

1. Optimization

$$J = \frac{1}{n} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{t}\|_2^2, \quad \mathbf{x} \in \mathbb{R}^{n \times d}, \quad \hat{\mathbf{w}} \in \mathbb{R}^{d \times 1}, \quad \mathbf{t} \in \mathbb{R}^{n \times 1}$$

1.1. Stochastic Gradient Descent (SGD)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla_{\mathbf{w}_t} I_i(\mathbf{x}_i, \mathbf{w}_t), \quad 1 \leq i \leq n.$$

1.1.1 Minimum Norm Solution

result from HW1, Q3.4: $\frac{\partial L}{\partial \hat{\mathbf{w}}} = \bar{\mathbf{w}} = \frac{\partial}{\partial \hat{\mathbf{w}}} \left(\frac{1}{n} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{t}\|_2^2 \right)$

$$= \frac{2}{n} \mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{t}), \quad \text{in the SGD case, } \bar{\mathbf{w}}_t = \frac{2}{n} \mathbf{x}_i^T (\mathbf{x}_i^T \hat{\mathbf{w}}_t - t_i)$$

Starting from zero initialization,

$$\bar{\mathbf{w}}_t(0) = -2 \mathbf{x}_i \mathbf{t}_i, \quad \text{gradient is in the direction of } \mathbf{x}_i$$

Being a row in $\mathbf{x} \in \mathbb{R}^{n \times d}$, \mathbf{x}_i^T is in the row space of \mathbf{x} , so \mathbf{x}_i is in the column space/span of \mathbf{x} .
 Therefore, gradient is in the span of \mathbf{x} .

Update steps of SGD never leaves the span of \mathbf{x} . Since for $\bar{\mathbf{w}}_t$ always in direction of \mathbf{x}_i , no matter what $(\mathbf{x}_i^T \hat{\mathbf{w}}_t - t_i)$ is.

With the assumption that $\mathbf{X}\hat{\mathbf{w}} = \mathbf{t}$ eventually,

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \bar{\mathbf{w}}_{t,i}, \quad \text{so } \mathbf{w}_{t+1} \in \text{span } \mathbf{x}, \quad \text{let } \mathbf{C} \in \mathbb{R}^{n \times 1} \text{ denote the fixed coefficients}$$

$$\mathbf{x} \mathbf{C}^T = \mathbf{t}$$

$$\mathbf{C} = (\mathbf{x} \mathbf{x}^T)^{-1} \mathbf{t}$$

$$\mathbf{w}^{**} = \mathbf{x}^T (\mathbf{x} \mathbf{x}^T)^{-1} \mathbf{t} = \mathbf{w}^*, \quad \text{same as before!}$$

1.1.2 Mini-batch SGD

$\mathbf{B} \in \mathbb{R}^{b \times d}$, the loss at each timestep is: $I(\mathbf{B}, \mathbf{w}_t)$

the question boils down to whether the update $\nabla_{\mathbf{w}_t} I(\mathbf{B}, \mathbf{w}_t) \in \text{span } \mathbf{x}$

$$\frac{\partial L}{\partial \hat{\mathbf{w}}_t} = \bar{\mathbf{w}}_t = \frac{2}{b} \mathbf{B}^T (\mathbf{B} \hat{\mathbf{w}}_t - \mathbf{t}_b),$$

$$\bar{\mathbf{w}}_t(0) = -\frac{2}{b} \mathbf{B}^T \mathbf{t}_b, \quad \text{derivative in the span of } \mathbf{B}_i,$$

$\mathbf{B}_i^T \mathbf{t}_i$ is a linear combination of columns in \mathbf{x} , so still in span \mathbf{x} ,

using the same reasoning as before, yes.

1.2 Adaptive Methods

$$G_{i,t} = G_{i,t-1} + (\nabla_{\mathbf{w}_{i,t}} I(\mathbf{w}_{i,t}))^2$$

$$\mathbf{w}_{i,t+1} = \mathbf{w}_{i,t} - \frac{\eta}{\sqrt{G_{i,t}} + \epsilon} \nabla_{\mathbf{w}_{i,t}} I(\mathbf{w}_{i,t})$$

1.2.1 Minimum Norm Solution

Similar as before, $\bar{\mathbf{w}} = \frac{2}{b} \begin{bmatrix} 2 \\ 1 \end{bmatrix} (0 - 2) = \begin{bmatrix} -8 \\ -4 \end{bmatrix}$ at zero initialization,

update is different for the two weights,

$$\left. \begin{aligned} \text{for first weight: } w_{1,1} &= w_{1,0} - \frac{\eta}{\sqrt{G_{1,1}} + \epsilon} (-8) \\ \text{for second weight: } w_{2,1} &= w_{2,0} - \frac{\eta}{\sqrt{G_{2,1}} + \epsilon} (-4) \end{aligned} \right\} \text{this descent roughly moves} \\ \text{in the } \begin{bmatrix} 2 \\ 1 \end{bmatrix} \text{ direction, which} \\ \text{is different from before, } \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

So, no longer a minimum since it's not perpendicular to the empirical minimizer line.

1.2.2

No, both RMSprop and Adam try to rescale the gradient to have norm 1 on average to improve the slow update problem in low curvature dimension.

