

$$\begin{cases} x: n \times 1 \\ y: m \times 1 \end{cases} \quad \begin{cases} y = Ax \\ A: m \times n \end{cases}$$

سوال ①

$$1) \quad y_i = \sum_{k=1}^n a_{ik} x_k \quad j \in [1, n] \Rightarrow \frac{\partial y_i}{\partial x_j} = a_{ij} \Rightarrow \frac{\partial y}{\partial x} = A_{m \times n}$$

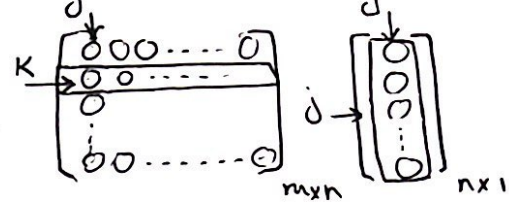
$$\frac{\partial y}{\partial x_j} = \begin{bmatrix} \frac{\partial y_1}{\partial x_j} \\ \frac{\partial y_2}{\partial x_j} \\ \vdots \\ \frac{\partial y_m}{\partial x_j} \end{bmatrix} \Rightarrow \frac{\partial y}{\partial x} = A = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

$$2) \quad \frac{\partial y}{\partial z} = \frac{\partial (Ax)}{\partial z} = \frac{\partial A}{\partial z} x + A \frac{\partial x}{\partial z} = A \frac{\partial x}{\partial z}$$

$$3) \quad \alpha = y^T A x \Rightarrow \alpha = \sum_{i=1}^m y_i \sum_{j=1}^n a_{ij} x_j = \sum_{i=1}^m \sum_{j=1}^n y_i a_{ij} x_j$$

$$\frac{\partial \alpha}{\partial x_k} = \sum_{i=1}^m y_i a_{ik} = \sum_{i=1}^m \begin{bmatrix} y_i \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{m \times 1} \begin{bmatrix} a_{i1} & a_{i2} & \dots & a_{in} \end{bmatrix}_{1 \times n} \Rightarrow$$

$$\frac{\partial \alpha}{\partial x} = y^T A \quad [1 \times n]$$

$$3) \quad \frac{\partial \alpha}{\partial y_k} = \sum_{j=1}^n a_{kj} x_j = \sum_{j=1}^n$$


$$\frac{\partial \alpha}{\partial y} = x^T A^T \quad [1 \times m]$$

$1 \times n \quad n \times m$

$$4) \quad \alpha = y^T x \quad [1 \times 1] \text{ scalar}$$

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

$$\otimes \quad \frac{\partial \alpha}{\partial x} = y^T, \quad \frac{\partial \alpha}{\partial y} = x^T \quad \text{"taken from previous part"}$$

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

5) if A is non-singular then $\det(A) \neq 0$ and A^{-1} exists.

$$\text{So: } AA^{-1} = I \quad \Rightarrow \quad \frac{\partial(AA^{-1})}{\partial \alpha} = \frac{\partial I}{\partial \alpha} = 0$$

$$\frac{\partial(AA^{-1})}{\partial \alpha} = \frac{\partial A}{\partial \alpha} A^{-1} + A \frac{\partial A^{-1}}{\partial \alpha} = 0 \quad \Rightarrow$$

$$A \frac{\partial A^{-1}}{\partial \alpha} = - \frac{\partial A}{\partial \alpha} A^{-1} \quad \xrightarrow{A^{-1} \times} \quad I \frac{\partial A^{-1}}{\partial \alpha} = - A^{-1} \frac{\partial A}{\partial \alpha} A^{-1}$$

$$\Rightarrow \quad \frac{\partial A^{-1}}{\partial \alpha} = - A^{-1} \frac{\partial A}{\partial \alpha} A^{-1}$$

$$y = \psi(u, v, z)$$

$$\nabla \psi = \left(\frac{\partial \psi}{\partial u}, \frac{\partial \psi}{\partial v}, \frac{\partial \psi}{\partial z} \right)_{3 \times 1}$$

$$J(\nabla \psi) = \begin{bmatrix} \frac{\partial}{\partial u} \left(\frac{\partial \psi}{\partial u} \right) & \frac{\partial}{\partial u} \left(\frac{\partial \psi}{\partial v} \right) & \frac{\partial}{\partial u} \left(\frac{\partial \psi}{\partial z} \right) \\ \frac{\partial}{\partial v} \left(\frac{\partial \psi}{\partial u} \right) & \frac{\partial}{\partial v} \left(\frac{\partial \psi}{\partial v} \right) & \frac{\partial}{\partial v} \left(\frac{\partial \psi}{\partial z} \right) \\ \frac{\partial}{\partial z} \left(\frac{\partial \psi}{\partial u} \right) & \frac{\partial}{\partial z} \left(\frac{\partial \psi}{\partial v} \right) & \frac{\partial}{\partial z} \left(\frac{\partial \psi}{\partial z} \right) \end{bmatrix}$$

