

经常让 vanishing gradients 发生的 **Activation func**: 1. **ReLU** saturates for negative values, giving zero gradient. 2. **Leaky ReLU**:  $\max(0.01x, x)$ . Does not saturate. Is not differentiable at  $x = 0$ . Does not approximate identity near the origin 3. **Tanh**, like the **sigmoid** unit, saturates when  $x$  very positive or negative. Is differentiable everywhere as it has no discontinuities. Approximates identity near the origin and no additional learning occurs. 4. **Sigmoid**. At extremes, the unit saturates and thus has zero gradient. This results in no learning with gradient descent 5. **Gaussian Error Linear Unit (GELU)**: it is popular in recent transformer structures like BERT. Saturates when  $x$  is very negative. Is differentiable everywhere as it has no discontinuities. Does not approximate identity near the origin. ^^^

$f(z) = \begin{cases} 1, & z \geq 0 \\ -1, & z < 0 \end{cases}$  这个 activation func has 0 gradients for all inputs (except at 0 where it's non-differentiable)意味着没 learning  
 $f(x) = x^3$

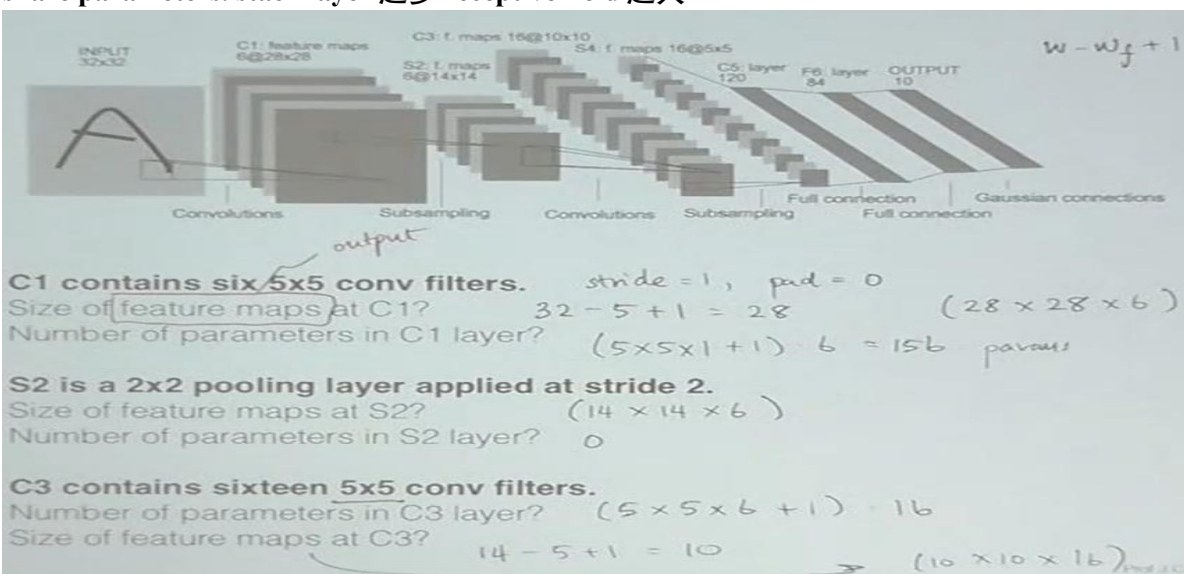
This activation func is nonlinear and differentiable everywhere, 这满足 good activation func 条件. 但很可能造成 exploding gradients input value 大, gradient 大. 而且, small inputs 会导致 vanishing gradients, e.g. inputs that are close to zero.

**CNN** 有 5 个  $3 \times 3$  filter, stride=1, padding = 1, 再  $2 \times 2$  max pooling with stride 2. Input image size  $256 \times 256 \times 3$ . First conv-pool layer 的 trainable parameters (including bias):  $(3 \times 3 (\text{filterSize}) \times 3 (\text{RGB}) + 1 (\text{biasTerm})) \times 5 (\text{filterNum}) = 140$  ^^^ CNN  $32 \times 32 \times 3$  第一层是 CNN, output =  $26 \times 26 \times 16$  filters number in the first convolutional = depth of the output volume (16),

Output size =  $\frac{\text{Input size} - \text{Filter size} + 2 \times \text{Padding}}{\text{Stride}} + 1$   $26 = \frac{32 - \text{Filter size} + 0}{1} + 1$  each filter has dimensions of  $7 \times 7 \times 3$ .

**Trainable parameters** 数量是  $32 \times (32 \times 32 \times 3)$  (第二层 filter 数)  $\times 3 \times 3 \times 16$  (filter). ^^^ input  $32 \times 32 \times 3$ , fully connected architecture 有 500 output neuron; convolutional layer 有 4 filter ( $4 \times 4$ ), 则 convolutional layer 的 output:  $29(32-4+1) \times 29 \times 4$ . FC 的 param:  $(32 \times 32 \times 3 + 1) \times 500$ , convolutional param  $(4 \times 4 \times 3 + 1) \times 4 (\text{filterNum})$

**Convolutional layers** 比 **fully connected layers** 有更少的 param 因为 **Convolutional layers have sparse connectivity and share parameters**. stack layer 越多 receptive field 越大



First layer:  $\mathbf{h}_1 = f(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{W}_s \mathbf{x}$

Second (Output) layer:  $\mathbf{z} = \mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2$

^^^ **forward propagation** 有第二个 convolutional layer (有 32 filters,  $3 \times 3$ ) depth of each filter in the second convolutional layer must match the depth of the input volume, 所以是 16 ^^^

where  $\mathbf{b}_1 \in \mathbb{R}^m$  and  $\mathbf{b}_2 \in \mathbb{R}^l$ .  $\mathbf{W}_s$  is the linear mapping for the "shortcut" connection  
 $f$  is the non-linear activation function like ReLU.  $\mathbf{W}_1 \in \mathbb{R}^{m \times n}$ ,  $\mathbf{W}_s \in \mathbb{R}^{m \times n}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{l \times m}$ . Total number of trainable parameters:  $= mn + m + mn + ml + 1 = 2mn + ml + m + 1$ . ^^^

**Backpropagation** 1. **Binary 分类**, 1wX 光没病 1k 有病数据集 (很不平衡). 用 **Data Augmentation** (能解决 image classification 的 class imbalance 问题) 解决 (对 train data/label 进行 Translation, **random Cropping** (flipping/mirroring 改 input pixel 但 label 不变能增加 train data point), **label smooth** (add noise to network, 把 label 概率 1 分给其他 label 一些

$$\mathcal{L}^{(i)} = \sum_{c=1}^C -y_c \log p_c \quad \mathcal{L}^{(i)} = y_c (1-\alpha) \cdot \log p_c + \alpha/(K-1)$$

$$\begin{aligned} \mathbf{z}_1 &= \mathbf{W}_1 \mathbf{x}^{(i)} + \mathbf{b}_1 \\ \mathbf{a}_1 &= \text{ReLU}(\mathbf{z}_1) \\ \mathbf{z}_2 &= \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2 \\ \hat{y}^{(i)} &= \sigma(\mathbf{z}_2) \end{aligned}$$

$$\mathcal{L}^{(i)} = \alpha \cdot y^{(i)} \cdot \log(\hat{y}^{(i)}) + \beta \cdot (1 - y^{(i)}) \cdot \log(1 - \hat{y}^{(i)})$$

$$J = -\frac{1}{m} \sum_{i=1}^m \mathcal{L}^{(i)}$$

$$\text{where } \hat{y}^{(i)} \in \mathbb{R}, y^{(i)} \in \mathbb{R}, \mathbf{x}^{(i)} \in \mathbb{R}^{D_x \times 1}, \mathbf{W}_1 \in \mathbb{R}^{D_{a1} \times D_x}, \mathbf{W}_2 \in \mathbb{R}^{1 \times D_{a1}}$$

