

1 Optimization

1-1 Stochastic Gradient Descent

1-1-1

$$\begin{aligned}
 L_i(x_i, w_t) &= \frac{1}{2} \|x_i w_t - t_i\|^2 \\
 &= \frac{1}{2} (x_i w_t - t_i)^T (x_i w_t - t_i) \\
 &= \frac{1}{2} (w_t^T x_i^T - t_i^T) (x_i w_t - t_i) \\
 &= \frac{1}{2} (w_t^T x_i^T x_i w_t - w_t^T x_i^T t_i - t_i^T x_i w_t + t_i^T t_i).
 \end{aligned}$$

$$\Rightarrow \frac{\partial L_i}{\partial w_t} = x_i^T x_i w_t - x_i^T t_i$$

$$\text{Initialization: } w_0 = 0 \Rightarrow \frac{\partial L_i}{\partial w_0} = -x_i^T t_i$$

$$\begin{aligned}
 \text{first iteration: } w_1 &\leftarrow w_0 - \eta \frac{\partial L_i}{\partial w_0} = 0 + \eta x_i^T t_i = \eta x_i^T t_i = \alpha x_i^T \quad \alpha \in \mathbb{R} \\
 \frac{\partial L_i}{\partial w_1} &= x_i^T x_i (\alpha x_i^T) - x_i^T t_i
 \end{aligned}$$

$$\begin{aligned}
 \text{second iteration: } w_2 &\leftarrow w_1 - \eta \frac{\partial L_i}{\partial w_1} = \alpha x_i^T - \eta (x_i^T x_i (\alpha x_i^T) - x_i^T t_i) \\
 &= \alpha x_i^T - \eta (\alpha x_i^T x_i x_i^T - t_i x_i^T) \\
 &= \underbrace{[\alpha - \eta (\alpha x_i^T x_i^T - t_i)]}_{\text{constant}} x_i^T \\
 &= b x_i^T \quad \text{for } b \in \mathbb{R}
 \end{aligned}$$

\Rightarrow By definition of stochastic gradient descent, x_i is sampled from X , therefore, x_i is contained in row space of X ;

\Rightarrow As we can see from above gradient updates, each weight update always lies in the span X since it's a linear combination of row vectors within X

\Rightarrow therefore, assume the final weight be $\hat{w} = X^T c$ for $c \in \mathbb{R}^n$

$$\Rightarrow X \hat{w} - t = X X^T c - t = 0 \Rightarrow c = (X X^T)^{-1} t$$

$$\Rightarrow \hat{w} = \underbrace{X^T (X X^T)^{-1} t}_{\text{Some answer as } w^* \text{ in HW1}} \quad ①$$

Some answer as w^* in HW1

• proof of min-norm:

\Rightarrow still, assume the optimization problem has a solution $\hat{w} = X^T c \in \text{span } X$ and let w_t be any zero-loss solution to the model

\Rightarrow therefore, we have: $Xw_t = t \Rightarrow w_t = X^{-1}t \Rightarrow w_t^T = t^T(X^{-1})^T$

According to ①: $\hat{w}^T = t^T [(X(X^T)^{-1})^T]^T X$

$$\begin{aligned} ②: \quad \hat{w}^T \hat{w} &= t^T [(X(X^T)^{-1})^T]^T X X^T (X(X^T)^{-1})^T t \\ &= t^T (X(X^T)^{-1})^T (X X^T) (X(X^T)^{-1})^T t = t^T (X(X^T)^{-1})^T t. \end{aligned}$$

$$③: \quad w_t^T \hat{w} = t^T (X^{-1})^T X^T (X(X^T)^{-1})^T t = t^T (X(X^T)^{-1})^T t.$$

$$② - ③: \quad (\hat{w} - w_t)^T \hat{w} = 0$$

By Pythagorean theorem: $\|w_t\|^2 = \|\hat{w} - w_t\|^2 + \|\hat{w}\|^2 \geq \|\hat{w}\|^2$



\Rightarrow therefore, \hat{w} has the min norm solution of all solutions $\Rightarrow \hat{w} = w^*$

1.1.2

Initialization: $w_0 = 0, \delta_0 = 0 \Rightarrow \frac{\partial L_i}{\partial w_0} = -X_{i:}^T t_i \Rightarrow \delta_1 = -\eta \frac{\partial L_i}{\partial w_0} = \eta X_{i:}^T t_i$.