$$\Rightarrow J(\nabla \psi) = \begin{bmatrix} \frac{\partial^2 \psi}{\partial u^2} & \frac{\partial^2 \psi}{\partial u \partial v} & \frac{\partial^2 \psi}{\partial u \partial z} \\ \frac{\partial^2 \psi}{\partial v \partial u} & \frac{\partial^2 \psi}{\partial v^2} & \frac{\partial^2 \psi}{\partial v \partial z} \\ \frac{\partial^2 \psi}{\partial z \partial u} & \frac{\partial^2 \psi}{\partial z \partial v} & \frac{\partial^2 \psi}{\partial z^2} \end{bmatrix}$$

$$H = \nabla^2 \psi = \begin{bmatrix} \frac{\partial^2 \psi}{\partial u^2} & \frac{\partial^2 \psi}{\partial u \partial v} & \frac{\partial^2 \psi}{\partial u \partial z} \\ \frac{\partial^2 \psi}{\partial v \partial u} & \frac{\partial^2 \psi}{\partial v^2} & \frac{\partial^2 \psi}{\partial v \partial z} \\ \frac{\partial^2 \psi}{\partial z \partial u} & \frac{\partial^2 \psi}{\partial z \partial v} & \frac{\partial^2 \psi}{\partial z^2} \end{bmatrix}$$

So Hessian of ψ function equals to
the Jacobian of its gradient.

برخی از انواع داده‌ها دارای تقارن ذاتی هستند و این تقارن به اشکال مختلفی می‌تواند وجود داشته باشد. برای مثال تقارن فضایی در تصاویر صورت وجود دارد و یک چهره ممکن است در امتداد محور عمودی تقارن داشته باشد. یعنی یک طرف صورت، طرف مقابل را می‌تواند منعکس کند. همچنین داده‌های سری زمانی می‌توانند دارای تقارن زمانی باشند؛ یعنی الگوهای (اُس) باشند که در طول زمان تکرار شوند. (به صورت ماهیانه، فصلی و...)

این تقارن به صورت‌های مختلفی مثل افزایش داده‌ها، مهندسی Feature، انتخاب بهتر معماری شبکه و انتخاب regularization بهتر باعث بهبود عملکرد شبکه Classification شود. افزایش داده: وجود تقارن به ما امکان تولید داده‌های جدید با اعمال تبدیلات بازتاب و چرخش می‌دهد.

مهندسی Feature: وجود تقارن کمک می‌کند که در مهندسی ویژگی، ویژگی‌های مبتنی بر تقارن را برای ثبت ویژگی‌های داده در نظر گرفت.

معنی مدل: می‌توان از شبکه‌های عصبی ای که دانش مربوط به تقارن را می‌تواند بدست بیاورد مثل شبکه‌های CNN استفاده کرد. چون این شبکه‌ها امکان بهره‌گیری از تقارن‌های translational را دارند.

regularization: از تکنیک‌های منظم‌سازی که وزن‌های شبکه را تسویه یا یکدلی فینچ‌های متقارن در داده می‌کنند، می‌توان استفاده کرد. ویژگی weight-sharing در CNN باعث یکدلی تقارن‌هایی در پارامترها می‌شود.

چون ترم regularization به تابع هزینه اضافی شود و شبکه به دنبال کاهش این ترم خواهد بود. پس باید این ترم regularization را طوری انتخاب کنیم تا وزن هایمان اختلاف زیادی نداشته باشند.

$$x = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \quad (1,2)$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad (2,1)$$

$$R(w) = w^T S w$$

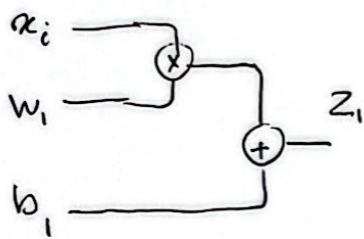
اگر بخواهیم w_1 و w_2 نزدیک هم باشند و نامتقارن نشوند، باید S را طوری انتخاب کنیم تا $(w_1 - w_2)^2$ را در ترم رگولاریزیشن ایجاد کنیم.

