

▷ Linear Regression on \mathbb{R}^+

①

• Problem:

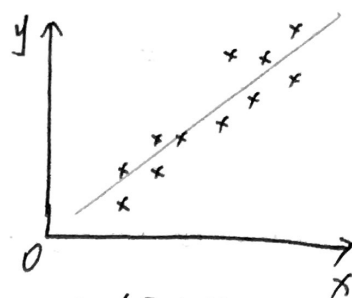
$\mathcal{D} = \{(x_i, y_i) \mid i=1, \dots, N, x_i, y_i \in \mathbb{R}^+\}$ is the ground truth data.

we want to use linear function, more specifically, Rectified linear function, to approximate the data. That is,

$$\hat{y} = \max(wx + b, 0) = \varphi(f(x))$$

Here $\varphi = \max(x, 0)$ is called Rectified Linear Unit (ReLU)

$f = wx + b$ is the linear function



To solve this problem, a "best" w and b are needed to be found. Rewrite this problem, we have

$$(w^*, b^*) = \arg \min_{w, b} L(w, b)$$

where $L(w, b)$ is a measure of how "bad" w, b are, a simple choice of which is mean squared error, i.e.,

$$L(w, b) = \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i(w, b))^2$$

Note:

• In process of finding optimal w and b , (x_i, y_i) is known, hence $\hat{y}_i = \hat{y}_i(w, b)$.

• In process of using founded w^* and b^* to predict y , (w^*, b^*) is known, hence $\hat{y} = \hat{y}(x)$

• Optimization with Gradient Descent with momentum

Init: $v_w, v_b = 0$.
choose θ, γ

Iter:

$$\begin{cases} v_w^+ = (1-\theta) \frac{\partial L}{\partial w} + \theta v_w \\ w^+ = w - \gamma \cdot v_w^+ \end{cases}$$

$$\begin{cases} v_b^+ = (1-\theta) \frac{\partial L}{\partial b} + \theta v_b \\ b^+ = b - \gamma \cdot v_b^+ \end{cases}$$

• Calculation of $\frac{\partial L}{\partial w}$, $\frac{\partial L}{\partial b}$

(2)

recall that

$$L(w, b) = \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i(w, b))^2$$

$$\text{let } \tilde{L} = (y_i - \hat{y}_i(w, b))^2$$

$$= \frac{1}{2N} \sum_{i=1}^N (y_i - \varphi(f_i(w, b)))^2$$

By chain rule,

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w} = \frac{1}{2N} \sum_{i=1}^N \frac{\partial \tilde{L}(w, b)}{\partial \varphi(f_i(w, b))} \cdot \frac{\partial \varphi(f_i(w, b))}{\partial f_i(w, b)} \cdot \frac{\partial f_i(w, b)}{\partial w} = \frac{1}{2N} \sum_{i=1}^N \frac{\partial \tilde{L}}{\partial \varphi} \cdot \frac{\partial \varphi}{\partial f_i} \cdot \frac{\partial f_i}{\partial w} \\ \frac{\partial L}{\partial b} = \frac{1}{2N} \sum_{i=1}^N \frac{\partial \tilde{L}(w, b)}{\partial \varphi(f_i(w, b))} \cdot \frac{\partial \varphi(f_i(w, b))}{\partial f_i(w, b)} \cdot \frac{\partial f_i(w, b)}{\partial b} = \frac{1}{2N} \sum_{i=1}^N \frac{\partial \tilde{L}}{\partial \varphi} \cdot \frac{\partial \varphi}{\partial f_i} \cdot \frac{\partial f_i}{\partial b} \end{array} \right.$$

where

$$\frac{\partial \tilde{L}}{\partial \varphi} = 2(\varphi(f_i(w, b)) - y_i)$$

$$\frac{\partial f_i}{\partial w} = x_i$$

$$\frac{\partial \varphi}{\partial f_i} = \begin{cases} 1 & , f_i(w, b) > 0 \\ 0 & , f_i(w, b) < 0 \end{cases}$$

$$\frac{\partial f_i}{\partial b} = 1$$

so (*) becomes

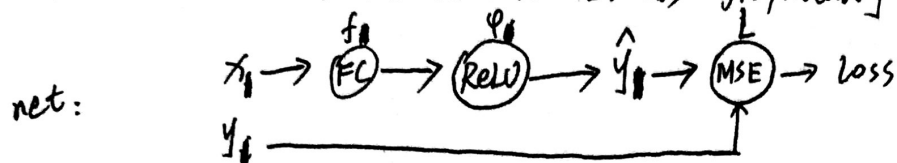
$$\frac{\partial L}{\partial w} = \begin{cases} \frac{1}{N} \sum_{i=1}^N x_i (wx_i + b) & , wx_i + b > 0 \\ 0 & , wx_i + b < 0 \end{cases}$$

$$\frac{\partial L}{\partial b} = \begin{cases} \frac{1}{N} \sum_{i=1}^N (wx_i + b) & , wx_i + b > 0 \\ 0 & , wx_i + b < 0 \end{cases}$$