First iteration: $w_1 \leftarrow w_0 + \delta_1 = 0 + \eta X_{i:}^T t_i = \eta X_{i:}^T t_i = a X_{i:}^T \quad a \in \mathbb{R}$

$$\frac{\partial L_i}{\partial w_1} = a X_{i:}^T X_{i:} X_{i:}^T - X_{i:}^T t_i, \quad \delta_2 = a \eta X_{i:}^T X_{i:} X_{i:}^T - (\eta - a\eta) X_{i:}^T t_i.$$

Second iteration: $w_2 \leftarrow w_1 - \eta \frac{\partial L_i}{\partial w_1} = a X_{i:}^T - a \eta X_{i:}^T X_{i:} X_{i:}^T + (\eta - a\eta) X_{i:}^T t_i.$

$$= a X_{i:}^T - [a \eta X_{i:}^T X_{i:}^T - (\eta - a\eta) t_i] X_{i:}^T$$

$$= [a - a \eta X_{i:}^T X_{i:}^T + (\eta - a\eta) t_i] X_{i:}^T$$

constant

$$= b X_{i:}^T$$

\Rightarrow the weight updates still follow the span of X which means momentum does not influence the linearity of SGD update steps, therefore, SGD with momentum will still converge to the min norm solution.

1.2 Adaptive Method

1.2.1 take a counter-example: $x_i = [2, 1]$, $w_0 = [0, 0]$, $t = 2$
 $G_{i-1} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$.

$$\frac{\partial L_i}{\partial w_i} = X_i^T X_i w_i - X_i^T t_i$$

initialization: $\frac{\partial L_i}{\partial w_0} = -X_i^T t_i = \begin{bmatrix} -4 \\ -2 \end{bmatrix} \cdot G_0 = \begin{bmatrix} 16 \\ 4 \end{bmatrix}$

first iteration: $w_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} \frac{-4\eta}{\sqrt{16} + \epsilon} \\ \frac{-2\eta}{\sqrt{4} + \epsilon} \end{bmatrix} \approx \begin{bmatrix} \eta \\ \eta \end{bmatrix}$ not a span of X_i

$$\frac{\partial L_i}{\partial w_1} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \end{bmatrix} \begin{bmatrix} \eta \\ \eta \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \times 2 = \begin{bmatrix} 6\eta - 4 \\ 3\eta - 2 \end{bmatrix}$$

$$G_{i-1} = \begin{bmatrix} 16 \\ 4 \end{bmatrix} + \begin{bmatrix} (6\eta - 4)^2 \\ (3\eta - 2)^2 \end{bmatrix} = \begin{bmatrix} (6\eta - 4)^2 + 16 \\ (3\eta - 2)^2 + 4 \end{bmatrix}$$

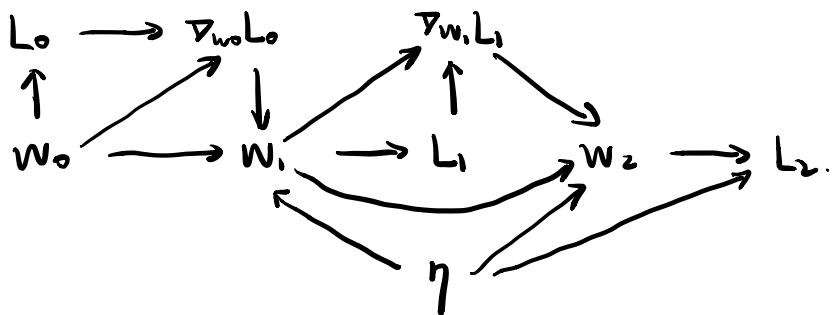
second iteration: $w_2 = \begin{bmatrix} \eta \\ \eta \end{bmatrix} - \begin{bmatrix} \frac{(6\eta - 4)\eta}{\sqrt{(6\eta - 4)^2 + 16} + \epsilon} \\ \frac{(3\eta - 2)\eta}{\sqrt{(3\eta - 2)^2 + 4} + \epsilon} \end{bmatrix}$ also, not a span of X_i

\Rightarrow therefore, the weight update no longer follows the span of X_i , which may due to the adaptive gradient influences the linearity, and won't follow the minimum norm in the end.

2. Gradient-based hyperparameter optimization.

2.1 Computation Graph

(2.1.1)



(2.1.2)

Forward: $O(1)$

Backward: $O(t)$ since η goes into every weight update
and it needs to store all variables for t iterations.

