

1.a

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \stackrel{(*)}{=} \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} = A$$

$$y_i = \sum_{k=1}^n a_{ik} x_k \rightarrow \frac{\partial y_i}{\partial x_j} = a_{ij} \quad (*)$$

1.b

مسئله را در حالت کلی می‌کنیم که z بردار باشد، تمام اعضای x و تمام اعضای z وابسته باشند.

$$y_i = \sum_{k=1}^n a_{ik} x_k \rightarrow \frac{\partial y_i}{\partial z_j} = \sum_{k=1}^n a_{ik} \frac{\partial x_k}{\partial z_j} \quad (*)$$

$$\begin{aligned} \frac{\partial y}{\partial z} &= \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \dots & \frac{\partial y_1}{\partial z_\ell} \\ \vdots & & \vdots \\ \frac{\partial y_m}{\partial z_1} & \dots & \frac{\partial y_m}{\partial z_\ell} \end{bmatrix}_{m \times \ell} \stackrel{(*)}{=} \begin{bmatrix} \sum_{k=1}^n a_{1k} \frac{\partial x_k}{\partial z_1} & \dots & \sum_{k=1}^n a_{1k} \frac{\partial x_k}{\partial z_\ell} \\ \vdots & & \vdots \\ \sum_{k=1}^n a_{mk} \frac{\partial x_k}{\partial z_1} & \dots & \sum_{k=1}^n a_{mk} \frac{\partial x_k}{\partial z_\ell} \end{bmatrix}_{m \times \ell} \\ &= \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}_{m \times n} \times \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \dots & \frac{\partial x_1}{\partial z_\ell} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial z_1} & \dots & \frac{\partial x_n}{\partial z_\ell} \end{bmatrix}_{n \times \ell} \\ &= A \times \frac{\partial x}{\partial z} \end{aligned}$$

1.c

اگر فرض کنیم $y = Ax$ ، داریم:

$$\begin{aligned} y^T A x &= \|y\|^2 = \left(\sum_{k=1}^n a_{1k} x_k \right)^2 + \dots + \left(\sum_{k=1}^n a_{mk} x_k \right)^2 \\ \rightarrow \frac{\partial y^T A x}{\partial x} &= \begin{bmatrix} \sum_{j=1}^m 2 a_{j1} \sum_{k=1}^n a_{jk} x_k \\ \vdots \\ \sum_{j=1}^m 2 a_{jn} \sum_{k=1}^n a_{jk} x_k \end{bmatrix} = 2 y^T A \end{aligned}$$

در کدورتا گفته شد که x, y را مستقل از هم در نظر بگیریم. در این صورت می‌توان $y^T A$ را ماتریسی مستقل از x دانست. بنابراین آنرا S می‌نامیم و طبق بخش 1.a می‌توان نوشت:

$$\frac{\partial \left[\overset{y^T A}{S} \right]_{1 \times n} \left[x \right]_{n \times 1}}{\partial \left[x \right]_{n \times 1}} \stackrel{a_{ij}}{=} S = y^T A$$

⊗ البته معمولاً متوجه اسکالر نسبت به بردار را به صورت column vector می‌نویسند اما در این سوال به صورت row vector در نظر گرفته شده.

1.c ^{ادامه} $\alpha = y^T A x$ را می توان به صورت dot product دوبعدار y و Ax در نظر گرفت و می دانیم که ترتیب در dot product بی اثر است

$$\alpha = y^T A x = (Ax)^T y = x^T A^T y$$

مستطاب قبل $x^T A^T$ را یک ماتریس مستقل از y در نظر می گیریم و طبق بخش a داریم:

$$\frac{\partial \alpha}{\partial y} = x^T A^T$$

البته می دانیم که در dot product دو بردار مستقل، مستقل نسبت به یکی برابر دیگری است اما فوایدی از آن استفاده نکنیم

