

### 1.1 What is a gate in LSTM

- ☒ Sigmoid multiplication to real value
- ☐ Sigmoid multiplication to positive real value
- ☐ ReLU multiplication to real value
- ☐ ReLU multiplication to positive real value

问的是 tanh 的取值范围  $[-1, 1]$

kernel size  
↓

### 1.2 How many weights does an nn. Conv2D(3, 20, 5) layer have?

- ☐ 60
- ☐ 300
- ☐ 900
- ☒ 1500

$$3 \times 20 \times 5 \times 5$$

$$\frac{25}{60} = 1500$$

### 1.3 Using nn.avgpool(2, stride=4) on a 200x200 input results in

- ☐ 0% inputs not used
- ☐ 25% inputs not used
- ☐ 50% inputs not used
- ☒ 75% inputs not used

$$\frac{200-2}{4} + 1 = 50.5$$

$$\frac{2500}{6000} \times 4$$

#### Regression Losses

##### Systematic comparison

	optimizable	outlier-robust	$\epsilon$ -tolerant
0/1 loss	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
squared loss $(y - t)^2$	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
absolute loss $ y - t $	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
log-cosh loss	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

##### Note:

$$l = \frac{1}{\beta} \log \cosh(\beta(y - t))$$

► Many further loss functions have been proposed in the literature (e.g. Huber's loss,  $\epsilon$ -sensitive loss, etc.). They often implement similar desirable properties as the log-cosh loss.

##### Systematic comparison

	optimizable	mislabeling-robust	builds margin
0/1 loss	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
perceptron loss $\max(0, -yt)$	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
log loss $\log(1 + \exp(-yt))$	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

### 1.6 The 0/1 loss $1_{\{y \neq t\}}$ is

- ☒ Not hard to optimize
- ☐ Not outlier robust
- ☐ ... to small variations
- ☐

### 1.7 The perceptron loss can be written as

- ☒  $\max(0, -yt)$
- ☐  $\max(0, -yt)^2$
- ☐  $1_{\{y \neq t\}}$
- ☐  $1_{\{y \neq t\}}$  for some  $\epsilon > 0$

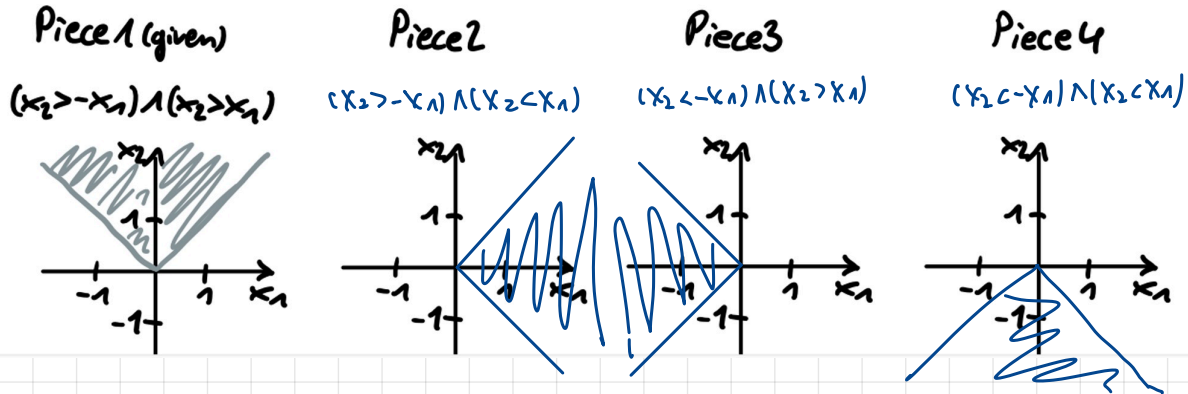
### 1.8 Small local invariances can be addressed by using

- ☐ Activation layers
- ☒ Pooling layers
- ☐ Convolution layers
- ☐ Linear layers