$\mathbf{b}_1 \in \mathbb{R}^{D_{a1}}, \mathbf{b}_2 \in \mathbb{R}$ . Hyperparameter  $\alpha$  和  $\beta$  的用处: Weigh 每个 class 给 loss func 贡献多少 能帮 GC 因为 network 会 take larger steps 当从 underrepresented class instances 学习.  $\alpha$  取 1,  $\beta$  取 10, 因为 sample 比是 10, 这样每个 class 的 samples 才会 contribute

equally.  $\nabla_{\hat{y}^{(i)}} \mathcal{L}^{(i)} = \alpha \frac{y^{(i)}}{\hat{y}^{(i)}} - \beta \frac{1-y^{(i)}}{1-\hat{y}^{(i)}} = \delta_{\hat{y}^{(i)}}$  求  $\nabla_{b_2} \mathcal{L}^{(i)}$  and denote it as  $\delta_{b_2}$ . Hint:

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1-\sigma(z)) \quad \nabla_{z_2} \hat{y}^{(i)} = \sigma(z_2)[1-\sigma(z_2)]. \quad \text{backpropagate to } z_2: \nabla_{z_2} \mathcal{L}^{(i)} = \sigma(z_2)[1-\sigma(z_2)] \nabla_{\hat{y}^{(i)}} \mathcal{L}^{(i)}. \quad \text{因为是+号 (pass$$

$$\nabla_{b_2} \mathcal{L}^{(i)} = \sigma(z_2)[1-\sigma(z_2)] \nabla_{\hat{y}^{(i)}} \mathcal{L}^{(i)}$$

gradient) 所以

$$\delta_{b_2} = \sigma(z_2)[1-\sigma(z_2)] \delta_{\hat{y}^{(i)}}$$

$$\nabla_{\mathbf{W}_2} \mathcal{L}^{(i)} = \mathbf{a}_1^T \delta_{b_2} = \delta_{\mathbf{W}_2} \quad \nabla_{\mathbf{a}_1} \mathcal{L}^{(i)} = \mathbf{W}_2^T \nabla_{b_2} \mathcal{L}^{(i)} = \mathbf{W}_2^T \delta_{b_2} \quad \text{back 到 } z_1 \quad \nabla_{z_1} \mathcal{L}^{(i)} = \mathbb{I}(\mathbf{z}_1 > 0) \odot \nabla_{\mathbf{a}_1} \mathcal{L}^{(i)}.$$

因为  $\mathbf{W}_2$  是加号所以直接 pass  $\nabla_{b_1} \mathcal{L}^{(i)} = \nabla_{z_1} \mathcal{L}^{(i)} = \delta_{b_1}$ .  $\delta_{\mathbf{W}_1} = \delta_{b_1} \mathbf{x}^{(i)T}$  因为  $\mathbf{W}_1$  是  $\times$  号 (switch). Add L2 regu (len=1) to

$\hat{J} = J + \|\mathbf{W}_2\|_2^2$ . 用 vanilla GC (learning rate  $\epsilon$ ), 则 update rule:

$$\mathbf{W}_2^{(k+1)} = \mathbf{W}_2^{(k)} - \epsilon [\nabla_{\mathbf{W}_2} J^{(k)} + 2\mathbf{W}_2^{(k)}]$$

$$\mathbf{h}_1 = \text{ReLU}(\mathbf{W}_1 \mathbf{x})$$

$$\hat{\mathbf{h}}_1 = \text{ReLU}(\mathbf{W}_1 \hat{\mathbf{x}})$$

$$\mathbf{z} = \mathbf{W}_2 \mathbf{h}_1$$

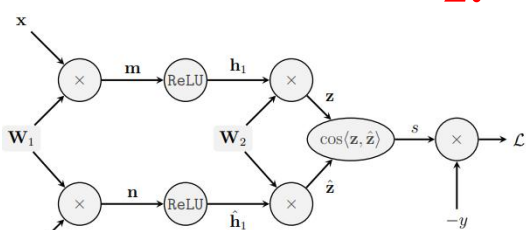
$$\hat{\mathbf{z}} = \mathbf{W}_2 \hat{\mathbf{h}}_1$$

$$s = \cos \langle \mathbf{z}, \hat{\mathbf{z}} \rangle = \frac{\mathbf{z}^T \hat{\mathbf{z}}}{\|\mathbf{z}\|_2 \|\hat{\mathbf{z}}\|_2}$$

$$\mathcal{L} = -y \cdot s$$

$$\mathbf{W}_2^{(k+1)} = \mathbf{W}_2^{(k)} - \epsilon [\nabla_{\mathbf{W}_2} J^{(k)} + 2\mathbf{W}_2^{(k)}]$$

2.



If the  $i^{th}$  pair of input  $(\mathbf{x}^{(i)}, \hat{\mathbf{x}}^{(i)})$  is composed of signature images both of which are genuine, then the label for the  $i^{th}$  example is +1 ( $y^{(i)} = +1$ ).

If the  $i^{th}$  pair of input  $(\mathbf{x}^{(i)}, \hat{\mathbf{x}}^{(i)})$  is composed of signature images both of which are forged, then the label for the  $i^{th}$  example is -1 ( $y^{(i)} = -1$ ).

If the  $i^{th}$  pair of input  $(\mathbf{x}^{(i)}, \hat{\mathbf{x}}^{(i)})$  is composed of signature images one of which is genuine and the other is forged, then the label for the  $i^{th}$  example is -1 ( $y^{(i)} = -1$ ).

$$\frac{d}{dz} \left( \frac{f(z)}{g(z)} \right) = \frac{f'(z)g(z) - g'(z)f(z)}{g(z)^2} \quad \nabla_{\mathbf{z}} \mathcal{L} = -y \frac{\partial s}{\partial \mathbf{z}} \quad \nabla_{\hat{\mathbf{z}}} \mathcal{L} = -y \frac{\partial s}{\partial \hat{\mathbf{z}}}$$

$$\text{for } f'(z) = \frac{df(z)}{dz} \text{ and } g'(z) = \frac{dg(z)}{dz}.$$

$$\frac{\partial s}{\partial \mathbf{z}} = \frac{(\|\mathbf{z}\|_2 \|\hat{\mathbf{z}}\|_2 \hat{\mathbf{z}} - (\mathbf{z}^T \hat{\mathbf{z}}) \frac{\|\hat{\mathbf{z}}\|_2}{\|\mathbf{z}\|_2} \mathbf{z})}{\|\mathbf{z}\|_2^2 \|\hat{\mathbf{z}}\|_2^2} \quad \delta_{\mathbf{z}} = -y \cdot \frac{(\|\mathbf{z}\|_2 \|\hat{\mathbf{z}}\|_2 \hat{\mathbf{z}} - (\mathbf{z}^T \hat{\mathbf{z}}) \frac{\|\hat{\mathbf{z}}\|_2}{\|\mathbf{z}\|_2} \mathbf{z})}{\|\mathbf{z}\|_2^2 \|\hat{\mathbf{z}}\|_2^2}$$

$$\frac{\partial s}{\partial \hat{\mathbf{z}}} = \frac{(\|\mathbf{z}\|_2 \|\hat{\mathbf{z}}\|_2 \mathbf{z} - (\mathbf{z}^T \hat{\mathbf{z}}) \frac{\|\mathbf{z}\|_2}{\|\hat{\mathbf{z}}\|_2} \hat{\mathbf{z}})}{\|\mathbf{z}\|_2^2 \|\hat{\mathbf{z}}\|_2^2} \quad \delta_{\hat{\mathbf{z}}} = -y \cdot \frac{(\|\mathbf{z}\|_2 \|\hat{\mathbf{z}}\|_2 \mathbf{z} - (\mathbf{z}^T \hat{\mathbf{z}}) \frac{\|\mathbf{z}\|_2}{\|\hat{\mathbf{z}}\|_2} \hat{\mathbf{z}})}{\|\mathbf{z}\|_2^2 \|\hat{\mathbf{z}}\|_2^2}$$

