



Workshop 7

COMP90051 Statistical Machine Learning
Semester 2, 2024

Learning Outcomes

By the end of this workshop you should be able to:

1. Understand and implement different **optimization algorithms**
2. Understand the Pros and Cons of different common used optimization algorithms
3. Be able to use optimizer in Pytorch
4. Understand and implement **autoencoder**

Optimization Algorithms

A common pipeline in ML is:

1. Define the model

2. Define the loss function

3. **Optimize** the model parameter w.r.t the loss function

How to optimize/update the model parameters?

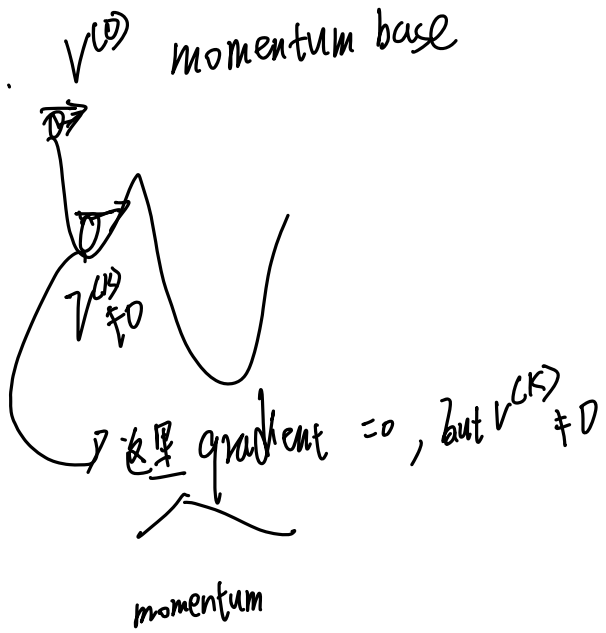
Optimization algorithms: GD, SGD, AdaGrad, Adam, ...

↓
adaptive

Handwritten notes:
ridge, linear \Rightarrow analytic solution
不用 opt Alg 的

non convex

update momentum depend on gradient.



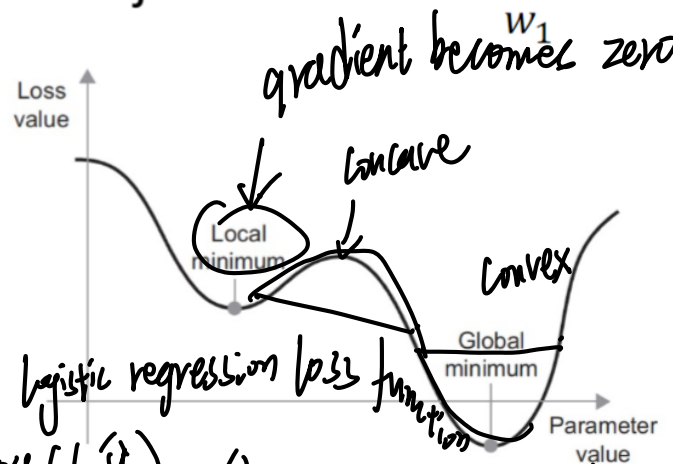
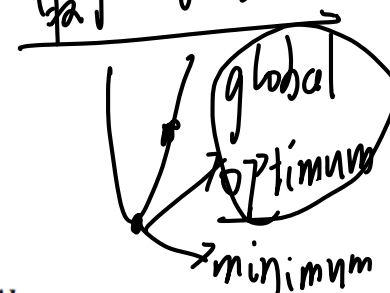
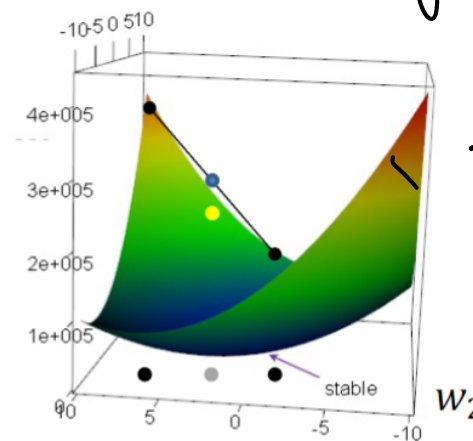
Messy Loss Function

In general, the loss function is complicated (not convex nor concave), especially for Deep Learning models.

