

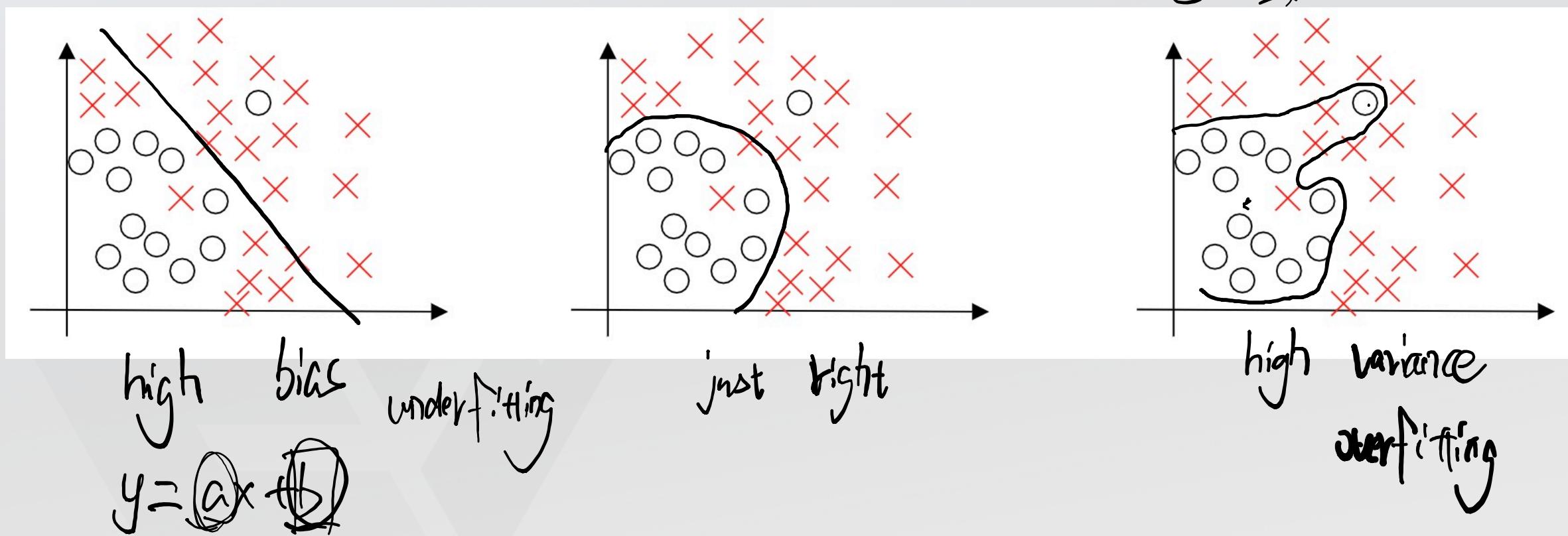
- ① 初始化参数
- ② 选择优化算法

(3)

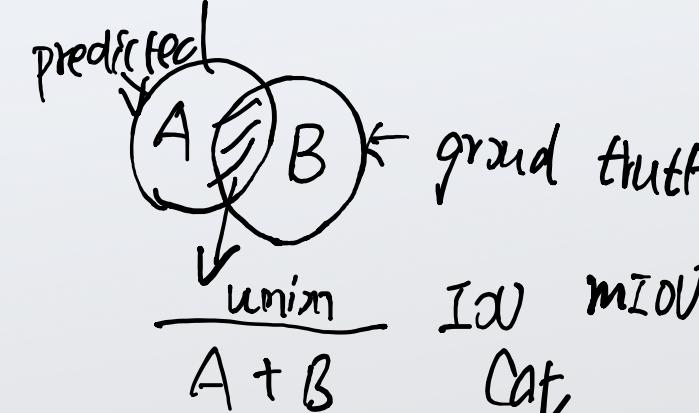
- {
- (1.) forward
- (2.) loss
- (3.) SGD
- (4.) update weight.