$$\begin{aligned} \nabla_{\mathbf{h}_1} \mathcal{L} &= \frac{\partial \mathbf{z}}{\partial \mathbf{h}_1} \frac{\partial \mathcal{L}}{\partial \mathbf{z}} & \nabla_{\mathbf{m}} \mathcal{L} &= \frac{\partial \mathbf{h}_1}{\partial \mathbf{m}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_1} & \delta_{\mathbf{m}} &= \mathbb{I}(\mathbf{m} \geq 0) \odot \delta_{\mathbf{h}_1} \\ \nabla_{\hat{\mathbf{h}}_1} \mathcal{L} &= \frac{\partial \hat{\mathbf{z}}}{\partial \hat{\mathbf{h}}_1} \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{z}}} & \frac{\partial \mathbf{z}}{\partial \mathbf{h}_1} &= \frac{\partial \hat{\mathbf{z}}}{\partial \hat{\mathbf{h}}_1} = \mathbf{W}_2^T & \delta_{\hat{\mathbf{h}}_1} &= \mathbf{W}_2^T \delta_{\hat{\mathbf{z}}} & \nabla_{\mathbf{n}} \mathcal{L} &= \frac{\partial \hat{\mathbf{h}}_1}{\partial \mathbf{n}} \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{h}}_1} & \delta_{\mathbf{n}} &= \mathbb{I}(\mathbf{n} \geq 0) \odot \delta_{\hat{\mathbf{h}}_1} & \nabla_{\mathbf{W}_1} \mathcal{L} &= \frac{\partial \mathbf{m}}{\partial \mathbf{W}_1} \frac{\partial \mathcal{L}}{\partial \mathbf{m}} + \frac{\partial \mathbf{n}}{\partial \mathbf{W}_1} \frac{\partial \mathcal{L}}{\partial \mathbf{n}} \end{aligned}$$

(9 points) In the similarity network architecture,  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  represents the embedding vectors for input signature images  $\mathbf{x}$  and  $\hat{\mathbf{x}}$  respectively. Suppose we are given a training sample,  $\{(\mathbf{x}^{(g)}, \hat{\mathbf{x}}^{(g)}), +1\}$ .

(b)

计算 train sample 的

loss 如果  $\mathbf{z}^{(g)} = \hat{\mathbf{z}}^{(g)}$ ,  $s^{(g)} = \cos\langle \mathbf{z}^{(g)}, \hat{\mathbf{z}}^{(g)} \rangle = 1$ ,  $\mathcal{L}^{(g)} = -1 \cdot s^{(g)} = -1$ . 计算 loss 如果  $\mathbf{z}^{(g)}$  and  $\hat{\mathbf{z}}^{(g)}$  are orthogonal to each other  $s^{(g)} = \cos\langle \mathbf{z}^{(g)}, \hat{\mathbf{z}}^{(g)} \rangle = 0$ ,  $\mathcal{L}^{(g)} = -1 \cdot s^{(g)} = 0$ . 计算 loss 如果  $\mathbf{z}^{(g)} = -\hat{\mathbf{z}}^{(g)}$ ,  $s^{(g)} = \cos\langle \mathbf{z}^{(g)}, \hat{\mathbf{z}}^{(g)} \rangle = -1$ ,  $\mathcal{L}^{(g)} = -1 \cdot s^{(g)} = 1$ . (c) 解释是否 loss function

让 embedding vectors 到 right direction: From part (b), 通过最小化 loss func, embedding vectors for the input pair of images that are genuine are getting forced to be similar to each other (loss for b(i) is the least) which is what we want to achieve for the

$$\mathcal{L} = \frac{1}{2} \|\mathbf{f} - \mathbf{y}\|^2$$

$$\mathbf{f} = \max(\alpha \cdot \mathbf{g}, \mathbf{g}), 0 < \alpha < 1$$

$$\mathbf{g} = \mathbf{h} + \mathbf{b}$$

$$\mathbf{h} = \mathbf{W}_k \mathbf{x}$$

signature forgery detection application. 3.

$$\frac{\partial \mathcal{L}}{\partial f_i} = \frac{\partial}{\partial f_i} \left[ \frac{1}{2} (\mathbf{f} - \mathbf{y})^T (\mathbf{f} - \mathbf{y}) \right]$$

$$= \frac{\partial}{\partial f_i} \left[ \frac{1}{2} (\mathbf{f}^T \mathbf{f} - 2\mathbf{y}^T \mathbf{f} + \mathbf{y}^T \mathbf{y}) \right]$$

$$= \frac{1}{2} (2f_i - 2y_i)$$

$$= f_i - y_i.$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{g}} = \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathcal{L}}{\partial \mathbf{f}}.$$

$$\frac{\partial f_j}{\partial g_i} = \frac{\partial}{\partial g_i} [\max(\alpha g_j, g_j)]$$

$$= \begin{cases} 0 & \text{if } i \neq j \\ \alpha & \text{if } i = j \text{ and } \alpha g_j > g_j \\ 1 & \text{if } i = j \text{ and } \alpha g_j \leq g_j \end{cases}$$

$$\frac{\partial g_j}{\partial h_i} = \frac{\partial}{\partial h_i} (h_j + b_j)$$

$$= \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{else} \end{cases}$$

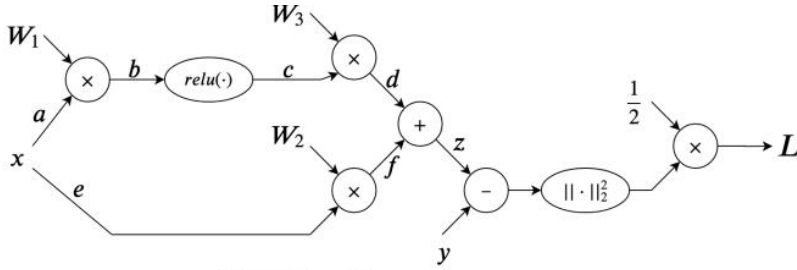
Therefore  $\frac{\partial \mathbf{g}}{\partial \mathbf{h}} = \mathbf{I}$ , and analogously  $\frac{\partial \mathbf{g}}{\partial \mathbf{b}} = \mathbf{I}$ . Therefore  $\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}} = \frac{\partial \mathcal{L}}{\partial \mathbf{g}}$ .

Finally we can use the provided formula to derive that  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_k} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}} \mathbf{x}^T$ .  $\mathbf{W}_1 \mathbf{x} = \mathbf{W}_1 (\mathbf{x} + \lambda \mathbf{z})$ . 在 train set 把 x 换成 x+λz, 再 train

new dataset. 哪些 param 会变? Fix a given input x, and consider its obfuscated variant  $\mathbf{x} + \lambda \mathbf{z}$ . Because the NN activations 在 linear transform 之后会一样, 所以所有哪些 gradient 都一样. The only 被影响的 gradients 是直接由 input 计算的. If the input is x, then the gradient  $\partial \mathcal{L} / \partial \mathbf{W}_1 = \partial \mathcal{L} / \partial \mathbf{h} \mathbf{x}^T$ ; if the input is  $\mathbf{x} + \lambda \mathbf{z}$ , then the gradient is  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}} \mathbf{x}^T$ . This gradient 会不同, so the rest of the network will train well. It's worth mentioning that this is only valid on the first step of training, 因为每次改变 w 都要重新

计算 z 4.

$$\mathbf{z} = \mathbf{W}_3 \text{relu}(\mathbf{W}_1 \mathbf{x}) + \mathbf{W}_2 \mathbf{x}. \quad \mathcal{L} = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{z} - \mathbf{y})^T (\mathbf{z} - \mathbf{y}), \text{ with } \mathbf{x} \in \mathbb{R}^n, \mathbf{y}, \mathbf{z} \in \mathbb{R}^m, \mathbf{W}_1 \in \mathbb{R}^{p \times n}, \mathbf{W}_2 \in \mathbb{R}^{m \times n}, \text{ and } \mathbf{W}_3 \in \mathbb{R}^{m \times p},$$