How to update the parameter is important!

- Recall linear regression, convex '**Bowl shaped**' objective
 - * gradient descent finds a **global** optimum
- In contrast, most DNN objectives are **not convex**

- * gradient methods get trapped in **local optima** or **saddle points**



always convex for logistic regression loss function

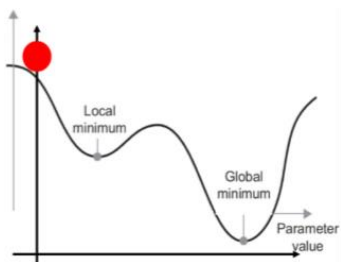
$$J = -\frac{1}{N} \sum_{i=1}^N [y \log \hat{y} + (1-y) \log (1-\hat{y})], \quad \hat{y} = \sigma(w^T x) \rightarrow \text{sigmoid} \Rightarrow \text{convex}$$

Better than GD/SGD?!

Momentum

“How to escape local minimum?”

Consider a ball with mass rolling down the loss function surface.



- * $\theta^{(t+1)} = \theta^{(t)} - v^{(t)}$
- * $v^{(t)} = \alpha v^{(t-1)} + \eta \nabla L(\theta^{(t)})$
- * α decays the velocity, e.g., 0.9

Momentum decays over time

AdaGrad

“Why update all parameter with the same rate?”

Frequent features should update less, while infrequent features should update more.

$$* g_i^{(t)} = g_i^{(t-1)} + \nabla L(\theta^{(t)})_i^2$$

$$* \theta_i^{(t+1)} = \theta_i^{(t)} - \frac{\eta}{\sqrt{g_i^{(t)} + \epsilon}} \nabla L(\theta^{(t)})_i$$

梯度平方 (handwritten note pointing to the square term in the first equation)

防止 local optimum (handwritten note pointing to the denominator in the second equation)

Adam

“How about we combine these two together?”

Combining elements from momentum and adaptive learning rate, most common used optimizer.

- * $m^{(t)} = \beta_1 m^{(t-1)} + (1 - \beta_1) \nabla L(\theta^{(t)})$
- * $v^{(t)} = \beta_2 v^{(t-1)} + (1 - \beta_2) \nabla L(\theta^{(t)})^2$
- * $\theta^{(t+1)} = \theta^{(t)} - \frac{\eta}{\sqrt{v^{(t)} / (1 - \beta_2) + \epsilon}} \frac{m^{(t)}}{1 - \beta_1}$
- * $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

g(t) 和 momentum 的结合
前两种的结合

Autoencoder: Why reconstruct the data itself??

build encoder to shrink the data

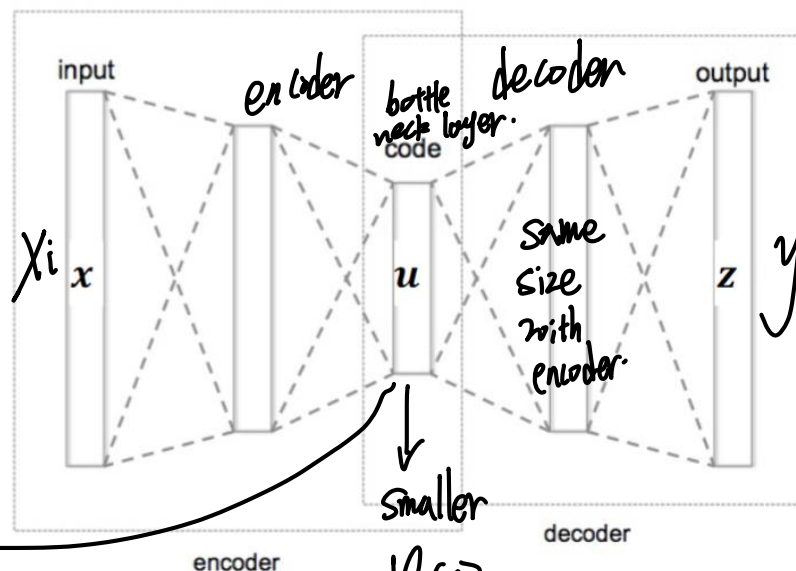
The encoder gives good **low-dimensional representation** of the data which could be used for learning tasks afterwards.