$$R(w) = \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = a w_1^2 + b w_1 w_2 + c w_1 w_2 + d w_2^2$$

$$\text{if } R(w) = (w_1 - w_2)^2 = w_1^2 - 2 w_1 w_2 + w_2^2$$

$$\text{then } a=1, b=-1, c=-1, d=1$$

$$S = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$



$$z_1 = w_1 x^{(i)} + b_1$$

$$x_i : [D_x, 1]$$

$$z_1 : [D_{a_1}, 1]$$

$$[D_{a_1}, 1] = [D_{a_1}, D_x] \cdot [D_x, 1] + [D_{a_1}, 1]$$

$$1) \quad w_1 : [D_{a_1}, D_x], \quad b_1 : [D_{a_1}, 1]$$

$$a_1, z_1 : [D_{a_1}, 1] \quad z_2 = w_2 a_1 + b_2$$

$$[D_{a_2}, 1] = [D_{a_2}, D_{a_1}] \cdot [D_{a_1}, 1] + [D_{a_2}, 1]$$

$$w_2 : [D_{a_2}, D_{a_1}], \quad b_2 : [D_{a_2}, 1]$$

$$2) \quad \frac{\partial J}{\partial \hat{y}^{(i)}} = \frac{-1}{m} \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} = \frac{-1}{m} \left(\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{1 - y^{(i)}}{1 - \hat{y}^{(i)}} \right) = \delta_1$$

$$3) \quad \frac{\partial \hat{y}^{(i)}}{\partial z_2} = \sigma(z_2) (1 - \sigma(z_2)) = \delta_2$$

$$4) \quad \frac{\partial z_2}{\partial a_1} = w_2 = \delta_3$$

$$5) \quad \frac{\partial a_1}{\partial z_1} = \begin{cases} 1 & z_1 \geq 0 \\ 0 & z_1 < 0 \end{cases} = \delta_4$$

$$6) \quad \frac{\partial z_1}{\partial w_1} = x^T = \delta_5$$

$$7) \frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial y^{(i)}} \times \frac{\partial y^{(i)}}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\Rightarrow \frac{\partial J}{\partial w_1} = \delta_1 \times \delta_2 \times \delta_3 \times \delta_4 \times \delta_5 \Rightarrow$$

$$\frac{\partial J}{\partial w_1} = \frac{-1}{n} \left(\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{1-y^{(i)}}{1-\hat{y}^{(i)}} \right) \sigma(z_2) (1-\sigma(z_2)) w_2 x^T$$

$$\text{for } z_1 \geq 0$$

$$\frac{\partial J}{\partial w_1} = 0 \quad \text{for } z_1 < 0$$

(1) تابع beale

$$f(x) = (1.5 - x_1 + x_1 x_2^2)^2 + (2.25 - x_1 + x_1 x_2^2)^2 + (2.625 - x_1 + x_1 x_2^3)^2$$

تابع beale غیر محدب است زیرا ماتریس Hessian آن در همه جا positive semi-definite نیست.

همچنین بطور خاص تر در نقطه $(0,1)$ دارای saddle point است

و در این نقطه دارای مقادیر مثبت و منفی در ماتریس Hessian است.

این موقع نشان می دهد تابع در برخی جهات به سمت بالا منحنی می شود و در برخی

جهات به سمت پایین. پس بصورت کلی Convex نمی باشد.

به دلایل مختلفی Convexity در optimization اهمیت دارد:

- توابع محدب مشکل local minima را ندارند و هر local minima در توابع محدب

یک global minima نیز هست.

- بهینه سازی توابع محدب ساده تر از توابع غیر محدب است و می توان راحت تر global minima را یافت.

- پایداری و robustness توابع محدب نسبت به انحراف کوچک بسطی از توابع

غیر محدب است.

$$\nabla f = \begin{bmatrix} 2xy^6 + 2xy^4 + 5.25y^3 - 4xy^3 + 4.5y^2 - 2xy^2 + 3y - 4xy + 6x - 12.75 \\ 6x^2y^5 + 4x^2y^3 - 6x^2y^2 - 2x^2y - 2x^2 + 15.75xy^2 + 9xy + 3x \end{bmatrix}$$

$$\nabla f(0,1) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} \quad \begin{aligned} (x^{(k+1)}, y^{(k+1)}) &= (x^{(k)}, y^{(k)}) - \alpha \nabla f(x^{(k)}, y^{(k)}) \\ &= (0,1) - \alpha (0,0) = (0,1) \end{aligned}$$