1.d

$$\alpha = \sum_{i=1}^n x_i y_i \rightarrow \frac{\partial \alpha}{\partial z_j} = \sum_{i=1}^n \frac{\partial x_i}{\partial z_j} y_i + \sum_{i=1}^n \frac{\partial y_i}{\partial z_j} x_i \quad (*)$$

z را مستطاب بخش b فرض می کنیم و داریم:

$$\frac{\partial \alpha}{\partial z} = \begin{bmatrix} \frac{\partial \alpha}{\partial z_1} \\ \vdots \\ \frac{\partial \alpha}{\partial z_\ell} \end{bmatrix}^T = \begin{bmatrix} \sum_{i=1}^n \frac{\partial x_i}{\partial z_1} y_i \\ \vdots \\ \sum_{i=1}^n \frac{\partial x_i}{\partial z_\ell} y_i \end{bmatrix}^T + \begin{bmatrix} \sum_{i=1}^n \frac{\partial y_i}{\partial z_1} x_i \\ \vdots \\ \sum_{i=1}^n \frac{\partial y_i}{\partial z_\ell} x_i \end{bmatrix}^T = [y_1 \dots y_n] \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \dots & \frac{\partial x_1}{\partial z_\ell} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial z_1} & \dots & \frac{\partial x_n}{\partial z_\ell} \end{bmatrix} + [x_1 \dots x_n] \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \dots & \frac{\partial y_1}{\partial z_\ell} \\ \vdots & & \vdots \\ \frac{\partial y_n}{\partial z_1} & \dots & \frac{\partial y_n}{\partial z_\ell} \end{bmatrix}$$

$$= y^T \frac{\partial x}{\partial z} + x^T \frac{\partial y}{\partial z}$$

در این سوال مستقل اسکالر نسبت به بردار به صورت row vector فرض شده است

1.e

$$A A^{-1} = I_m \xrightarrow{\frac{\partial}{\partial \alpha}} \frac{\partial A}{\partial \alpha} A^{-1} + A \frac{\partial A^{-1}}{\partial \alpha} = 0 \rightarrow A \frac{\partial A^{-1}}{\partial \alpha} = -\frac{\partial A}{\partial \alpha} A^{-1}$$

$$x A^{-1} \xrightarrow{\frac{\partial}{\partial \alpha}} \frac{\partial A^{-1}}{\partial \alpha} = -A^{-1} \frac{\partial A}{\partial \alpha} A^{-1}$$

2

گرایان: بردار مشتق‌های جزئی تابع یا فرجهی اسکالر، نسبت به هر کدام از ورودی‌های سی‌دهد

ماتریس گرادیان: ماتریسی از مشتق‌های جزئی که الان از آن مشتق جزئی فرجهی نام نسبت به ورودی نام است (برای تابع با ورودی و فرجهی برداری)

$$\nabla y = \begin{bmatrix} \frac{\partial y}{\partial u} \\ \frac{\partial y}{\partial v} \\ \frac{\partial y}{\partial z} \end{bmatrix}$$

$$\text{Jacobian}(\nabla y) = \begin{bmatrix} \frac{\partial^2 y}{\partial u^2} & \frac{\partial^2 y}{\partial u \partial v} & \frac{\partial^2 y}{\partial u \partial z} \\ \frac{\partial^2 y}{\partial v \partial u} & \frac{\partial^2 y}{\partial v^2} & \frac{\partial^2 y}{\partial v \partial z} \\ \frac{\partial^2 y}{\partial z \partial u} & \frac{\partial^2 y}{\partial z \partial v} & \frac{\partial^2 y}{\partial z^2} \end{bmatrix}$$

ماتریس سی‌دهد که دقیقاً به تعریف

Hessian رسیدیم

3

در شبکه‌های متفاوت همواره می‌توان بر اساس domain آن شبکه فرض‌های انجام داد تا فضای حالت را محدود کرده و inductive bias مناسبی را بر اساس شرایط آن مسئله، در راه حل‌مان قوی‌تر کنیم. به عنوان مثال اگر بدانیم که در صورت مسئله ما همواره در داده‌ها نوعی تقارن وجود دارد، می‌توانیم این فرض را در راه حل‌مان استفاده کنیم. استفاده از تقارن در شبکه‌های مصنوعی مانند metric learning (یافتن representation مناسبی که similarity در آن معنا داشته باشد)، خوشه‌بندی، classification و ... می‌تواند کاربرد داشته باشد.

