGD ما در هر تکرار یک تکرار به دنبال یافتن نقطه ضعیف تر هستیم که تابع هزینه به ازای آن جاکسیم کند. همانطور که در شکل مشاهده می شود فضای جستجو محدود به $\| \delta \|_{\infty} \leq \epsilon$ است.

شب قدر: توجه! اگر منظور سوال چه PGD معرّفی شده در مقاله باشد، تفاوت آن با این روش گفته شده این است که شروع تعداد ندارد.

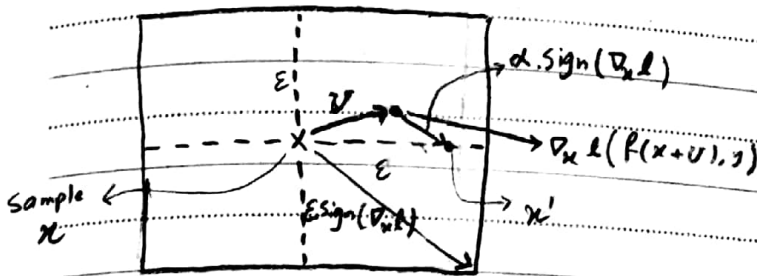
رابطه برضای مندرجۀ فضا نه ساخته شده باشد. f_{GSM-RS} به صورت نوید است:

$$\delta = U(-\epsilon, \epsilon)$$

$$\delta = \delta + \alpha \cdot \text{sign}(\nabla_x \text{loss}(f(x+\delta), y))$$

$$\delta = \max(\min(\delta, \epsilon), -\epsilon)$$

$$x' = x + \delta$$



تفاوت روش f_{GSM-RS} با روش f_{GSM} در
موضوع تغییرات آن است. این موضوع تغییرات
درین روش در جهت $\text{sign}(\Delta L)$ با ضرب α
عملیات تکمیل که نتایج به تمام جود
عکس ۱۱۶۱ دای ساخت نموده فضا نه بررسی
داشته باشیم. بنابرین عملکرد این روش
در حد عملکرد روش گزینی PhD می باشد.