$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial \mathbf{x}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{a}} + \frac{\partial \mathcal{L}}{\partial \mathbf{e}} \\
 &= \frac{\partial \mathbf{b}}{\partial \mathbf{x}} \frac{\partial \mathcal{L}}{\partial \mathbf{b}} + \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \frac{\partial \mathcal{L}}{\partial \mathbf{f}} \\
 &= \frac{\partial \mathbf{b}}{\partial \mathbf{x}} \frac{\partial \mathcal{L}}{\partial \mathbf{b}} + \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \frac{\partial \mathbf{z}}{\partial \mathbf{f}} \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \\
 &= \mathbf{W}_1^T \left( \mathbb{1}(\mathbf{b}) \odot (\mathbf{W}_3^T \frac{\partial \mathcal{L}}{\partial \mathbf{z}}) \right) + \mathbf{W}_2^T \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \\
 &= \mathbf{W}_1^T \left( \mathbb{1}(\mathbf{W}_1 \mathbf{x}) \odot (\mathbf{W}_3^T \frac{\partial \mathcal{L}}{\partial \mathbf{z}}) \right) + \mathbf{W}_2^T \frac{\partial \mathcal{L}}{\partial \mathbf{z}}
 \end{aligned}$$

with

$$\begin{aligned}
 \frac{\partial \mathbf{z}}{\partial \mathbf{f}} &= \mathbf{I} \\
 \frac{\partial \mathbf{f}}{\partial \mathbf{z}} &= \mathbf{W}_2^T
 \end{aligned}$$

5.

$$\begin{aligned}
 \mathcal{L} &= \frac{1}{2} \|h_3(h_2(h_1(\mathbf{x}))) - \mathbf{y}\|^2 \\
 h_1(\mathbf{x}) &= \text{PReLU}(\mathbf{W}_1 \mathbf{x}) \\
 h_2(\mathbf{x}) &= \text{ELU}(\mathbf{W}_2 \mathbf{x}) \\
 h_3(\mathbf{x}) &= \mathbf{W}_3 \mathbf{x}
 \end{aligned}$$

- $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{n \times n}$  and  $\mathbf{W}_3 \in \mathbb{R}^{m \times n}$
- For the PReLU unit,  $f(x) = \max(\alpha x, x)$ .
- For the ELU unit,  $f(x) = \max(\alpha(e^x - 1), x)$ .

$$\begin{aligned}
 e &= \mathbf{W}_1 x \\
 d &= \text{ReLU}(e) \\
 c &= \mathbf{W}_2 d \\
 b &= \text{ELU}(c) \\
 a &= \mathbf{W}_3 b
 \end{aligned}$$

Therefore,  $y = \|a\|^2$ .

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial W_3} &= \frac{a}{\partial W_3} \cdot \frac{\partial \mathcal{L}}{\partial a} = (a - y) b^T \frac{\partial a}{\partial b} = W_3^T \\
 \frac{\partial \mathcal{L}}{\partial e} &= \text{diag}(\mathbb{1}\{c_1 < 0\} \cdot \alpha e^{c_1} + \mathbb{1}\{c_1 > 0\}, \dots, \mathbb{1}\{c_n < 0\} \cdot \alpha e^{c_n} + \mathbb{1}\{c_n > 0\}) \\
 \frac{\partial \mathcal{L}}{\partial d} &= \text{diag}(\mathbb{1}\{e_1 < 0\} \cdot \alpha + \mathbb{1}\{e_1 > 0\}, \dots, \mathbb{1}\{e_n < 0\} \cdot \alpha + \mathbb{1}\{e_n > 0\}) \\
 \frac{\partial \mathcal{L}}{\partial e} &= \frac{\partial e}{\partial W_1} \cdot \frac{\partial d}{\partial e} \cdot \frac{\partial c}{\partial d} \cdot \frac{\partial \mathcal{L}}{\partial c} \\
 &= \frac{\partial d}{\partial e} \cdot W_2^T \cdot \frac{\partial b}{\partial c} \cdot W_3^T \cdot (a - y) \cdot x^T
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial W_2} &= \frac{\partial c}{\partial W_2} \cdot \frac{\partial b}{\partial c} \cdot \frac{\partial a}{\partial b} \cdot \frac{\partial \mathcal{L}}{\partial a} \\
 &= \frac{\partial b}{\partial c} \cdot W_3^T \cdot (a - y) \cdot d^T
 \end{aligned}$$



$$h_p = W_1 x_p^{(i)} + b_1$$

$$z_1 = \text{ReLU}(h_p)$$

$$h_q = W_1 x_q^{(i)} + b_1$$

$$z_2 = \text{ReLU}(h_q)$$

$$z = z_1 - z_2$$

$$z_3 = W_2 z + b_2$$

$$\hat{y}^{(i)} = \sigma(z_3)$$

$$L^{(i)} = L_{CE}(y^{(i)}, \hat{y}^{(i)})$$

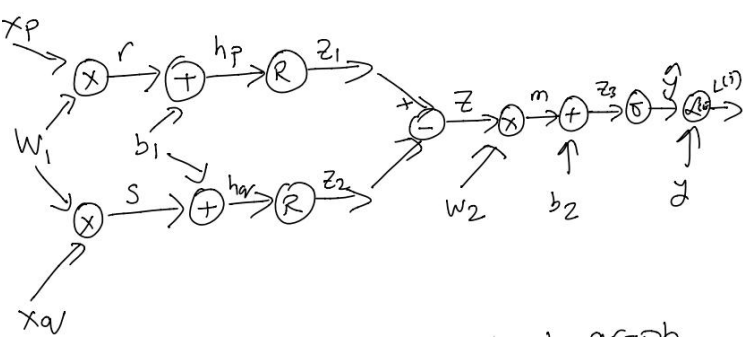
$$L = -\frac{1}{m} \sum_{i=1}^m L^{(i)}$$

Lce=Cross Entropy Loss 先写下 dimension

$$x_p^{(i)} \in \mathbb{R}^{D_x}, x_q^{(i)} \in \mathbb{R}^{D_x}, \hat{y}^{(i)} \in \mathbb{R}$$

$$W_1 \in \mathbb{R}^{D_z \times D_x}, b_1 \in \mathbb{R}^{D_z}$$

$$W_2 \in \mathbb{R}^{1 \times D_z}, b_2 \in \mathbb{R}$$



$$\hat{y}^{(i)} = \sigma(z_3)$$

$$\sigma(z_3) = \frac{1}{1 + e^{-z_3}} \quad \frac{d\hat{y}^{(i)}}{dz_3} = \sigma(z_3)[1 - \sigma(z_3)]$$

$$\frac{dL^{(i)}}{db_2} = \frac{dL}{dz_3} = \delta_{z_3} \quad \frac{dL^{(i)}}{dw_2} = \frac{dm}{dz} \frac{dL^{(i)}}{dm}$$

$$\delta_{b_2} = \delta_{z_3} \quad \frac{dm}{dz} \in \mathbb{R}^{1 \times D_z}$$

因为 w2 是 row vector, 所以结果是  $z^T \frac{dm}{dz} = z^T$

$$\frac{dL^{(i)}}{dw_2} = \delta_{z_3} z^T \frac{dL^{(i)}}{dz} = \frac{dm}{dz} \frac{dL^{(i)}}{dm} \quad \therefore \delta_z = \delta_{z_3} w_2^T \delta_{z_1} = \delta_z$$

$$\frac{dL^{(i)}}{dh_q} = \frac{dz_2}{dh_q} \frac{dL^{(i)}}{dz_2}$$

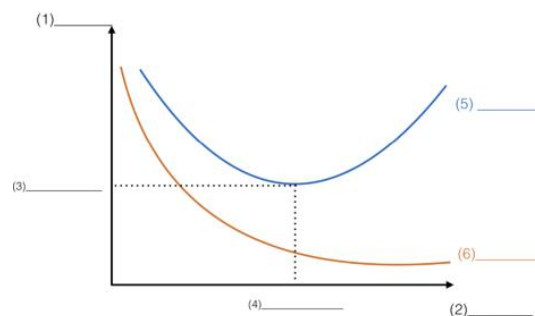
$$\frac{dL^{(i)}}{dr} = \delta_{h_p}$$

$$\delta_{z_2} = -\delta_z = \mathbb{I}(h_q > 0) \odot \delta_{z_2} \quad \delta_{b_1} = \delta_{h_p} + \delta_{h_q} \frac{dL^{(i)}}{ds} = \delta_{h_q}$$

$$\left. \frac{dL^{(i)}}{dw_1} \right|_{\text{path 1}} = \frac{dr}{dw_1} \frac{dL^{(i)}}{dr} = \frac{dL^{(i)}}{dr} x_p^T \quad \left. \frac{dL^{(i)}}{dw_1} \right|_{\text{path 2}} = \frac{ds}{dw_1} \frac{dL^{(i)}}{ds} = \frac{dL^{(i)}}{dh_q} x_q^T$$

$$\delta_{w_1} = \delta_{h_p} x_p^T + \delta_{h_q} x_q^T$$

**Regularization: 特征:** L1 regularization 学到的 weight 比 L2 的更 sparse. + L1 经常导致 weights 变 0 + L1 imposes a penalty = sum of model weights 的绝对值. This encourages the optimization of certain feature weights to 0, promoting simpler and sparser patterns in the model. **By pushing less important feature weights to zero, L1 aids in preventing overfitting.** Additionally, examining the trained weights allows for **feature selection**, simplifying the model and conserving computational resources. + **Dropout** acts as regularization(防止 overfitting due to a layer's "over-reliance" on a few of its inputs. 因为