Weight Initialization

Bias and Variance



high bias \Rightarrow Bigger network
longer/depth mIoU
training data error



Bias and Variance

high var \Rightarrow more data

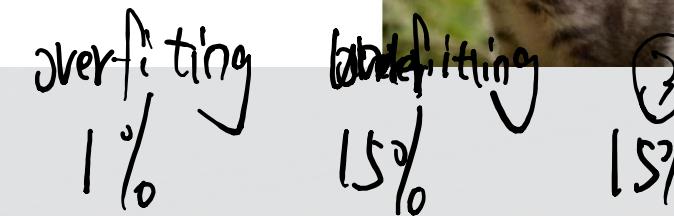
regularization

Test/validation data error.

Cat classification

Train data error :

Test(validation/dev) data error:



4
2.5%
1%
low bias
low var.

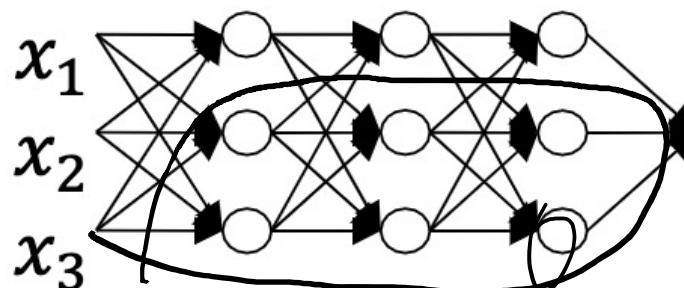
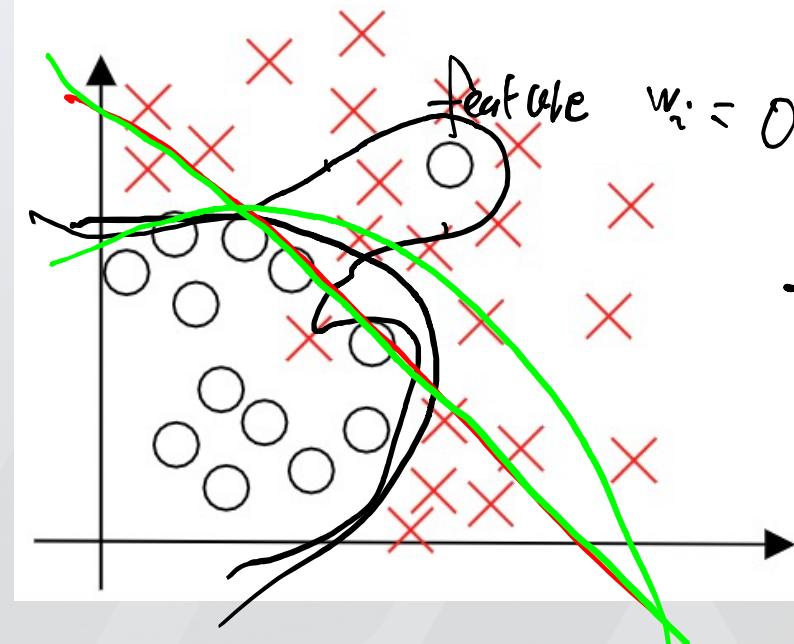
30%
high bias
high var.

16%

15%

平衡歸納度 $\hat{y} \rightarrow y$ (w)

Regularization



$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|^2$$

hyper parameter
houchangtech.com

$$L_1 = \frac{\lambda}{2m} \sum_{i=1}^m |w_i| = \frac{1}{2m} \|w\|_1 + \frac{\lambda}{2m} \|b\|^2$$