Autoencoder is **unsupervised**! No need for labelling, save time and money!

- Given data without labels x_1, \dots, x_n , set $y_i \equiv x_i$ and train a DNN to predict $z(x_i) \approx x_i$

- Set **bottleneck** layer u in middle "thinner" than input, and/or

- * corrupt input x with noise
- * regularise s.t. u is sparse
- * regularise to contract inputs



$$y_i = x_i \quad \|y - x\|^2 = \sum \|y_i - x_i\|^2$$

autoencoder

半监督PLA
降低维数还想让
他与原数据有相
同表达能力。

adapted from Chervinskii at
Wikimedia Commons (CC4)

↓ smaller
 $\psi(x)$
feature vector.

Worksheet 7

$$g^k = \nabla_n(L) \quad \text{gradient vector at iteration } k$$

$$(g^k)^T$$

$$G^i = \left(\sum_{k=1}^i (g_1^k)^2, \sum_{k=1}^i (g_2^k)^2, \dots \right)$$

> diagonal version

Sum of square of gradient

$$w^{it} = w^i - \frac{\eta}{\sqrt{G^i}} \cdot g^i$$

$$g^k \times g^k \leftarrow$$

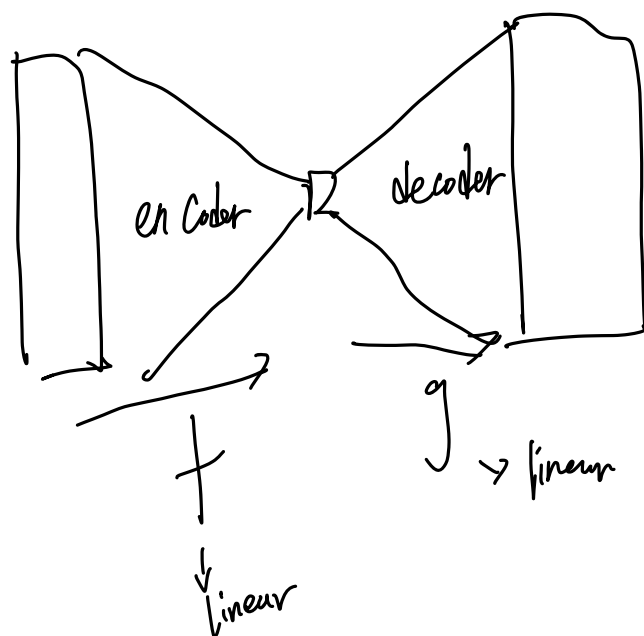
$$G^i = \begin{bmatrix} \vdots \\ g^k \\ \vdots \end{bmatrix} \begin{bmatrix} \dots & g^k \end{bmatrix} = \begin{bmatrix} \text{matrix } G \end{bmatrix} \rightarrow \text{对称矩阵}$$

~~$n \times n$~~

matrix 对称

momentum
 w is updated by $\begin{cases} \text{velocity} \\ \text{gradient} \end{cases}$ together.

Autoencoder



$$h(x) = \cancel{f \circ f(x)} \quad g \circ f(x)$$