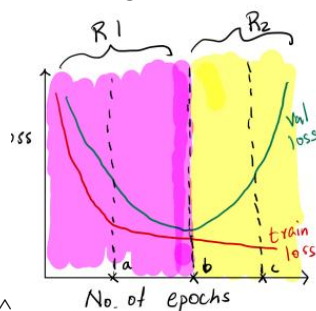


这些 input 不总是 present during training (i.e. they are dropped at random), the layer learns to 用所有 input, improving generalization.实现: Dropout applies a random mask to activations during training to simulate training a large ensemble of models. Each ensemble comprises units that passed through a mask, and at test time, predictions are averaged across these subnetworks. This approach approximates bagging, where multiple models are trained on different subsets of data. However, in dropout, all models share a significant portion of weights, distinguishing it slightly from traditional bagging.) **Implement:** inverted drop.train 的时候产生 mask\*values after affine transform/p(prob of keep), test 的时候 nothing(用 hidden layer 所有的 unit) + 失败的 regularization attempts (such as having too large a weight on a parameter norm penalty) 会导致 model underfitting + 用太多 regu 会 underfit train set 会导致 worse performance + Helpful for generalization/ test set performance 而不是 training set + optimizer 用 Adagrad 会有 equal or smaller learning rate in successive iterations.(因为 gradient history stays the same or increases in each dimension for every iteration of Adagrad.) + (不对:) Adding a regularization penalty 总是减少 training loss. ^^why **Ensemble** 一些 NN 会增加 generalization Perform(If NN make independent errors->they 可能 to make errors on different trials. We only get trial wrong when several NN make 一样 errors. So, model 之间 correlation 更低, ensemble model accuracy 更好)^^ **Transfer Learning** 可以 freeze most parameters of the original network^^(不对)**Multitask learning** is not applicable 如果一个 task 的数据少^^**Early stopping** 是一种 regu 方法, 在每个 iteration evaluate train 和 validation error loss 来返回有最低 validation error 的 model(1) Loss(2) Number of iterations (3) best validation error (4) Optimal stopping point (5) Validation loss (6) Training loss^^^During training, training error 减少 validation error 先减再加. After training completes, the testing error 比 train 高很多, 因为 **overfitting**, 可能的原因: The model is overtrained. The training data is not enough. The training data is enough but highly correlated. The model is too complex for the dataset. The data is too noisy.

NN 有高 train 准确但 validation 准确低, 则用 dropout 或 dataset augmentation 用在 input 层, 增加 layer 数没用

**Training, Validation, and Testing.** Training set 用于 adjust parameters to fit the input. Validation set 用于选 optimal hyperparameters. Testing set 用于作为 proxy to see how well the model (the architecture and hyperparameters) would generalize to new data. The test set 只能用 1 次, after the cross-validation procedure, where each fold will be used k times for training/validation. 增加 trials 数 can allow more complicated model to be learned. 也能 against overfitting and make the model more robust.

**Validation accuracy 高因为:** The group was effectively training on the validation set, since the subsampled trials were highly correlated, and hence the parameters were adjusted to fit the highly correlated trials in the training set, which were effectively in the validation set as well. They should have first split the training and validation set, and then subsampled. 多次在 train set 测试会导



R1 → underfitting  
R2 → overfitting  
b → stopping point

致 overfit to the testing set, 不会产生 new data ^^  
error keep go down, 但 validation error go up, 在 b 停止因为 lowest validation

**L-∞ regularization**

overfitting: training

$$\|x\|_{\infty} = \max_i |x_i|$$

$$LSE(x) = \ln \left( \sum_{i=1}^n e^{x_i} \right)$$

Show that the following inequality holds for  $n \geq 1$   $\|x\|_{\infty} \leq LSE(x) \leq \|x\|_{\infty} + \ln(n)$  (4)

suppose  $P = \max_i \{ |x_1|, |x_2|, \dots, |x_n| \}$  lower bound:  $e^P \leq \sum_{i=1}^n e^{x_i}$  Upper

$$\sum_{i=1}^n e^{x_i} \leq n e^P$$

bound  $\sigma$  Combine: taking natural algorithm:

$$\ln(e^P) \leq \ln \left( \sum_{i=1}^n e^{x_i} \right) \leq \ln(n e^P)$$

$$\Rightarrow P \ln(e) \leq \ln \left( \sum_{i=1}^n e^{x_i} \right) \leq \ln(n) + P \ln(e)$$

$$\|x\|_{\infty} \leq LSE(x) \leq \|x\|_{\infty} + \ln(n)$$

Is the lower bound in (4) strict for  $n > 1$ ?

Clearly for  $n > 1$ ,

$$\ln(e^P) < \ln \left( \sum_{i=1}^n e^{x_i} \right)$$

$$e^P < \sum_{i=1}^n e^{x_i}$$

Taking natural algorithm

$$P \ln(e) < \ln \left( \sum_{i=1}^n e^{x_i} \right)$$

$$\|x\|_{\infty} < LSE(x)$$

Under what

condition on  $x$ , will the upper bound in (4) be satisfied with equality?

for all  $i, j \in n$

Suppose  $|x_i| = |x_j|$

$$\sum_{i=1}^n e^{x_i} = \sum_{i=1}^n e^P$$

$$= n e^P$$

Hence, the upper bound will be satisfied with equality

$$LSE(x) = \|x\|_{\infty} + \ln(n)$$

(d) Use

$$\|x\|_{\infty} \leq \frac{1}{t} LSE(tx) \leq \|x\|_{\infty} + \frac{\ln(n)}{t}$$

for some scaling constant  $t >$

the result from (4) to show that the following inequality holds,

Let  $t > 0$  be some scaling

constant. Substituting  $x$  with

$tx$  in (2), we get

$$0 \|tx\|_{\infty} \leq LSE(tx) \leq \|tx\|_{\infty} + \ln(n)$$

$$t \|x\|_{\infty} \leq LSE(tx) \leq t \|x\|_{\infty} + \ln(n)$$

Dividing by  $t$ , we get the required inequality,

Now,

$$\|tx\|_{\infty} = t \|x\|_{\infty} \quad \|x\|_{\infty} \leq \frac{1}{t} LSE(tx) \leq \|x\|_{\infty} + \frac{\ln(n)}{t}$$

$$\mathcal{L}_{reg}(\theta) = \mathcal{L}(\theta) + R(W)$$

on  $W_i$  (也叫 Chebyshev norm)

$R(W) = \frac{\lambda}{n} \sum_{i=1}^n \|W_i\|_{\infty}$ .  $\|A\|_{\infty} = \max_{1 \leq k \leq p} \sum_{j=1}^m |A_{kj}|$  For  $A \in \mathbb{R}^{p \times m}$ , the  $L - \infty$  norm returns the maximum row sum of  $A$ . Mathematically,

Assume gradient of  $|w|$  with respect to  $w$  is  $\text{sign}(w)$ . 算  $R(w)$  对  $w$  求导

Based on the dimensions of input, hidden layer and output dimensions, the individual weight matrix dimensions will be known to us. Without loss of generality, let the weight matrix of the  $i^{th}$  hidden layer  $W_i \in \mathbb{R}^{m \times p}$  comprising of row vectors  $w_r^i$  and elements represented as  $w_{rc}^i$  then

$$R(W) = \frac{\lambda}{n} \sum_{i=1}^n \left( \max_{1 \leq r \leq m} \sum_{c=1}^p |w_{rc}^i| \right)$$