این تقارن در حالت کلی می‌تواند یک group structure را تعیین کند و یا در حالت صورت سوال ما، تقارن در یک تصویر باشد.

هدف آن است که تأثیر فرض تقارن را در loss شبکه ایجاد کنیم. برای این کار می‌توان regularization term را به نوری تغییر داد که در نظر نگرفتن تقارن ورودی را penalize کند.

طبق شرایط صورت سوال فرض می‌کنیم بخواهیم یک شبکه binary classification انجام دهیم و ورودی‌مان یک تصویر 1x2 باشد. در این حالت باید بردار وزن‌ها نیز طول 2 داشته باشد.

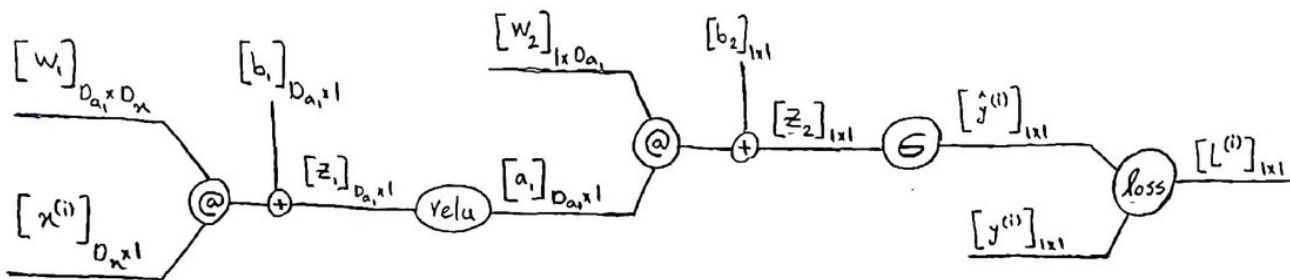
مطلوب ما آن است که در بردار وزن $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ مقادیر w_1 و w_2 به هم نزدیک باشند زیرا ورودی دارای تقارن است و بنابراین تأثیر 2 الان آن در فرجهی شبکه نیز سبب به هم باشد.

اگر می‌کنیم اگر S را $\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ در نظر بگیریم، w_1 و w_2 نزدیک به هم می‌مانند:

$$R(w) = w^T S w = \underbrace{\begin{bmatrix} w_1 & w_2 \end{bmatrix}}_{\begin{bmatrix} w_1 - w_2 & w_2 - w_1 \end{bmatrix}} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = w_1(w_1 - w_2) + w_2(w_2 - w_1) = (w_1 - w_2)(w_1 - w_2) = (w_1 - w_2)^2$$

بنابراین هرچه فاصله w_1 و w_2 بیشتر شود لاس بیشتر خواهیم داشت ← مدل در جهتی وزن‌ها را اهمیت می‌دهد که تقارن ایجاد شود

4.a



$$J = \frac{1}{m} \sum_{i=1}^m L^{(i)}$$

دری تصویر تمام اجزای مشخص شده است

4.b

$$\frac{\partial J}{\partial \hat{y}^{(i)}} = \delta_1 = \frac{1}{m} \left(\frac{y^{(i)}}{\hat{y}^{(i)}} + \frac{(y^{(i)} - 1)}{1 - \hat{y}^{(i)}} \right) = [\delta_1]_{1 \times 1}$$

4.c

$$\frac{\partial \hat{y}^{(i)}}{\partial z_2} = \delta_2 = \underbrace{\sigma(z_2)}_{\text{sigmoid}} \times (1 - \sigma(z_2)) = [\delta_2]_{1 \times 1}$$