2.2 Optimal Learning Rates

(2.2.1)

$$\begin{aligned} L_0 &= \frac{1}{2} \|Xw_0 - t\|_2^2 \\ &= \frac{1}{2} (Xw_0 - t)^T (Xw_0 - t) \\ &= \frac{1}{2} (w_0^T X^T - t^T) (Xw_0 - t) \\ &= \frac{1}{2} (w_0^T X^T X w_0 - w_0^T X^T t - t^T X w_0 + t^T t) \end{aligned}$$

$$\Rightarrow \frac{\partial L_0}{\partial w_0} = X^T X w_0 - X^T t$$

$$\begin{aligned} \Rightarrow w_1 &\leftarrow w_0 - \eta \frac{\partial L_0}{\partial w_0} = w_0 - \eta X^T (Xw_0 - t) \\ &= w_0 - \eta X^T a \quad a \in \mathbb{R}^n \end{aligned}$$

$$\Rightarrow L_1 = \frac{1}{2} \|Xw_1 - t\|_2^2$$

$$= \frac{1}{2} \|Xw_0 - \eta X^T a - t\|_2^2$$

$$= \frac{1}{2} \|a - \eta X^T a\|^2$$

$$\begin{aligned}
&= \frac{1}{2} (a - \eta X X^T a)^T (a - \eta X X^T a) \\
&= \frac{1}{2} (a^T - \eta a^T X X^T) (a - \eta X X^T a) \\
&= \frac{1}{2} (a^T a - \eta a^T X X^T a - \eta a^T X X^T a + \eta^2 a^T X X^T X X^T a) \\
&= \frac{1}{2} (a^T a - 2\eta a^T X X^T a + \eta^2 a^T X X^T X X^T a) \\
&= \frac{1}{2} a^T (-\eta X X^T + I)^2 a
\end{aligned}$$

$I \in R^{n \times n}$

(2.2.2)

$$\frac{\partial L_1}{\partial \eta} = a^T (-\eta X X^T + I) (-X X^T) a$$

$$\frac{\partial^2 L_1}{\partial \eta^2} = a^T (-X X^T)^2 a \Rightarrow \text{convex, since second derivative is non-negative.}$$

(2.2.3)

$$\frac{\partial L_1}{\partial \eta} = a^T (-\eta X X^T + I) (-X X^T) a$$

$$\Rightarrow a^T (-\eta X X^T + I) (-X X^T) a = 0.$$

$$(-\eta a^T X X^T + a^T) (-X X^T) a = 0.$$

$$(\eta a^T X X^T X X^T - a^T X X^T) a = 0.$$

$$\eta a^T X X^T X X^T a = a^T X X^T a$$

$$\eta (X X^T a)^2 = (X^T a)^2$$

$$\Rightarrow \eta = \frac{(X^T a)^2}{(X X^T a)^2}$$

2.3 Multiple Inner-loop Iterations

(2.3.1) $w_1 = w_0 - \eta X^T a \quad L_1 = \frac{1}{2} a^T (-\eta X X^T + I)^2 a$

$$\Rightarrow \frac{\partial L}{\partial w_1} = X^T X w_1 - X^T t = X^T X (w_0 - \eta X^T a) - X^T t$$

$$w_2 \leftarrow w_1 - \eta \frac{\partial L}{\partial w_1} = w_1 - \eta X^T X (w_0 - \eta X^T a) + \eta X^T t$$

$$= w_0 - \eta X^T a - \eta X^T X (w_0 - \eta X^T a) + \eta X^T t$$

$$= w_0 - \eta X^T X w_0 - \eta X^T a + \eta^2 X^T X X^T a + \eta X^T t$$

Finding pattern:

$$\begin{aligned}
L_2 &= \frac{1}{2} \|xw_2 - t\|_2^2 \\
&= \frac{1}{2} \left\| \underbrace{xw_0}_{\text{constant}} - \eta \underbrace{xx^T}_{\text{constant}} \underbrace{xw_0}_{\text{constant}} - \eta xx^T a + \eta^2 xx^T xx^T a + \eta \underbrace{xx^T t - t}_{\text{constant}} \right\|_2^2 \\
&= \frac{1}{2} \left\| a - \eta xx^T a - \eta xx^T a + \eta^2 xx^T xx^T a \right\|_2^2 \\
&= \frac{1}{2} \left\| a - 2\eta xx^T a + \eta^2 xx^T xx^T a \right\|_2^2 \\
&= \frac{1}{2} (a^T - 2\eta a^T xx^T + \eta^2 a^T xx^T xx^T) (a - 2\eta xx^T a + \eta^2 xx^T xx^T a) \\
&= \frac{1}{2} (a^T a - 2\eta a^T xx^T a + \eta^2 a^T xx^T xx^T a - 2\eta a^T xx^T a + 4\eta^2 a^T xx^T xx^T a \\
&\quad - 2\eta^3 a^T xx^T xx^T xx^T a + \eta^2 a^T xx^T xx^T a - 2\eta^3 a^T xx^T xx^T xx^T a + \eta^4 a^T xx^T xx^T xx^T xx^T a) \\
&= \frac{1}{2} (a^T a - 4\eta a^T xx^T a + 6\eta^2 a^T xx^T xx^T a - 4\eta^3 a^T xx^T xx^T xx^T a + \eta^4 a^T xx^T xx^T xx^T xx^T a \\
&\quad + \eta^4 a^T xx^T xx^T xx^T xx^T a) \\
&= \frac{1}{2} a^T (I - 4\eta xx^T + 6\eta^2 xx^T xx^T - 4\eta^3 xx^T xx^T xx^T + \eta^4 xx^T xx^T xx^T xx^T) a. \\
\Rightarrow L_t &= \frac{1}{2} a^T (-\eta xx^T + I)^{2t} a.
\end{aligned}$$

By observation, the constants for each term in the brackets are binomial coefficients
a common pattern can be found $\Rightarrow L_t = \frac{1}{2} a^T (-\eta xx^T + I)^{2t} a.$

proof by induction:

- Base: $t=0 \Rightarrow L_0 = \frac{1}{2} a^T a = \frac{1}{2} \|xw_0 - t\|^2$

- Assume $L_t = \frac{1}{2} a^T (-\eta xx^T + I)^{2t} a$

- Induction: need to prove $L_{t+1} = \frac{1}{2} a^T (-\eta xx^T + I)^{2(t+1)} a.$

from the assumption: $L_t = \frac{1}{2} \|xw_t - t\|^2 = \frac{1}{2} a^T (-\eta xx^T + I)^{2t} a$
 $\Rightarrow xw_t - t = (-\eta xx^T + I)^t a$

$$\frac{\partial L_t}{\partial w_t} = x^T x w_t - x^T t = x^T (xw_t - t) = x^T (-\eta xx^T + I)^t a$$

$$\Rightarrow w_{t+1} \leftarrow w_t - \eta \frac{\partial L_t}{\partial w_t} = w_t - \eta x^T (-\eta xx^T + I)^t a.$$

$$\begin{aligned}
\Rightarrow L_{t+1} &= \frac{1}{2} \| Xw_{t+1} - t \|^2 \\
&= \frac{1}{2} \| X(Xw_t - \eta X^T(-\eta X X^T + I)^{-1} a) - t \|^2 \\
&= \frac{1}{2} \| Xw_t - \eta X^T(-\eta X X^T + I)^{-1} a - t \|^2 \\
&= \frac{1}{2} \| (-\eta X X^T + I)^{-1} a - \eta X^T(-\eta X X^T + I)^{-1} a - t \|^2 \\
&= \frac{1}{2} \| (-\eta X X^T + I)^{t+1} a \|^2 = \frac{1}{2} a^T (-\eta X X^T + I)^{2(t+1)} a
\end{aligned}$$

QED!

2.3.2 By spectral decomposition $\Rightarrow X X^T = Q \Lambda Q^T$

$$\Rightarrow L_t = \frac{1}{2} a^T (-\eta Q \Lambda Q^T + I)^{2t} a = \frac{1}{2} a^T Q (-\eta \Lambda + I)^{2t} Q^T a$$

since Λ is a diagonal matrix and L_t will result in a scalar value:

$$\Rightarrow L_t = \frac{1}{2} \sum_{i=1}^n c_i^2 (-\eta \lambda_i + 1)^{2t} \quad \text{where } c_i \text{ is the } i^{\text{th}} \text{ element of } a^T Q \\ \lambda_i \text{ is the } i^{\text{th}} \text{ element in the diagonal matrix.}$$

$$\Rightarrow \frac{\partial L_t}{\partial \eta} = \sum_{i=1}^n t c_i^2 (-\eta \lambda_i + 1)^{2t-1}$$

$$\Rightarrow \frac{\partial^2 L_t}{\partial \eta^2} = \sum_{i=1}^n t(2t-1) c_i^2 (-\eta \lambda_i + 1)^{2t-2} \geq 0 \quad \text{since} \left\{ \begin{array}{l} t(2t-1) \geq 0 \\ c_i^2 \geq 0 \\ (-\eta \lambda_i + 1)^{2t-2} \geq 0 \text{ because of the square.} \end{array} \right.$$