$j = \text{index of row}$

$$R(\mathbf{W}) = \frac{\lambda}{n} \sum_{i=1}^n \sum_{c=1}^p |w_{jc}^i| \quad \nabla_{\mathbf{W}_i} R(\mathbf{W}) = \frac{\lambda}{n} \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \text{sign}(w_{j1}^i) & \text{sign}(w_{j2}^i) & \dots & \text{sign}(w_{jp}^i) \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

with maximum row-sum

$$= \frac{\lambda}{n} \begin{bmatrix} 0 \\ 0 \\ \text{sign}(\mathbf{w}_j^i) \\ 0 \end{bmatrix}$$

^^^ 加上 L2 的 updated loss func:  $\hat{J} = J + \|\mathbf{W}_2\|_2^2$ . 用了 vanilla gradient descent with learning rate  $\epsilon$  后 update

$$\mathbf{W}_2^{(k+1)} = \mathbf{W}_2^{(k)} - \epsilon [\nabla_{\mathbf{W}_2} J^{(k)} + 2\mathbf{W}_2^{(k)}]$$

role

$$\Omega(\theta) = \|\mathbf{W}\|_1 = \sum_{i,j} |W_{i,j}|$$

**Shrink size of model (加 L1--impose a sparsity constraint on the model. some parameters 有 optimal value 0)** .Loss

func:  $\tilde{J}(\mathbf{W}) = \alpha \|\mathbf{W}\|_1 + J(\mathbf{W})$  encourage a subset of the **weights to become zero**. During training, **dropout** make the model seem smaller 因为一些 activations are zeroed out. 但 多个 training epochs 之后, each unit is likely included in at least one subnetwork, making all units integral to the model. 所以测试时候 a model with dropout still 仍有一样数量的 parameters as a vanilla network. Additionally, weights corresponding to dropped units are not zeroed out but are skipped during backpropagation. Scaling activations during test time does not reduce model size since the number of nonzero weights remains the same. **L2 regularization** shrinks the magnitude of the weights, but the overall size of the model stays the same. Specifically, L2 does not cause parameters to

Mathematically show that  $\ell_2$  regularization shrinks the weight in gradient descent.

Hint: start with  $\tilde{\mathcal{L}}(\theta; \mathbf{X}, \mathbf{y}) = \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \frac{\alpha}{2} \|\theta\|_2^2$  and derive the gradient descent step for  $\theta$ .

be sparse.

$$(c) \quad \tilde{\mathcal{L}}(\theta; \mathbf{X}, \mathbf{y}) = \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \frac{\alpha}{2} \|\theta\|_2^2$$

$\theta$  is a vector

$$\nabla_{\theta} \tilde{\mathcal{L}}(\theta; \mathbf{X}, \mathbf{y}) = \frac{\partial}{\partial \theta} \tilde{\mathcal{L}}(\theta; \mathbf{X}, \mathbf{y})$$

$$\frac{\partial}{\partial \theta} \|\theta\|_2^2 = \frac{\partial}{\partial \theta} \theta^T \theta = 2\theta$$

$$= \frac{\partial}{\partial \theta} \left[ \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \frac{\alpha}{2} \|\theta\|_2^2 \right]$$

$$= \nabla_{\theta} \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \alpha \theta$$

$$\theta \leftarrow \theta - \epsilon \nabla_{\theta} \tilde{\mathcal{L}}(\theta; \mathbf{X}, \mathbf{y})$$

$$\leftarrow \theta - \epsilon \left[ \nabla_{\theta} \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \alpha \theta \right]$$

$$\leftarrow \underbrace{(1 - \epsilon\alpha)}_{\text{weight decay}} \theta - \epsilon \nabla_{\theta} \mathcal{L}(\theta; \mathbf{X}, \mathbf{y})$$

weight decay

^^^ Why can L2 regularization “weight decay” 有 L2 的 loss function

$$\Omega(\theta) = \frac{1}{2} \|\mathbf{W}\|_2^2$$

$$= \frac{1}{2} \sum_{i,j} W_{i,j}^2$$

$$\tilde{J}(\theta) = \frac{\alpha}{2} \mathbf{W}^T \mathbf{W} + J(\theta) \quad \mathbf{W} \leftarrow \mathbf{W} - \epsilon (\alpha \mathbf{W} + \nabla_{\mathbf{W}} J(\theta))$$

weights 有 decay tendency  $\nabla_{\mathbf{W}} \tilde{J}(\theta) = \alpha \mathbf{W} + \nabla_{\mathbf{W}} J(\theta) \quad \leftarrow (1 - \epsilon\alpha) \mathbf{W} - \epsilon \nabla_{\mathbf{W}} J(\theta)$  By modifying the gradient update with a weight decay term (1-..), for nonzero .., the weight matrix shrinks by a constant factor on each training iteration, before applying the gradient update. Across training, this will lead to weight matrices with smaller, more diffuse entries than vanilla gradient update



**Batch vs Layer Normalization:** batch normal 比 random weight initial 让网络更 robust, 因为 Random w i 有 exploding/vanishing gradients + exacerbate internal covariate shift 问题. 对于一个合适的 batch size, random w i 不影响 Batch Normalization as strongly 因为 Batch Normal 的 output 的 mean 和 variance 被限制和被模型学习了///

**Layer 和 Batch Normal 不同点:** batch 的 normalization occurs on a per neuron per batch that is the mean and variance calculated for all instances in a batch per channel. 在 layer N 里 mean 和 variance 被计算 for a single instance across all channels that is across a layer.///用 layer 不用 batch 的情况: As layer N normalizes the activations across an entire layer, 它比 Batch N 表现更好当 batch size 被减少. 它表现 also 更好当 modelling sequential data. Layer N is extensively used in RNNs.///**Batch 特征:** Batch N introduces noise to a hidden layer's activation, 因为 mean 和 standard deviation are estimated with a mini-batch of data. + 加速 train by requiring fewer iterations to converge to a given loss value + 有 learnable parameter + is a linear transformation + 用 running mean and variance from training statistics ^^ **Batch N 把每个 artificial neuron normalize, 所以 mean=0 variance=1, 这有用** 因为当 gradient update for a given weight, matrix assume 其他不变. 这不对! So, gradient steps in earlier layers may dramatically change the statistics of the neurons for the next layer. Batch N normalizes these statistics, so that they aren't dramatically different after each gradient step. Batch N prevents the outputs of each layer from being 太大/太小. ^^ **Batch N 怎么帮助 NN more robust to initialization:** (BN) normalizes the outputs at each layer, preventing activations from decaying to zero or exploding to infinity when cascading through multiple layers. This robustness to vanishing or exploding gradients during backpropagation is crucial for effective training. Vanishing gradients lead to information loss, making it difficult to propagate gradients upstream, while exploding gradients cause training instability or overflow errors due to excessively large steps in the loss

$$\text{mean} \rightarrow \mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

$$\text{variance} \rightarrow \sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$

landscape. **实现:** In Training, 会有  $m(\text{batchSize}) * n(\text{Num Unit})$ , 步骤

$$\hat{x}_i^{(j)} = \frac{x_i^{(j)} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} \quad y_i^{(j)} = \underbrace{\sigma_i \hat{x}_i^{(j)}}_{\text{trainable}} + \beta_i$$

两个 trainable 参数 (scaling 和 shifting). 主要是算 minibatch 的 mean 和 variance, 再 BatchNorm normalizes the layer's input based on the mean and variance 得到 average mean 和 variance. 用 momentum 更新

$$\bar{\mu} \leftarrow \bar{\mu} \cdot \text{momentum} + (1 - \text{momentum}) \mu$$

$$\bar{\sigma}^2 \leftarrow \bar{\sigma}^2 \cdot \text{momentum} + (1 - \text{momentum}) \sigma^2$$

running average mean 和 variance, 用两个参数 (learnable during training) scaling 和 shifting 决定 optimal distribution. 也会 track avg mean 和 variance. **In Testing**, 用 avgs 去 normalize test data, 再把  $\gamma$  和  $\beta$  apply 到 normalized data.

$$v_t = \alpha v_{t-1} - \epsilon \nabla_{\theta} L(\theta_{t-1} + \alpha v_{t-1})$$