Q is positive semi definite

(2)

$$h(x) = \frac{1}{2} x^T Q x + x^T c + b$$

فرض کنیم تابع $h(x)$ بصورت زیر باشد:

$$h(x) = \frac{1}{2} x^T Q x \quad Q \text{ is PSD}$$

اعمال بردارین کاهشی:

$$x_k = x_{k-1} - \alpha \nabla h(x_{k-1})$$

$$\Rightarrow x_k = x_{k-1} - \alpha Q x_{k-1} = (I - \alpha Q) x_{k-1}$$

الگوریتم بازگشتی ادامه دهیم:

$$x_k = (I - \alpha Q)^k x_0$$

الگوریتم Q را بصورت تجزیه eigen value هایش بنویسیم:

$$Q = V \Lambda V^T \quad \begin{cases} V: \text{eigen vectors} \\ \Lambda: \text{eigen values} \end{cases}$$

$$(I - \alpha Q)^k x_0 = (I - \alpha (V \Lambda V^T))^k x_0$$

$$= (I - \alpha V \Lambda V^T)^k x_0 = (V (I - \alpha \Lambda) V^T)^k x_0$$

$$x_k = V (I - \alpha \Lambda)^k V^T x_0$$

این الگوریتم بهینه سازی بردارین کاهشی برای $h(x)$ فرضی در جهت بردارهای ویژه ماتریس Q است.

چون الگوریتم بہینہ سازی کارائی کا ہستی نسبت بہ کیفیت تغیرنا پذیر است۔
 ہیں سی تو انہم تابع $h(x)$ ، 1 بہ تابع رجبہ 2 بامدکز غیر صفر نیز تکمیل (ہم)۔

$$h(x) = \frac{1}{2} x^T Q x + a^T c + b$$

$$m^{(k)} = \beta_1 m^{(k-1)} + (1 - \beta_1) (\partial_w E)^{(k)}$$

(3) Adam الگوریتم

$$v^{(k)} = \beta_2 v^{(k-1)} + (1 - \beta_2) (\partial_w E)^{2(k)}$$

$$\hat{m}^{(k)} = \frac{\hat{m}^{(k)}}{1 - \beta_1^k}, \quad \hat{v}^{(k)} = \frac{\hat{v}^{(k)}}{1 - \beta_2^k}$$

$$w^{(k+1)} = w^{(k)} - \frac{\eta}{\sqrt{\hat{v}^{(k)} + \epsilon}} \hat{m}^{(k)}$$

$$\odot \partial_w E^{(k)} = g_k$$

$$v_k = \beta_2 v_{k-1} + (1 - \beta_2) g_k^2 \Rightarrow v_k = (1 - \beta_2) \sum_{j=1}^k \beta_2^{k-j} g_j^2$$

از 2-norm بہ تغیری (ہم) : ∞ -norm

$$v_k = (1 - \beta_2^p) \sum_{j=1}^k \beta_2^{p(k-j)} |g_j|^p$$

$$\lim_{p \rightarrow \infty} (v_k)^{1/p} = \lim_{p \rightarrow \infty} \left((1 - \beta_2^p) \sum_{j=1}^k \beta_2^{p(k-j)} |g_j|^p \right)^{1/p} =$$

$$\lim_{p \rightarrow \infty} \left(\sum_{j=1}^k (\beta_2^{k-j} |g_j|)^p \right)^{1/p} = \max_j (\beta_2^{k-j} |g_j|)$$

آلگوریتم ∞ -norm برای k در Step k برابر u_k به‌صورت زیر تعریف می‌شود:

$$u^{(k)} = \max(\beta_2 \cdot u^{(k-1)}, |g_k|)$$

پس الگوریتم جدیدی شکل می‌گیرد و خواهد بود.

$$m^{(k)} = \beta_1 m^{(k-1)} + (1 - \beta_1) (\partial_w E)^{(k)}$$

$$u^{(k)} = \max(\beta_2 u^{(k-1)}, |(\partial_w E)^k|)$$

$$\theta^{(k)} = \theta^{(k-1)} - \frac{\eta}{(1 - \beta_1^k) u^{(k)}} m^{(k)}$$

برگرفته از مقاله Adam: a method for stochastic optimization

• الگوریتم حاصل شده (Adamax) در زمان‌هایی که گرادینت‌ها بزرگ هستند stable تر است یا زمانی که مشخصه‌های داده در طول زمان تغییرپذیرند و نیاز است با توجه به آن learning rate را تغییر دهیم، مناسب است.

الگوریتم Adam بصورت عمومی optimizer مناسبی است. علی‌الحضوف زمانی که نیاز است بزرگ و پارامترهای زیادی داریم. این الگوریتم تغییرپذیر نسبت به rescale شدن گرادینت‌ها است و این ویژگی زمانی که گرادینت‌ها بین مسئله‌های متفاوت بصورت قابل توجه تغییر می‌کنند یا حتی در بخش‌های مختلف یک مسئله متغیر هستند، کمک می‌کند.