real number of sample number 0

$$L_2 = \sum_{i=1}^m w_i^2 = w^T w.$$

w_1, w_2

$$J(w, b) = L + \frac{\lambda}{2m} (|w_1| + |w_2|)$$

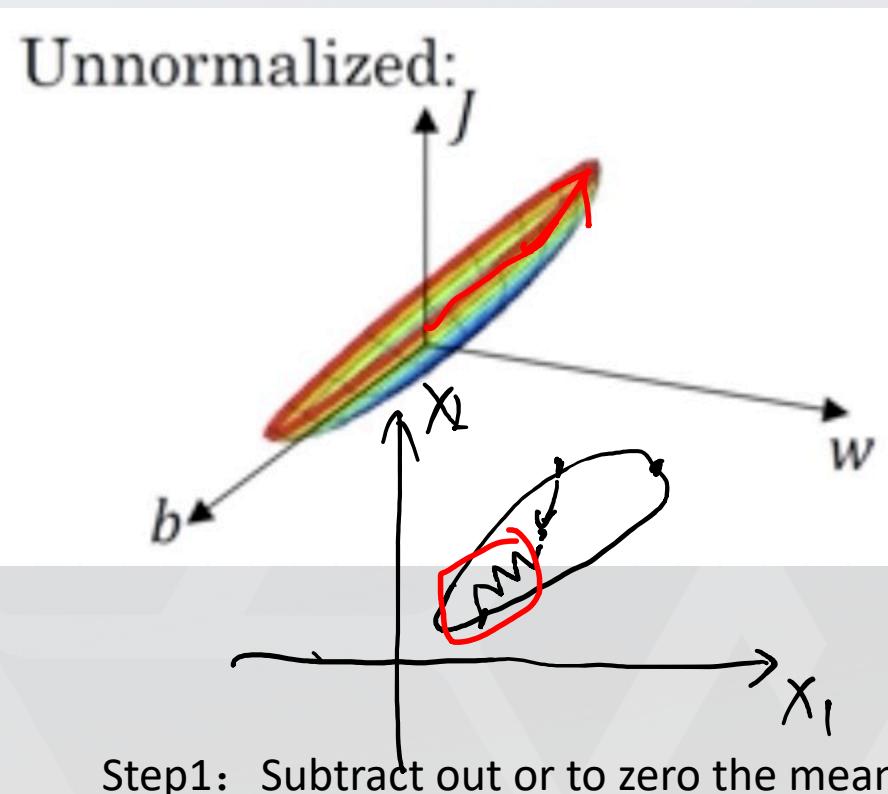
$$\text{update } w_1 : w_1 := w_1 - \alpha \frac{\partial L}{\partial w_1} - \frac{\alpha \lambda}{2m} \text{sgn}(w_1) - \frac{\partial L}{\partial w_1}$$

$$J(w, b) = L + \frac{\lambda}{2m} (w_1^2 + w_2^2)$$

$$\text{update } w_1 : (1 - \frac{\alpha}{m}) w_1 - \frac{\partial L}{\partial w_1} \Rightarrow \text{梯度下降法}$$

加快收敛速度.