**optimization:** Nesterov Momentum, update rule  $\theta_t = \theta_{t-1} + v_t$  (1)

$$v_{\text{new}} = \alpha v_{\text{old}} - \epsilon \nabla_{\tilde{\theta}_{\text{old}}} L(\tilde{\theta}_{\text{old}})$$

$$\tilde{\theta}_{\text{new}} = \tilde{\theta}_{\text{old}} + v_{\text{new}} + \alpha(v_{\text{new}} - v_{\text{old}}) \quad (2) \text{ one advantage of using the update rule in (2) over the update rule in (1)}$$

From ①, we have:

$$v_{\text{new}} = \alpha v_{\text{old}} - \epsilon \nabla_{\theta} L(\theta_{\text{old}} + \alpha v_{\text{old}})$$

⇒ let us assume that new parameter  $\tilde{\theta}$

$$\tilde{\theta} = \theta + \alpha v$$

$$\text{This is: } \left. \begin{array}{l} \tilde{\theta}_{\text{old}} = \theta_{\text{old}} + \alpha v_{\text{old}} \\ \tilde{\theta}_{\text{new}} = \theta_{\text{new}} + \alpha v_{\text{new}} \end{array} \right\} \begin{array}{l} \theta_{\text{old}} = \tilde{\theta}_{\text{old}} - \alpha v_{\text{old}} \\ \theta_{\text{new}} = \tilde{\theta}_{\text{new}} - \alpha v_{\text{new}} \end{array}$$

$$\begin{aligned} \nabla_{\tilde{\theta}} L(\tilde{\theta}_{\text{old}}) &= \frac{\partial L(\tilde{\theta}_{\text{old}})}{\partial \tilde{\theta}} = \frac{\partial \theta}{\partial \tilde{\theta}} * \frac{\partial L(\tilde{\theta}_{\text{old}})}{\partial \theta} \\ &= \frac{\partial \theta}{\partial \tilde{\theta}} * \frac{\partial L(\theta_{\text{old}} + \alpha v_{\text{old}})}{\partial \theta} = \nabla_{\theta} L(\theta_{\text{old}} + \alpha v_{\text{old}}) \\ \left[ \because \frac{\partial \theta}{\partial \tilde{\theta}} &= \frac{\partial}{\partial \tilde{\theta}} (\tilde{\theta} + \alpha v) = \mathbb{I} + 0 = \mathbb{I} \right] \end{aligned}$$

From ①, we have:  $\theta_{new} = \theta_{old} + v_{new}$

So, we have,  $\nabla_{\theta} L(\tilde{\theta}_{old}) = \nabla_{\theta} L(\theta_{old} + \alpha v_{old})$   $\tilde{\theta}_{old} - \alpha v_{new} = \theta_{old} - \alpha v_{old} + v_{new}$

$v_{new} = \alpha v_{old} - \epsilon \nabla_{\theta} L(\tilde{\theta}_{old}) \rightarrow$  This is ②  $\tilde{\theta}_{old} = \theta_{old} + v_{new} + \alpha(v_{new} - v_{old}) \rightarrow$  This is ② equation ① by a change in variable.

This representation helps us in implementation of nesterov momentum as this doesn't require us to calculate gradient at a different value of  $\theta + \alpha v$ .

Also, both  $\theta$  and  $\tilde{\theta}$  start from the same value of parameter initialization, as  $\theta = \tilde{\theta}$  initially as  $v=0$  at start

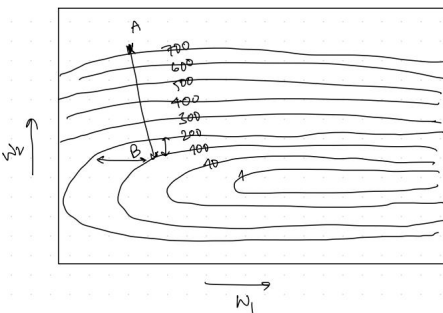
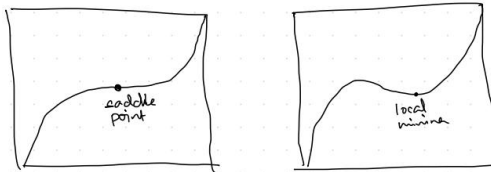
**L1, L2 saddle**

**point**  $L_1(x)$  has a saddle point. The curve decreases on one side of the saddle and increases on other side of saddle point. The gradient at a saddle point is equal to zero.

$L_2(x)$  has a poor local minima. This is a local minima as it is a minimum value of the function in it's surrounding until a certain limit in all the direction. The gradient is zero at a local minima.

We need momentum to get out of saddle point or local minima because the gradient becomes zero at these points.

It is also more likely to escape the saddle point than the local minima because after crossing the local minima, the momentum decreases due to gradient being in opposite direction to the momentum accumulated down the slope.

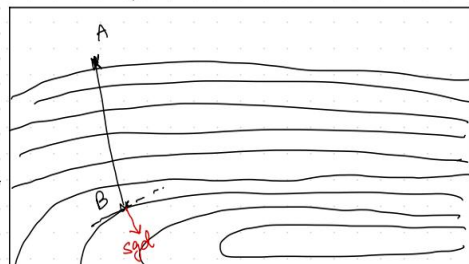


(1) The gradient along  $w_2$  direction is higher than the gradient along  $w_1$  direction.

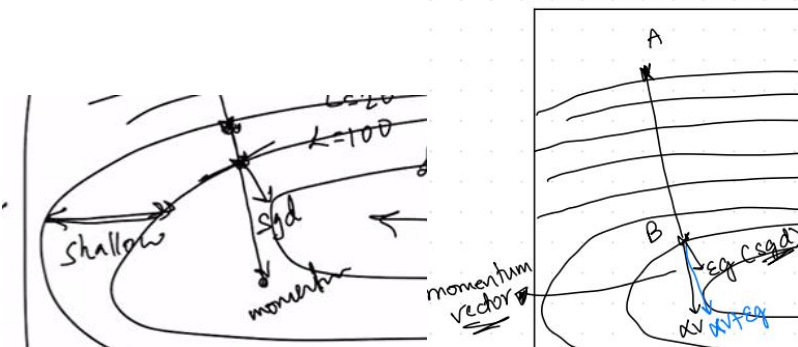
This is because the contour lines along  $w_2$  directions are very close to each other which indicates a steep curve in  $w_2$  direction.

The contour lines along  $w_1$  direction are further apart and hence has slower descent. (or low gradient)

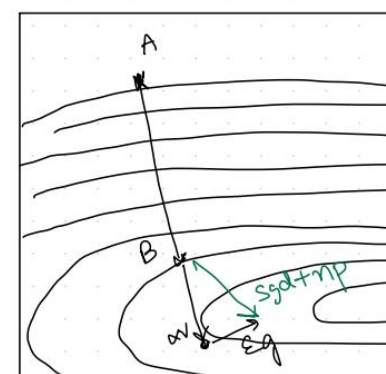
Vanilla gradient descent is in the direction perpendicular to a contour line



(3) SGD + momentum



sgd + nesterov momentum.

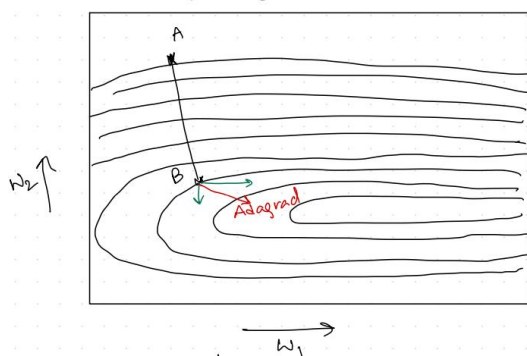


In this, first we take a step along the momentum and then calculate the gradient at that point in graph.

If gradient accumulated is less in a direction, then the step  $\frac{\epsilon}{\sqrt{a+v}}$  is more in that direction.

From plot, we can see that when it traversed from A to B, the gradient accumulated along  $w_2$  direction is more and gradient along  $w_1$  direction is lesser than it.

Hence, the step along  $w_1$  direction is more and step along  $w_2$  direction is less.



Also, in this particular case we can see that Adagrad is converging faster to minima without much zigzagging.