3. Convolutional Neural Networks

3.1 Convolutional filters

$$\Rightarrow I * J = \begin{bmatrix} 1 & 2 & 2 & -2 & 0 \\ -2 & 1 & 0 & 2 & 1 \\ 3 & 0 & 0 & 1 & -1 \\ -2 & 2 & 0 & 2 & -1 \\ 0 & -2 & 3 & -2 & 0 \end{bmatrix}$$

3.2 Size of conv nets

Output size:

$$\text{conv3-64} : (112 - 3 + 2 \times 1) + 1 = 112$$

$$\text{maxpool} : (112 - 2) / 2 + 1 = 56.$$

$$\text{conv3-128} : (56 - 3 + 2 \times 1) + 1 = 56.$$

$$\text{maxpool} : (56 - 2) / 2 + 1 = 28$$

$$\text{conv3-256} : (28 - 3 + 2 \times 1) + 1 = 28$$

$$\text{conv3-256} : (28 - 3 + 2 \times 1) + 1 = 28.$$

$$\text{maxpool} : (28 - 2) / 2 + 1 = \underline{\underline{14}}.$$

1) Number of parameters.

$$\text{Conv3-64} : 3 \times 3 \times 3 \times 64 + 64 = 1792$$

$$\text{Conv3-128} : 3 \times 3 \times 64 \times 128 + 128 = 73856.$$

$$\text{Conv3-256} : 3 \times 3 \times 128 \times 256 + 256 = 296148$$

$$\text{Conv3-256} : 3 \times 3 \times 256 \times 256 + 256 = 590080.$$

$$\text{FC-1024} : 14 \times 14 \times 256 \times 1024 + 1024 = 51381248$$

$$\text{FC-100} : 1024 \times 100 + 100 = 102500$$

$$\Rightarrow \underline{\underline{\text{Summ}}} = 52444644. \text{ parameters}$$

2) Number of neurons.

$$\text{conv3-64} : 112 \times 112 \times 64 = 802816$$

$$\text{maxpool} : 56 \times 56 \times 64 = 200704$$

$$\text{conv3-128} : 56 \times 56 \times 128 = 401408$$

$$\text{maxpool} : 28 \times 28 \times 128 = 100352$$

$$\text{conv3-256} : 28 \times 28 \times 256 = 200704$$

$$\text{conv3-256} : 28 \times 28 \times 256 = 200704$$

$$\text{maxpool} : 14 \times 14 \times 256 = 50176$$

$$\text{FC-1024} : 1024$$

$$\text{FC-100} : 100$$

$\Rightarrow \underline{\underline{1957988 \text{ neurons.}}}$

3) Number of connections.

$$\text{conv3-64} : 112 \times 112 \times 3 \times 3 \times 64 \times 3 = 21676032$$

$$\text{maxpool} : 56 \times 56 \times 2 \times 2 \times 64 = 802816$$

$$\text{conv3-128} : 56 \times 56 \times 3 \times 3 \times 128 \times 64 = 231211008$$

$$\text{maxpool} : 28 \times 28 \times 2 \times 2 \times 128 = 401408$$

$$\text{conv3-256} : 28 \times 28 \times 3 \times 3 \times 256 \times 128 = 231211008$$

$$\text{conv3-256} : 28 \times 28 \times 3 \times 3 \times 256 \times 256 = 462422016$$

$$\text{maxpool} : 14 \times 14 \times 2 \times 2 \times 256 = 200704$$

$$\text{FC-1024} : 14 \times 14 \times 256 \times 1024 = 51380224$$

$$\text{FC-100} : 1024 \times 100 = 102400$$

$\Rightarrow \underline{\underline{999407616 \text{ connections.}}}$

3.3 Receptive Field

1. Depth of convolutional architecture. If a CNN has deeper or more layers, the RF will be larger since the input one feature maps from previous layers and it gets downsampled every layer;
2. kernel size. A larger kernel size will "extract" more information from the input and make the RF become bigger in size.
3. stride. Higher stride will increase the size of RF since it indicates the "jump" between each region and bigger jumps will make the entire RF increase its size by certain proportions.