Normalizing inputs

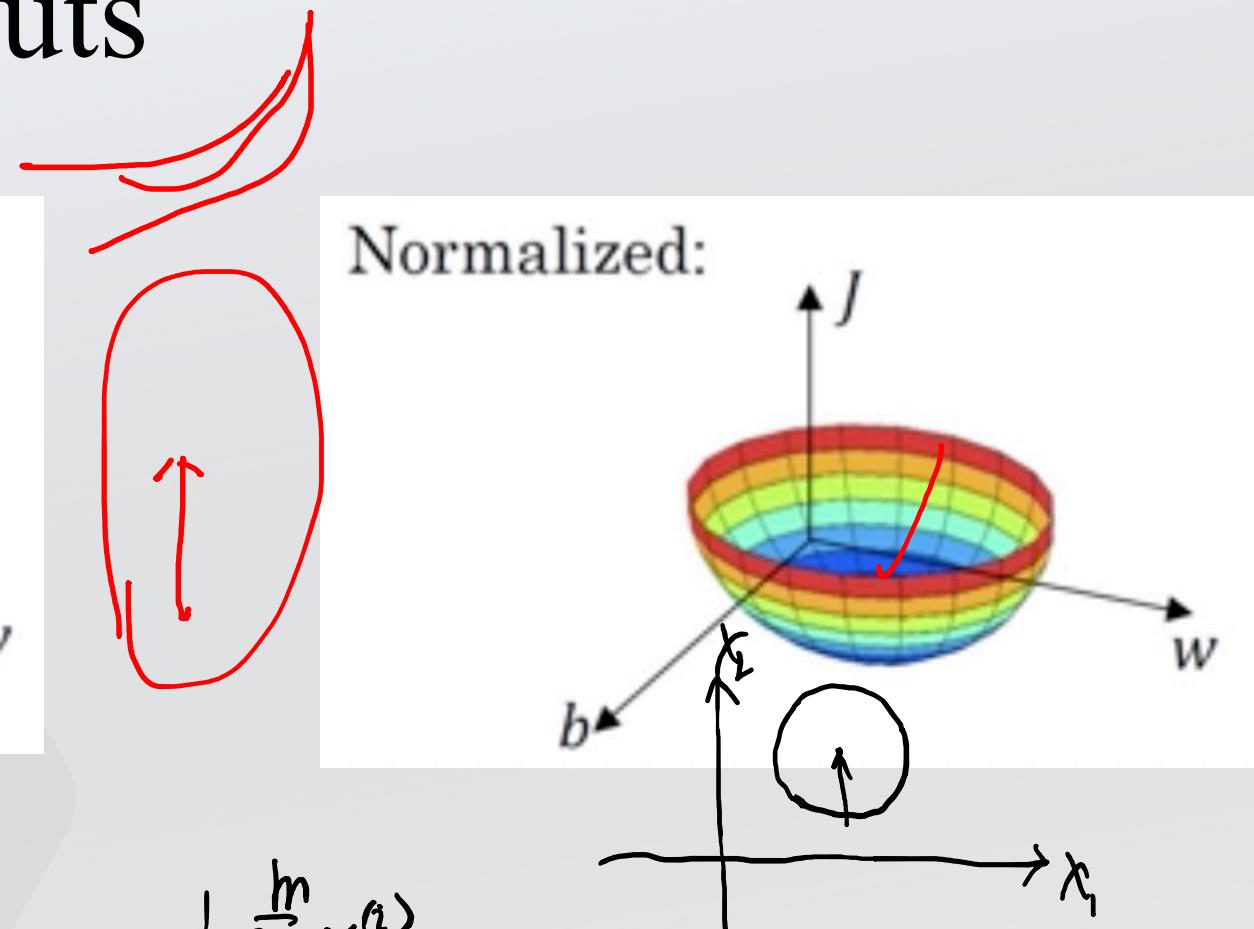


Step1: Subtract out or to zero the mean

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

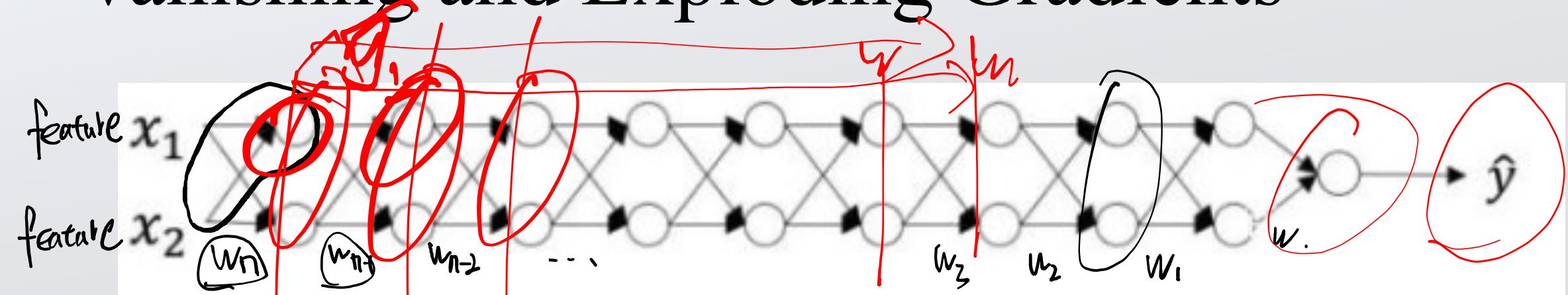
Step2: Normalize the variances

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$



$$x_i := \frac{x_i - \mu}{\sigma}$$

Vanishing and Exploding Gradients



active function $g(Q) = 2$ bias $b = 0$

$$\hat{y} = w^{(n)} \cdot w^{(n-1)} \cdot w^{(n-2)} \dots w^{(3)} \cdot w^{(2)} \cdot w^{(1)} \cdot x + b = [w] \cdot x$$