In the Gradient descent + momentum scheme, find a general expression of  $v_t$  in terms of gradients  $g_1, g_2, \dots, g_t$  and  $\epsilon$  (learning rate), considering an initial value of momentum  $v_0 = 0$ .

sgd + momentum update:

$$v_t = \alpha v_{t-1} - \epsilon g_t$$

$$\theta_t = \theta_{t-1} + v_t$$

$\rightarrow v_0 = 0$  [Given]

$$v_1 = \alpha v_0 - \epsilon g_1 = -\epsilon g_1$$

$$v_2 = \alpha v_1 - \epsilon g_2 = \alpha(-\epsilon g_1) - \epsilon g_2 = -\epsilon(g_2 + \alpha g_1)$$

$$v_3 = \alpha v_2 - \epsilon g_3 = \alpha(-\epsilon(g_2 + \alpha g_1)) - \epsilon g_3 = -\epsilon(g_3 + \alpha g_2 + \alpha^2 g_1)$$

From this, we can observe that

$$\Rightarrow v_t = -\epsilon(g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \dots + \alpha^{t-1} g_1)$$

$$\begin{aligned} E(v_t) &= E[-\epsilon(g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \dots + \alpha^{t-1} g_1)] \\ &= -\epsilon[E(g_t) + \alpha E(g_{t-1}) + \dots + \alpha^{t-1} E(g_1)] \end{aligned}$$

$$\begin{aligned}\theta_0 &= \theta_0 \\ \theta_1 &= \theta_0 + v_1 \\ \theta_2 &= \theta_1 + v_2 = \theta_0 + v_1 + v_2 \\ \theta_3 &= \theta_2 + v_3 = \theta_0 + v_1 + v_2 + v_3 \\ &\vdots \\ \text{So, for general } \theta_t &= \theta_0 + v_1 + v_2 + \dots + v_t\end{aligned}$$

Given,  $g_1, g_2, \dots, g_t$  are i.i.d r.v.s  
with mean  $\mu$  and var  $\sigma^2$ .  
From part (d), we have

$$E(\theta_3) = E(\theta) - \underbrace{EU[1 + (1+\alpha) + (1+\alpha+\alpha^2)]}$$

$$\begin{aligned} E(\theta_3) &= E(\theta_0) + \sum_{t=1}^3 \left[ -\varepsilon \mu \frac{[1 - \alpha^t]}{[1 - \alpha]} \right] \\ &= E(\theta_0) - \frac{\varepsilon \mu}{1 - \alpha} \left[ \sum_{t=1}^3 (1 - \alpha^t) \right] \\ &= E(\theta_0) - \frac{\varepsilon \mu}{1 - \alpha} \left[ 3 - \sum_{t=1}^3 \alpha^t \right] \\ &= E(\theta_0) - \frac{\varepsilon \mu}{1 - \alpha} \left[ 3 - \left( \frac{1 - \alpha^3}{1 - \alpha} \right) \right] \end{aligned}$$

**Classification vs Regression task:** Classification 用于当 target 是 categorical, regression 被用当 target 是 continuous. 他们俩都是 supervised machine learning algorithms. **Classification 例子:** 预测 yes or no + gender+Type of color. **Regression 问题例子:** 估计产品的 sale 和 price+预测球队分数^^

**Stochastic gradient descent (SGD) VS full-batch GD (GD):** 两个算法都是找 parameter that 最小化 loss func by evaluating parameters against data and then making adjustments **SGD 特征 vs Batch:** SGD 算的 gradient 会 noisier 比 Batch Gradient Descent 算+ SGD avoids trap of poor local minima(因为就 1 个 data point),converge to flat local minimum + SGD take more steps 但 converge faster + Batch GC 要 load entire data(mini-batch 用 some data, SG 用 1 data),所以 computationally more expensive



^^^用 SGDwith 正确算法,但 loss 一直不变,说明 learning rate 太小要增大^^^

**Initial small weight:** weights 小, activation magnitude 下降 with each successive layer. last layer's activations 接近 0. 这导致 backpropagated gradients 小^^^  
**Initial param same value:** Breaking symmetry is intention for the random initialization. If initialize all weights same value (e.g. 1), each hidden unit will get same signal. E.g. all weights initialized to 1, each unit gets signal equal to sum of inputs (and outputs sigmoid (sum(inputs))). The gradient of all neurons in one layer 会一样. 这导致 train 只会有效更新 neuron in different layers. .^^^FC network 用 tanh, no bias, only weights: **Initial weights to 0** result in 0 gradients + no learning.  
**Initial weights 大** 导致 vanishing gradients. **Vanishing gradients lead to information loss, making it difficult to propagate gradients upstream, while exploding gradients cause training instability or overflow errors due to excessively large steps in the loss landscape.**

(不对) Initial weights 一样值导致 gradient descent 后 weights 还 equal

$$\mathcal{L} = \sum_{i=1}^m \left( \log \sum_{j=1}^c e^{a_j(x^{(i)})} - a_{y^{(i)}}(x^{(i)}) \right) \quad \text{gradient:}$$

$$\frac{\partial \mathcal{L}}{\partial w_k} = \sum_{i=1}^m \left( \frac{e^{a_k(x^{(i)})}}{\sum_{j=1}^c e^{a_j(x^{(i)})}} x^{(i)} - \mathbb{1}(k = y^{(i)}) \cdot x^{(i)} \right)$$

$$\frac{\partial \mathcal{L}}{\partial b_k} = \sum_{i=1}^m \left( \frac{e^{a_k(x^{(i)})}}{\sum_{j=1}^c e^{a_j(x^{(i)})}} - \mathbb{1}(k = y^{(i)}) \right)$$

Softmax classifier

MTR:

$$\tilde{\mathcal{L}}(\theta) = \overset{\text{loss fn.}}{\mathcal{L}(\theta; \mathbf{x}, \mathbf{y})} + \alpha \overset{\text{norm penalty}}{\Omega(\theta)}$$

$$\alpha = 0? \quad \tilde{\mathcal{L}}(\theta) = \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})$$

$$\alpha \rightarrow \infty? \quad \tilde{\mathcal{L}}(\theta) \rightarrow \alpha \Omega(\theta)$$

Consider a model  $\tilde{\mathcal{L}}(\theta) = \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\theta)$  where  $\mathcal{L}(\theta; \mathbf{X}, \mathbf{y})$  is some loss function and  $\Omega(\theta)$  is some norm penalty. What are the effects on the model when  $\alpha = 0$  and  $\alpha \rightarrow \infty$ ?

有 loss func, model 更可能 overfit training data,  $=\infty$  no learning.^^^当 minimize the negative log likelihood for a classification problem with c class, 等于去做: Maximizing the likelihood of observing the training data. = Minimizing the Cross Entropy loss 不等于 Minimizing the Mean Squared Error.^^^ A 5-fold **cross-validation** fit 5-different model + (不对) A 5-fold cross-validation 导致 1 model instance being fitted over and over again 5 times + (不对) A 5-fold cross-validation 导致 5-different model instances being fitted over and over again 5 times. 1 fold 作为 validation, k-1 个 fold 作为 train+test, 重复决定哪些 folds 作为 train+test, 前提是有一个 Separate testing set, training set split into train 和 validation data. // **Never do validation on test set.**

**Hyperparameter**(任何在 GC 没有被 learnt 的东西): Loss Function+ Learning Rate+ Number of Layers+ Batch Size.

**Multi-layer Perceptron (MLP)** increasing the number of units within each layer can enhance the model's complexity: (i) 如果 model **underfitting** on the training data, 加更多 units allows the MLP to capture more complex patterns, 可能降低 training and testing errors. (ii) 但这可能也会导致 **overfitting**, as the increased model capacity can make it overly sensitive to the training data. Consequently, while the training loss may continue to decrease, the testing error could increase

**Role of the bias correction step in the Adam optimizer:** The Adam optimizer utilizes running averages to estimate the gradient and its square, which are initially biased towards 0 due to their initialization. As a result, the optimizer may take larger steps in the early stages of training, leading to unstable training and slower convergence. To address this, a bias correction step is introduced to adjust these estimations, making them more accurate. As training progresses and the estimations become more precise, the bias correction factor gradually approaches 1, diminishing the impact of bias correction.