- Motivation: modularize the calculation of $\frac{\partial L}{\partial w}$, $\frac{\partial L}{\partial b}$
(Engineering-oriented)

As we can see from (x), $\frac{\partial L}{\partial \varphi}$, $\frac{\partial \varphi}{\partial f_i}$ are used repeatedly. Hence things would be easier if we ~~set~~ ^{separate} the calculation of $\frac{\partial L}{\partial \varphi}$, $\frac{\partial \varphi}{\partial f_i}$, $\frac{\partial \varphi}{\partial w}$ and $\frac{\partial f_i}{\partial b}$, although not so obviously in this simplest example.

To achieve this, we rewrite $L(w, b)$ graphically.



Here $x, y, \hat{y} \in \mathbb{R}^{1 \times N}$

Roughly we define such a graph as a (training) network. (1) To compute $L(w, b)$, we call `net.forward(x, y)`, which in this case is as following

MATLAB

function loss = forward(net, x, y)

~~loss = L.forwardLoss(~~

$\hat{y} = \varphi.\text{forward}(f.\text{forward}(x))$

$\text{loss} = L.\text{forwardLoss}(\hat{y}, y)$

end

(2) To ~~be~~ update parameters w, b , we call ~~net.~~ `update_params($\gamma, \theta, y, \hat{y}$)`

MATLAB

function net = update_params(net, $\gamma, \theta, y, \hat{y}$)

$\frac{\partial L}{\partial \varphi} = L.\text{backwardLoss}(y, \hat{y})$ ~~$\in \mathbb{R}^{1 \times N}$~~ $y, \hat{y}, \frac{\partial L}{\partial \varphi} \in \mathbb{R}^{1 \times N}$

$\frac{\partial \varphi}{\partial f} = \varphi.\text{backward}(f(w, b), \frac{\partial L}{\partial \varphi})$, $f(w, b) \in \mathbb{R}^{1 \times N}$

$[\frac{\partial L}{\partial x}, \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}] = f.\text{backward}(x, \frac{\partial L}{\partial f})$, $x \in \mathbb{R}^{1 \times N}$

$f = f.\text{update_params}(\gamma, \theta, \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b})$

end

By doing this, we could separately define the forward and backward of each "layer", which is the circled part in the network: ~~FE~~, fully connected layer (FC), rectified linear unit layer (ReLU), and mean squared error (MSE). (4)

To make this simple, we give the one-variable version, which is simply a rewritten version of our previous formulas in an engineering way. Take FC layer as an example, ReLU and MSE are similar.

% FC Layer

```
classdef FC <handle
```

```
properties
```

```
    w  
    b
```

```
end
```

```
properties (private)
```

```
    v_w  
    v_b
```

```
end
```

```
methods
```

```
function layer = FC()
```

```
    layer.v_w = 0
```

```
    layer.v_b = 0
```

```
    layer.w = randn()
```

```
    layer.b = randn()
```

```
end
```

```
function z = forward(layer, x)
```

```
    z = layer.w * x + layer.b, x ∈ ℝ1×N
```

```
end
```

```
function edit dx dl
```

```
function [∂L/∂x, ∂L/∂w, ∂L/∂b] = backward(layer, x, ∂L/∂z)
```

$$\frac{\partial L}{\partial x} = \text{layer.w} \cdot x, x \in \mathbb{R}^{1 \times N}$$

$$\frac{\partial L}{\partial w} = \sum \left(\frac{\partial L}{\partial z} \right)^T \cdot x, \frac{\partial L}{\partial w} \in \mathbb{R}$$

$$\frac{\partial L}{\partial b} = \sum \frac{\partial L}{\partial z}, \frac{\partial L}{\partial b} \in \mathbb{R}$$

```
end
```

```
function new layer = update_params(layer, γ, θ, ∂L/∂w, ∂L/∂b)
```

```
layer
```

```
    layer.v_w = layer.v_w * θ + ∂L/∂w * (1-θ)
```

```
    layer.w = layer.w - γ * layer.v_w
```

```
    layer.v_b = layer.v_b * θ + ∂L/∂b * (1-θ)
```

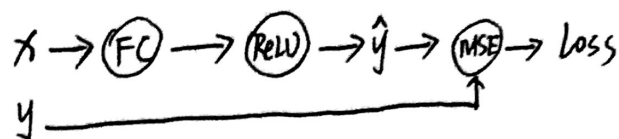
```
    layer.b = layer.b - γ * layer.v_b
```

```
end
```