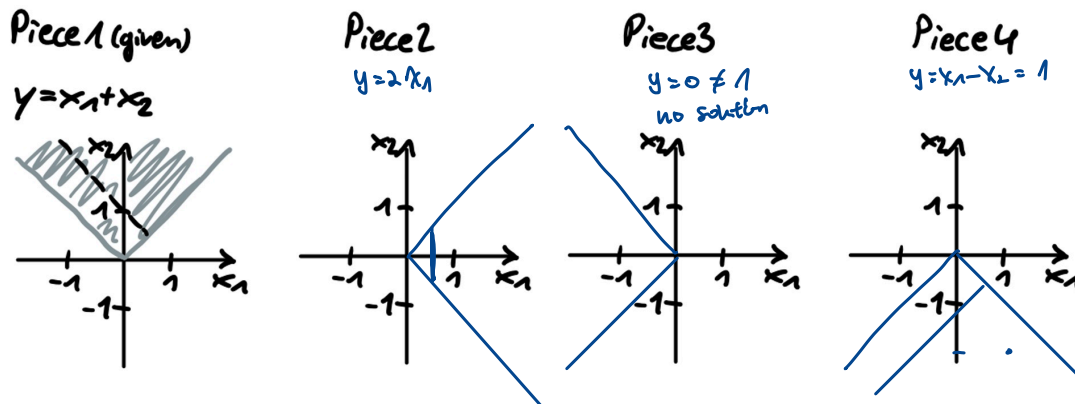
## 2. Task: Decision boundary

Given:  $a_3 = \max(x_1 + x_2)$   
 $a_4 = \max(x_1 - x_2)$   
 $y = a_3 + a_4$

a) Draw the four pieces on which  $y$  is linear



b) Draw the decision boundary on each piece for  $y = f(x) = 1$



## 3. Task: Gradient, bound and regularization

Given:  $z_{ij} = |x_i - c_{ij}|$ ,  $a_j = \exp(-\sum_{i=1}^d z_{ij})$ ,  $y = \sum_{j=1}^h w_j a_j$

a) Calculate  $\frac{\partial y}{\partial x_i}$

$$\begin{aligned} \frac{\partial y}{\partial x_i} &= \sum_{j=1}^h w_j \frac{\partial a_j}{\partial x_i} \frac{\partial a_j}{\partial z_{ij}} \frac{\partial z_{ij}}{\partial x_i} \\ &= \sum_{j=1}^h w_j \cdot \exp\left(-\sum_{i=1}^d z_{ij}\right) \cdot (-1) \cdot \text{sign}(x_i - c_{ij}) \\ &= - \sum_{j=1}^h w_j a_j \text{sign}(x_i - c_{ij}) \end{aligned}$$

b) Show that  $\left\| \frac{\partial y}{\partial x} \right\| \leq d \|w\|_1$  Hint:  $|\sum_i a_i| \leq \sum_i |a_i|$

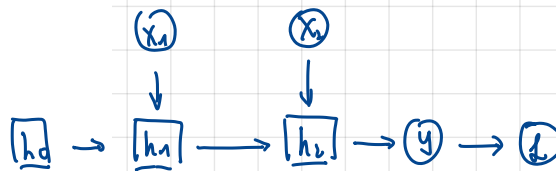
$$\begin{aligned} \left\| \frac{\partial Y}{\partial X} \right\| &= \sum_{i=1}^d \left| \frac{\partial Y}{\partial x_i} \right| \\ &= \sum_{i=1}^d \left| \sum_{j=1}^h w_{ij} a_j \right| \\ &\leq \sum_{i=1}^d \sum_{j=1}^h |w_{ij}| \quad \text{da } a_j \leq 1 \\ &= d \|W\| \end{aligned}$$

c) Explain in 1-2 sentences how this bound helps with choosing the regularization

If we regularize (penalize)  $\|W\|$  to be keep it small, then the Gradient  $\frac{\partial Y}{\partial X}$  is also small which helps avoid overfitting

#### 4. Task: RNN

Given:  $h_1 = \tanh(x_1^T W + h_0)$   
 $h_2 = \tanh(x_2^T W + h_1)$   
 $y = h_2$   
 $\mathcal{E} = (y - t)^2$



a) Draw the associated NN graph

$$\frac{\partial \mathcal{E}}{\partial y} = 2(y - t)$$

b) Calculate  $\frac{\partial \mathcal{E}}{\partial y}$

c) State the derivative of  $\frac{\partial \mathcal{E}}{\partial W}$  using the chain rule

d) Calculate  $\frac{\partial \mathcal{E}}{\partial W}$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial W} &= \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial h_2} \left( \frac{\partial h_1}{\partial W} + \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial W} \right) \\ &= 2(y - t) \cdot 1 \cdot \left( \tanh'(x_1^T W + h_0) \cdot x_1 + \tanh'(x_2^T W + h_1) \cdot \tanh'(x_1^T W + h_0) \cdot x_1 \right) \\ &= 2(y - t) \tanh'(x_2^T W + h_1) (x_2 + \tanh'(x_1^T W + h_0) x_1) \end{aligned}$$

#### 5. Task: Programming

a) Program an adversarial attack on a trained model using  $\min \|x - z\|^2 + \max(0, t \cdot f(z))$

```

1 import torch, import...
2 def adv_sarial(x, T, model):
3     # Create/modify z
4     z = x.clone().detach().requires_grad_(True)
5
6     optim = torch.optim.SGD([z], 0)
7     # create optimization
8     for i in range(1000):
9         optim.zero_grad()
10        loss = torch.norm(x - z, p=2) + torch.max(torch.zeros_like(T), T * model(z)).sum()
11        loss.backward()
12        optim.step()
13
14    return z
  
```

$x: 1 \times d$   
 $t: 1 \times d$   
 $model(x): 1 \times 1$   
 use 1000 iterations, stopping not needed

b) How would this code need to be modified to take  $x \in \mathbb{R}^{n \times d}$  as input? Also, it takes logs and provides line numbers!

c) How can the model be made more robust?