4.d

$$\frac{\partial z_2}{\partial a_1} = \delta_3 = w_2 = [\delta_3]_{1 \times D_{a_1}}$$

4.e

$$\frac{\partial a_1}{\partial z_1} = \delta_4 = \begin{cases} 0 & z_1 < 0 \\ 1 & z_1 \geq 0 \end{cases} \quad \begin{array}{l} \text{هر امان } z_1 \text{ که منفی یا صفر باشد} \leftarrow 0 \\ \text{هر امان } z_1 \text{ که مثبت باشد} \leftarrow 1 \end{array} = [\delta_4]_{D_{a_1} \times 1}$$

4.f

$$\frac{\partial z_1}{\partial w_1} = \delta_5 = x^{(i)T} = [\delta_5]_{1 \times D_x}$$

4.g

$$\frac{\partial J}{\partial w_1} = \delta_1 * \delta_2 * [\delta_3^T]_{D_{a_1} \times 1} \overset{\text{element wise}}{*} [\delta_4]_{D_{a_1} \times 1} \overset{\text{matmul}}{@} [\delta_5]_{1 \times D_x}$$

5.1

Beale

$$F(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2$$

$$\begin{aligned} \frac{\partial F}{\partial x} &= 2(1.5 - x + xy)(-1 + y) \\ &\quad + 2(2.25 - x + xy^2)(-1 + y^2) \\ &\quad + 2(2.625 - x + xy^3)(-1 + y^3) \end{aligned}$$

$$\begin{aligned} \frac{\partial F}{\partial y} &= 2(1.5 - x + xy)(x) \\ &\quad + 2(2.25 - x + xy^2)(2xy) \\ &\quad + 2(2.625 - x + xy^3)(3xy^2) \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 F}{\partial x^2} &= 2(-1 + y)(-1 + y) \\ &\quad + 2(-1 + y^2)(-1 + y^2) \\ &\quad + 2(-1 + y^3)(-1 + y^3) \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 F}{\partial y^2} &= 2x(x) \\ &\quad + 2[(2xy)(2xy) + (2x)(2.25 - x + xy^2)] \\ &\quad + 2[(3xy^2)(3xy^2) + (6xy)(2.625 - x + xy^3)] \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 F}{\partial x \partial y} &= 2 \left[(x)(-1 + y) + (1.5 - x + xy) \right] \\ &\quad + 2 \left[(2xy)(-1 + y^2) + (2y)(2.25 - x + xy^2) \right] \\ &\quad + 2 \left[(3xy^2)(-1 + y^3) + (3y^2)(2.625 - x + xy^3) \right] \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 F}{\partial y \partial x} &= 2 \left[(-1 + y)(x) + (1.5 - x + xy) \right] \\ &\quad + 2 \left[(-1 + y^2)(2xy) + (2y)(2.25 - x + xy^2) \right] \\ &\quad + 2 \left[(-1 + y^3)(3xy^2) + (3y^2)(2.625 - x + xy^3) \right] \end{aligned}$$

$$\text{gradient} = \begin{bmatrix} \frac{\partial F}{\partial x} \\ \frac{\partial F}{\partial y} \end{bmatrix}$$

$$\text{hessian} = \begin{bmatrix} \frac{\partial^2 F}{\partial x^2} & \frac{\partial^2 F}{\partial x \partial y} \\ \frac{\partial^2 F}{\partial y \partial x} & \frac{\partial^2 F}{\partial y^2} \end{bmatrix}$$

می دانیم که تابع $f(x,y)$ که دو بار مشتق پذیر است، محدب است اگر و تنها اگر ماتریس Hessian آن psd باشد (در تعریف 5.1)