2. Gradient-based Hyper-parameter

2.1 Computation Graph of Learning Rates.

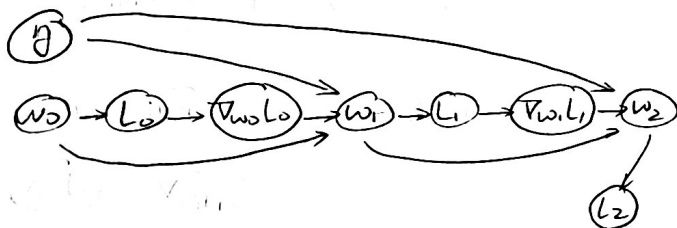
2.1.1

$$w_1 = w_0 - \eta \nabla_{w_0} L_0, \quad \nabla_{w_0} L_0 = \frac{2}{n} \lambda^T (\lambda w_0 - t)$$

$$w_2 = w_1 - \eta \nabla_{w_1} L_1, \quad \nabla_{w_1} L_1 = \frac{2}{n} \lambda^T (\lambda w_1 - t)$$

$$L_2 = \frac{1}{n} \|\lambda w_2 - t\|_2^2$$

Computation graph:



2.1.2

memory complexity for:

forward-propagation: $O(1)$.

back-propagation: $O(t)$.



t of such connections, all need to be stored in memory

2.1.3.

Since hyper-parameters are shared by all iterations, $\frac{\partial w_i}{\partial \eta}$ need to be stored for all previous iterations. For realistic training we can easily run out of memory.
 with many iterations

2.2 Learning Learning Rate.

$$w_1 = w_0 - \eta \nabla_{w_0} L_0$$

2.2.1

$$w_1 = w_0 - \eta \frac{2}{n} \lambda^T (\lambda w_0 - t)$$

$$L_1 = \frac{1}{n} \|\lambda w_1 - t\|_2^2$$

$$= \frac{1}{n} \|\lambda (w_0 - \eta \frac{2}{n} \lambda^T (\lambda w_0 - t)) - t\|_2^2 \quad \text{let } a_0 = (\lambda w_0 - t)_{n \times 1}$$

$$= \frac{1}{n} \|\lambda (w_0 - \eta \frac{2}{n} \lambda^T a_0) - t\|_2^2$$

2.2.2

$$\frac{\partial L_1}{\partial \eta} = \frac{\partial}{\partial \eta} \left[\frac{1}{n} \|\lambda w_0 - \eta \frac{2}{n} \lambda^T a_0 - t\|_2^2 \right]$$

$$= \left(\frac{4}{n^2} \lambda \lambda^T a_0 \right)^T (\lambda w_0 - \eta \frac{2}{n} \lambda^T a_0 - t)$$

$$\frac{\partial L_1}{\partial \eta^2} = \frac{4}{n^2} a_0^T \lambda \lambda^T \frac{2}{n} \lambda^T a_0$$

$$|x| = \frac{8}{n^3} (\lambda \lambda^T a_0)^T (\lambda^T a_0) > 0$$

therefore, L_1 is convex w.r.t η .

2.2.3

$$\frac{\partial L_1}{\partial \eta} = 0$$

$$(a_0^T \lambda \lambda^T) (\lambda w_0 - \eta \frac{2}{n} \lambda^T a_0 - t) = 0$$

$$a_0^T \lambda \lambda^T \lambda w_0 - \eta \frac{2}{n} a_0^T \lambda \lambda^T \lambda^T a_0 - a_0^T \lambda \lambda^T t = 0$$

$$\eta \frac{2}{n} a_0^T \lambda \lambda^T \lambda^T a_0 = a_0^T \lambda \lambda^T \lambda w_0 - a_0^T \lambda \lambda^T t$$

$$\eta \frac{2}{n} = \frac{a_0^T \lambda \lambda^T \lambda w_0 - a_0^T \lambda \lambda^T t}{a_0^T \lambda \lambda^T \lambda^T a_0}$$

$$\eta^* = \frac{n}{2} \cdot \frac{a_0^T \lambda \lambda^T \lambda w_0 - a_0^T \lambda \lambda^T t}{a_0^T \lambda \lambda^T \lambda^T a_0}$$

check for denominator $= 0$,
when $\lambda \lambda^T a_0 = 0$.

$\frac{\partial L_1}{\partial \eta} = 0$, already at optimum, trivial.

3. Convolutional Neural Networks.

3.1. Convolutional Filters.

$$\text{padding width} = \frac{k-1}{2} = \frac{3-1}{2} = 1.$$

$$I = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$J = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

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

3.2. Size of ConvNets.

| Layer | # Params. | |
|-----------|--|---|
| conv3-64 | $64 \times (3 \times 3 \times 3 + 1) \cdot \sqrt{= 1792}$ | RGB image - $112 \times 112 \times 3$. |
| max pool | 0 | kernel size 3×3 . |
| conv3-128 | $128 \times (3 \times 3 \times 64 + 1) \cdot \sqrt{= 73856}$ | |
| max pool | 0 | |
| conv3-256 | $256 \times (3 \times 3 \times 128 + 1) \cdot \sqrt{= 295168}$ | |
| conv3-256 | $256 \times (3 \times 3 \times 256 + 1) \cdot \sqrt{= 590080}$ | |
| max pool | 0 | |
| FC-1024 | $(256 \times 256) \times 1024 = 51380224$ | |
| FC-100 | $1024 \times 100 + 100$ | |
| Softmax | 0 | |
| Total. | $52493520 + 112p = 52494644$ | |