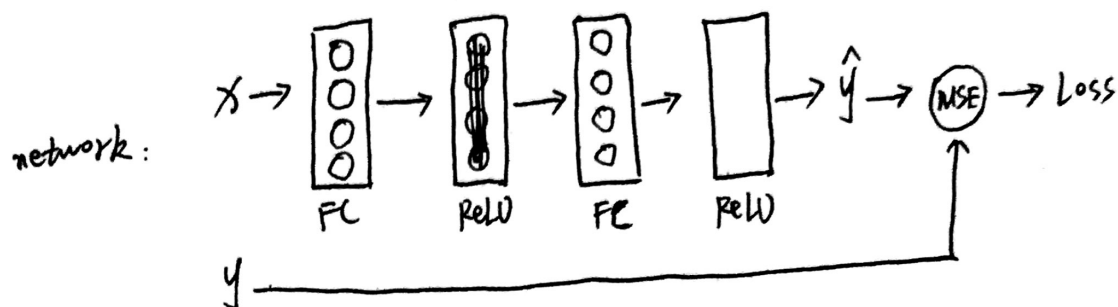
```
end
```

- Going deeper and wider

The network used in our example is width 1, depth 2 network.



we could make this ~~more~~ deeper and wider with previous work of modularization



The corresponding formular of this network is

$$y = \varphi(f_2(\varphi(f_1(x)))) \quad , x, y \in \mathbb{R}^{n \times 1}$$

$$f_i^{(0)} = w_i^T x + b \quad , w_i \in \mathbb{R}^{n \times 1} , b \in \mathbb{R} \quad , i=1,2.$$

$$\varphi(x) = \max(x, 0) \quad \text{pointwise max.}$$

$$L(w_1, b_1, w_2, b_2) = \frac{1}{2N} \sum_{i=1}^N \|y_i - \hat{y}_i\|_2^2$$

The forward ~~and~~ param-update and backward of the network could be easily done once we define the act of each type of layer.

- General framework

for any version of approximation by composition of simple functions.

$$y = \varphi_n(\varphi_{n-1}(\dots(\varphi_2(\varphi_1(x))))))$$

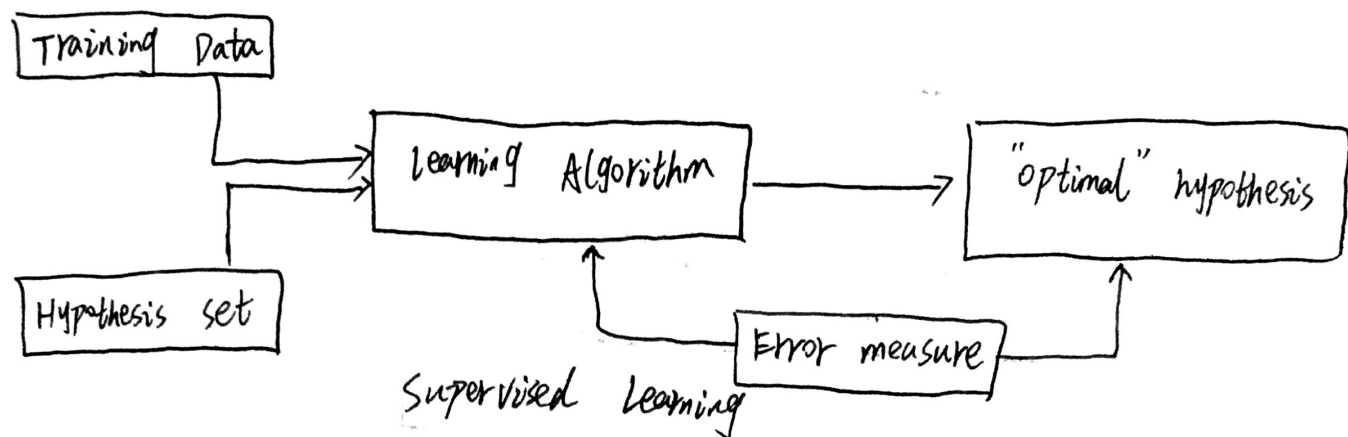
the only thing we need to do then is define:

φ_i . backward

φ_i . forward

φ_i . update_param \rightsquigarrow if necessary

Currently, linear combination and convolution are two main ~~and~~ simple function that requires parameter updates, where convolution layer is very suitable for image and computer vision tasks.



- 陆嘉生 "A primer of machine learning"

Building blocks of deep learning:

- Training data
- Validation data : not used in training, to test the performance
- Hypothesis set : the network
- Loss function : describe how "bad" the network performs
- Optimizer : learning Algorithm