$$\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \dots \times \text{内存溢出}$$

$$\begin{bmatrix} 2.5 & 0 \\ 0 & 0.5 \end{bmatrix} \rightarrow 0$$

$$\frac{z^{(0)} = w^{(0)} \cdot x}{a^{(0)} = g(z^{(0)}) = z^{(1)}}$$

$$a^{(1)} = g(z^{(0)}) = g(w^{(1)} \cdot a^{(0)}) = w^{(2)} \cdot w^{(1)} \cdot x$$

Vanishing and Exploding Gradients

To prevent this problem:

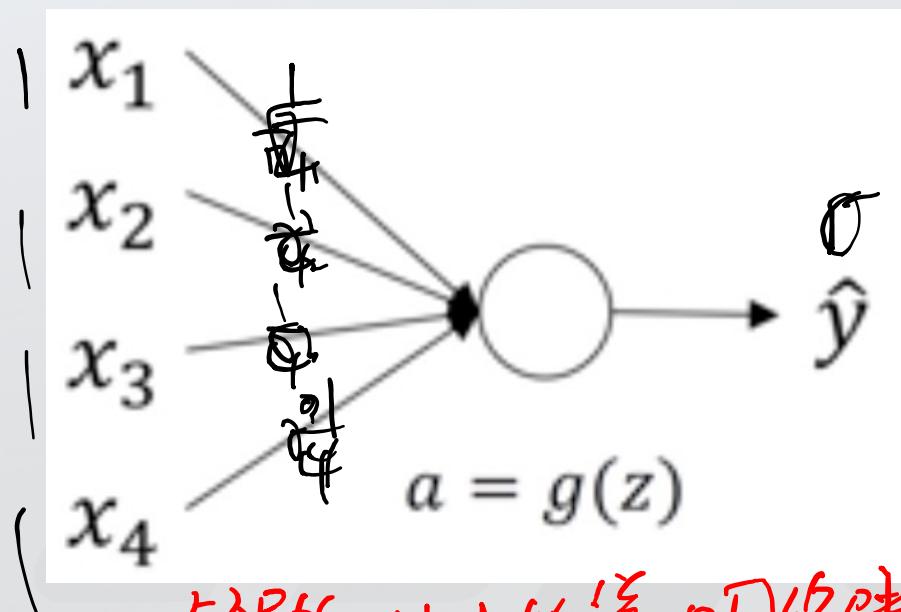
$$E[\cancel{a}^{[l-1]}] = E[a^{[l]}]$$

$$Var(a^{[l-1]}) = Var(a^{[l]})$$

$$\hat{y} = w \cdot x$$

$$a = L(\hat{y})$$

Weight initialization for deep networks



} 每个层的 output 的数量与网络层有关
 每个 \hat{y} 的数量与层与网络层有关
 各个层的 w 的精度与层数无关

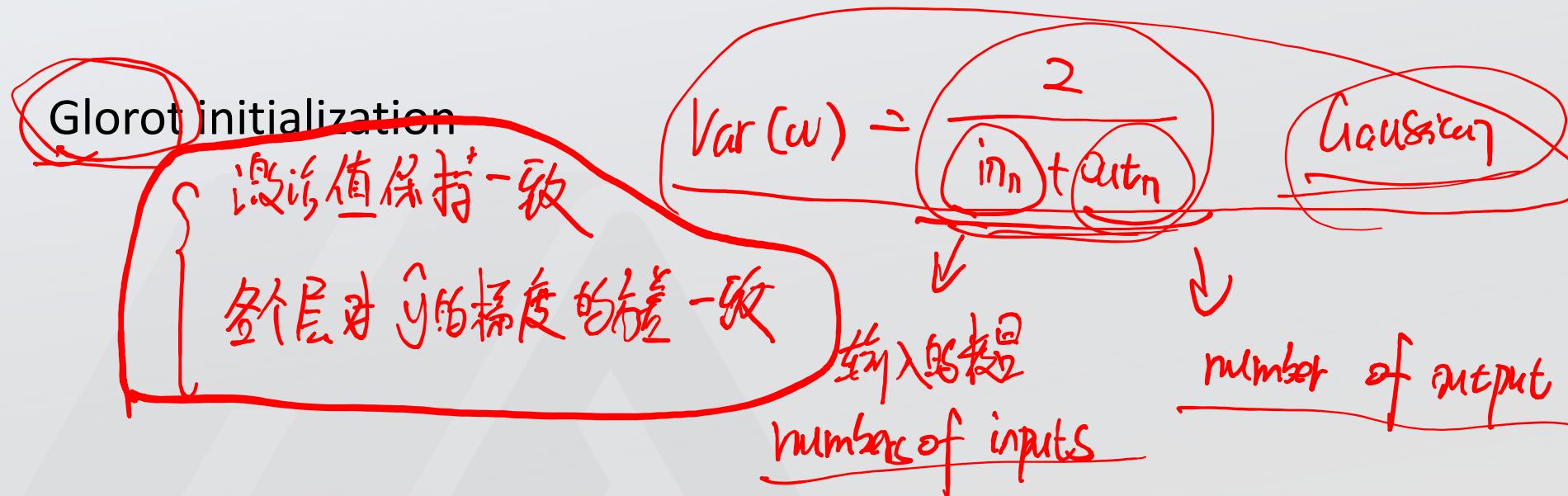
$$\hat{y} = \cancel{w_1}x_1 + \cancel{w_2}x_2 + \cancel{w_3}x_3 + \cancel{w_4}x_4 - \dots - \cancel{w_n}x_n.$$

$n \rightarrow$ $w \downarrow$

