

Deep Learning: Homework 1

Instructed by *Yi Wu*

Due on Mar 10, 2021

Runlong Zhou YaoClass 82 2018011309

1.1

Answer:

- Convolutional layer.

$$\frac{d(\text{weight}(j, k) \star z(k))(h, w)}{dz(a, b, c)} = \mathbb{I}\{k = a, h \leq b < h + k_H, w \leq c < w + k_W\} \text{weight}(j, k, b - h, c - w),$$

hence,

$$\begin{aligned} \frac{dy(j, h, w)}{dz(a, b, c)} &= \sum_{k=0}^{C_{\text{in}}-1} \mathbb{I}\{k = a, h \leq b < h + k_H, w \leq c < w + k_W\} \text{weight}(j, k, b - h, c - w) \\ &= \mathbb{I}\{h \leq b < h + k_H, w \leq c < w + k_W\} \text{weight}(j, a, b - h, c - w). \end{aligned}$$

Further,

$$\begin{aligned} \frac{dL}{dz(a, b, c)} &= \sum_{j, h, w} \frac{dL}{dy(j, h, w)} \cdot \frac{dy(j, h, w)}{dz(a, b, c)} \\ &= \sum_{j=0}^{C_{\text{out}}-1} \sum_{h=b-k_H+1}^b \sum_{w=c-k_W+1}^c \text{weight}(j, a, b - h, c - w) \frac{dL}{dy(j, h, w)}. \end{aligned}$$

$$\frac{dy(j, h, w)}{d\text{weight}(j, a, b, c)} = z(a, h + b, w + c),$$

hence,

$$\frac{dL}{d\text{weight}(j, a, b, c)} = \sum_{h, w} z(a, h + b, w + c) \frac{dL}{dy(j, h, w)}.$$

$$\frac{dL}{dbias(j)} = \sum_{h, w} \frac{dL}{dy(j, h, w)}.$$

- Max-pooling layer.

$$\frac{dy(j, h, w)}{dz(a, b, c)} = \mathbb{I}\left\{j = a, (b, c) = \arg \max_{(m, n) \leq (k_H-1, k_W-1)} z(j, \text{stride}[0] \times h + m, \text{stride}[1] \times w + n)\right\},$$

hence,

$$\begin{aligned} & \frac{dL}{dz(a, b, c)} \\ &= \sum_{j, h, w} \frac{dL}{dy(j, h, w)} \cdot \frac{dy(j, h, w)}{dz(a, b, c)} \\ &= \sum_{h, w} \mathbb{I} \left\{ (b, c) = \arg \max_{(m, n) \leq (k_H - 1, k_W - 1)} z(a, \text{stride}[0] \times h + m, \text{stride}[1] \times w + n) \right\} \frac{dL}{dy(a, h, w)}. \end{aligned}$$

• **Tanh.**

$$\frac{dy}{dz} = \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} = 1 - y^2,$$

hence,

$$\frac{dL}{dz} = \frac{dL}{dy} \cdot \frac{dy}{dz} = (1 - y^2) \frac{dL}{dy}.$$

2

Answer:

1. Three settings with smallest test losses are listed as follows:

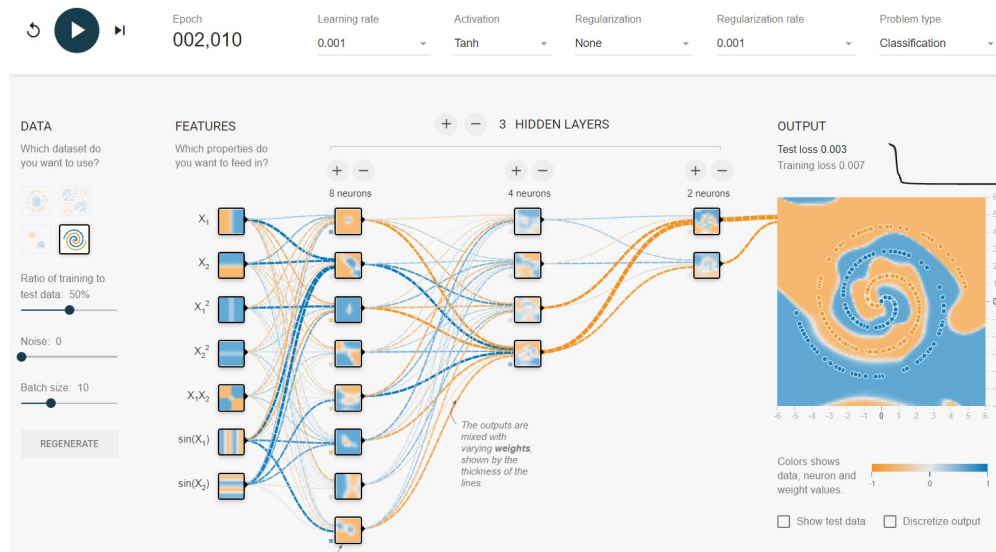


Figure 1: Test loss = 0.003

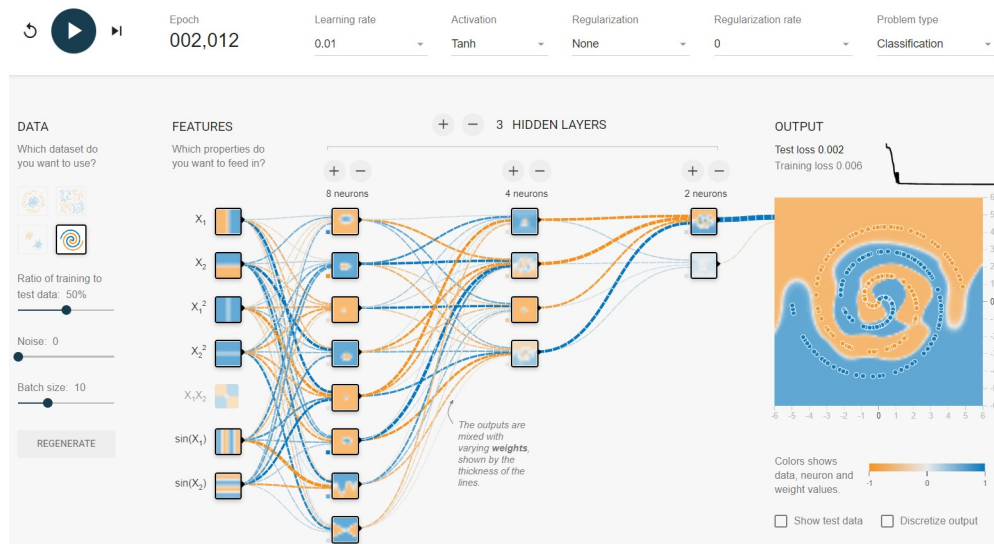


Figure 2: Test loss = 0.002

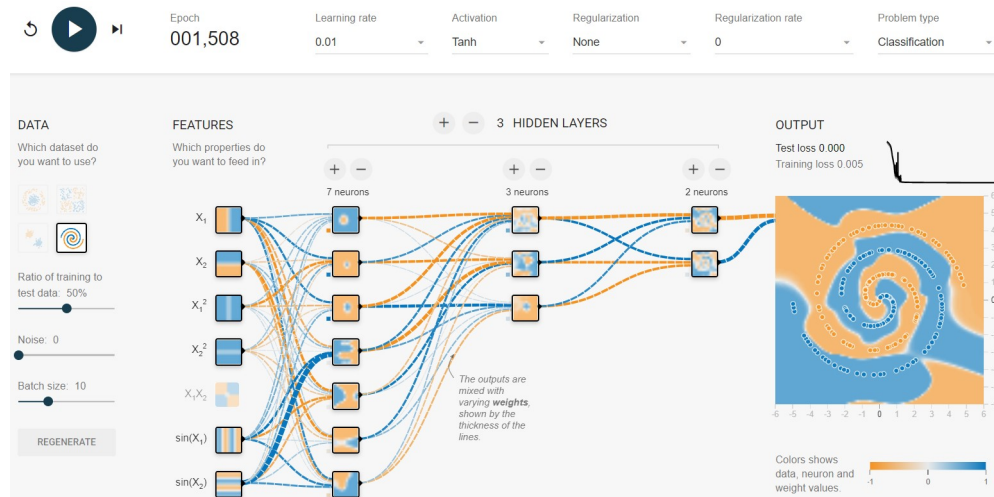


Figure 3: Test loss = 0.000

The first setting includes all the kernels so it works fairly well. By observing that x_1x_2 had diminishing weights during training, it was removed in the second setting. By cutting more neurons out, a lighter network (the last setting) is attained.

2. **Learning rate:** Higher learning rates make losses decrease faster in the beginning, but the losses are more likely to jitter afterwards. Lower learning rates give more smooth optimizing process, but the convergence rate is slower.

Hidden nodes: More hidden nodes make convergence slower and do not necessarily mean lower test loss. Extremely small number of nodes cannot learn well either.

Regularization: Higher regularization factor makes convergence faster, but may converge to a worse model since the some weights are suppressed to lower values, nullifying some nodes.