نشان می دهیم که به ازای نقطه $(0,0)$ ماتریس Hessian برابر است با

$$\text{Hessian} = \begin{bmatrix} 2+2+2 & 3+0+0 \\ 3+0+0 & 0+0+0 \end{bmatrix} = \begin{bmatrix} 6 & 3 \\ 3 & 0 \end{bmatrix} \xrightarrow{\text{eigen values}} \begin{cases} \lambda_1 = 3+3\sqrt{2} & \text{مثبت} \\ \lambda_2 = 3-3\sqrt{2} & \text{منفی} \end{cases}$$

تمام eigen val ها بزرگتر مساوی صفر نیستند
پس ماتریس Hessian به ازای نقطه $(0,0)$ psd نیست
پس تابع Beale تابع محدب نیست

۴) به طور کلی توابع غیر محدب برای optimization سخت تر هستند زیرا می توانند دارای local minima های مختلف، saddle point، نواحی Flat بزرگ و curvature های متفاوت باشند. این موارد برای الگوریتم های بهینه سازی چالش را هستند و باعث می شوند که الگوریتم نتواند global minimum را پیدا کند. اکثر مسائلی که شبکه های عصبی با آن رو به رو هستند از این جنس است و بنابراین optimizer های متوهمی در درس معرفی شد که سعی دارند این چالش ها را رفع کنند

$$\text{gradient}(0,1) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(0,1) \\ \frac{\partial f}{\partial x_2}(0,1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$x^1 = x^0 - \eta \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow \text{نقطه شش پس از یک بار به روز رسانی} = x^1 = x^0 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

5.2

الگوریتم gradient descent در هر قدم در جهت منفی گرادیان $h(n)$ حرکت می‌کند. می‌دانیم Q ماتریس pd و بنابراین متقارن است. پس:

$$\nabla_n h(n) = \frac{1}{2} (Q + Q^T) n + c = Qn + c$$

$$\text{میانگین } x^* \text{ را به جواب می‌دهیم} \rightarrow Qx^* + c = 0 \rightarrow c = -Qx^*$$

بردار r^t را به صورت $r^t = x^t - x^*$ تعریف می‌کنیم و داریم:

$$\left. \begin{aligned} x^{t+1} &= x^t - \eta (Qx^t + c) \\ r^t &= x^t - x^* \rightarrow x^t = r^t + x^* \end{aligned} \right\} \rightarrow r^{t+1} + x^* = r^t + x^* - \eta (Qx^t - Qx^*)$$

$$\rightarrow r^{t+1} = (I - \eta Q) r^t$$

$$r^{t+1} = U \underbrace{(I - \eta U^T Q U)}_{\Sigma} U^T r^t \quad \textcircled{I}$$

$$\Sigma = \begin{bmatrix} 1 - \eta \lambda_1 & & 0 \\ & \ddots & \\ 0 & & 1 - \eta \lambda_n \end{bmatrix}$$

از جبرخطی می‌دانیم که چون Q متقارن است پس diagonalizable است. پس می‌توان Q را به صورت مقابل decompose کرد:

$$Q = U \Sigma U^T \rightarrow \Sigma = U^T Q U$$

که Σ ماتریس قطری است که eigen val های Q اعضای قطر آن هستند و U ماتریس unitary است که ستون های آن eigen vector های Q هستند.

$$\rightarrow r^{t+1} = U \begin{bmatrix} (1 - \eta \lambda_1)^{t+1} & & 0 \\ & \ddots & \\ 0 & & (1 - \eta \lambda_n)^{t+1} \end{bmatrix} U^T r^0 \quad \textcircled{II}$$

طبق تعریف بردار r^t جهت حرکت به سمت نقطه بهینه را نشان می‌دهد.

می‌دانیم که ماتریس های unitary یا reflection اند یا rotation و اندازه ها

را تغییر نمی‌دهند. همچنین می‌دانیم که ماتریس قطری تنها در راستای محور ها

scale می‌کند و جهت را تغییر نمی‌دهد.