$$\frac{1}{4}x_1 + \frac{1}{4}x_1 + \frac{1}{4}x_1 + \frac{1}{4}x_1 - 1$$

$$\frac{1}{4}x_0 + \frac{1}{4}x_0 + \frac{1}{4}x_1 + \frac{1}{4}x_1 = 1$$

Weight initialization for deep networks



Keras: lecun_normal(seed=None), glorot_uniform(seed=None)

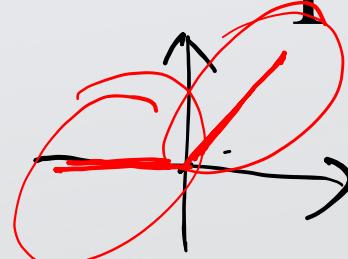
Link: <http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

Weight initialization for deep networks

He initialization

relu

$$\text{var}(w) = \frac{2}{n}$$



Keras: `he_normal(seed=None)` , `he_uniform(seed=None)`

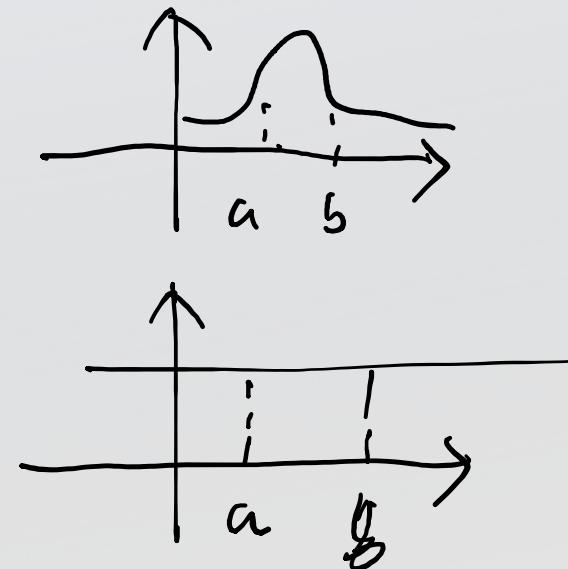
Link: <https://arxiv.org/pdf/1502.01852.pdf>

Weight initialization for deep networks

Lecun initialization

$$y = \alpha x \rightarrow$$

$$\text{Var}(w) = \frac{1}{n}$$



Keras: `lecun_normal(seed=None)`

`lecun_uniform(seed=None)`

Link: <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>



一所专注前沿互联网技术领域的创新实战大学