در حقیقت فرمول \textcircled{I} : بیانگر آن است که r^t در هر مرحله change of basis

روی آن اعمال شده و توسط ماتریس U^T به فضای می‌رود که eigen vector های Q

محور های مختصات اند، سپس در آن فضا توسط ماتریس قطری در راستای این

eigen vector scale می‌شود و سپس توسط U به فضای قبلی برمی‌گردد. بنابراین می‌توان گفت که جهت حرکت به سمت x^* در هر مرحله، در راستای eigen vector

های Q است. فرمول \textcircled{II} نشان می‌دهد که اگر تکرار باشد converge کنیم باید ماتریس قطری به مرور به ماتریس 0 میل کند پس $0 < \eta \lambda_i < 1 \rightarrow 0 < 1 - \eta \lambda_i < 1$

بدون از دست رفتن کلیت می‌توان فرض کرد که λ_n بزرگترین eigen val باشد و بنابراین $(1 - \eta \lambda_n)^{t+1}$ دیرتر از بقیه λ_n های ماتریس به صفر میل می‌کند و بنابراین در

کلام های نهایی تنها در جهت eigen vector ای حرکت می‌کنیم که کمترین eigen val را دارد

⊗ البته حرکت کردن در راستای یک eigen vector در کلام های نهایی معنای بهینه‌سازی دارد تا حرکت کردن در راستای تمام eigen vector ها زیرا n eigen vector داریم و

eigen space بعضی به اندازه بُعد فضای اولیه دارد

5.3

a Adam

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2 \rightarrow \text{element wise}$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

در الگوریتم Adam برای به دست آوردن v_t به جای دوم 2 می توان از حالت کلی تر دوم p استفاده کرد که به صورت زیر در می آید:

$$v_t = \beta_2^p v_{t-1} + (1-\beta_2^p) |g_t|^p \rightarrow \text{element wise}$$

حال p را به ∞ میل می دهیم و u_t را به صورت زیر تعریف می کنیم:

$$\begin{aligned} u_t &= \lim_{p \rightarrow \infty} (v_t)^{\frac{1}{p}} \\ &= \lim_{p \rightarrow \infty} \left((1-\beta_2^p) \sum_{i=1}^t \beta_2^{p(t-i)} \cdot |g_i|^p \right)^{\frac{1}{p}} \\ &= \lim_{p \rightarrow \infty} \left(\sum_{i=1}^t (\beta_2^{(t-i)}) \cdot |g_i|^p \right)^{\frac{1}{p}} \\ &= \max \left(\beta_2^{t-1} |g_1|, \dots, \beta_2 |g_{t-1}|, |g_t| \right) \\ &= \max \left(\beta_2 \cdot u_{t-1}, |g_t| \right) \end{aligned}$$

Adamax

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

$$u_t = \max(\beta_2 \cdot u_{t-1}, |g_t|)$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{u_t} \hat{m}_t$$

بنابراین به Varient ای از Adam به نام Adamax رسیدیم که به صورت قابل است و زیر اساس infinity norm عمل می کند:

b با توجه به اینکه در Adamax بین مقدار قبلی و گرادیان ماکسیمم گرفته می شود، این روش در مواردی که گرادیان noisy باشد robust تر عمل می کند و عملاً مقادیر کوچک گرادیان را ignore می کند. یک نمونه از این حالت زمانی است که آپدیت وزن ها به صورت sparse باشد. در این حالت ابعادی که در آن ها نباید آپدیتی انجام شود تحت تأثیر نویز گرادیان قرار نمی گیرند. به طور کلی برای برآورد کردن با max نسبت به L2 norm در برابر نویز robust تر است. همچنین در مواقعی که پارامترها با درگانی کمی آپدیت می شوند نیز adamax می تواند بهتر باشد.

البته در کل Adam معروف تر و مورد استفاده تر است زیرا به رنج مسائل بیشتری apply می شود و محاسبات efficient تری دارد.

در روش های momentum و adaptive بودن را استفاده می کنند اما Adamax برای adaptive کردن از max استفاده می کند و Adam از L2 norm