

Week4 Recitation

Zoey Song
CS190I Deep Learning

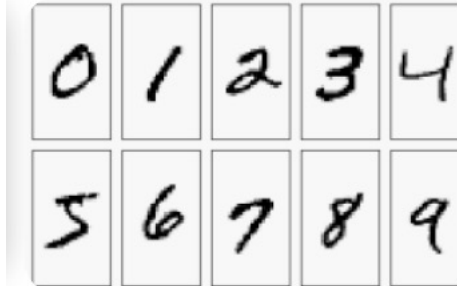


Outline

- **Training & Inference**
- **Gradient Descent**
- **Regularization**
- **Example for Forward & Backward propagation**

Classification Task

- E.g., hand-written digit classification
- Training
 - Cross-entropy
 - Model output: $\hat{y} = \text{Softmax}(Wh+b)$
 - Loss: $l = -y \log \hat{y}$
- Inference
 - $i = \underset{i \in \{1,2,\dots,k\}}{\text{argmax}} \hat{y}_i$



Regression Task

- E.g., predict stock price
- Training
 - MSE-loss
 - Model output: $\hat{y} = f(Wx+b)$
 - Loss: $l = \frac{1}{2} (y - \hat{y})^2$
- Inference
 - $\hat{y} = f(Wx+b)$

Gradient Descent

- Considering the optimization objective

$$\min_{\theta} F(x, \theta)$$

- Suppose we have N samples, at time step t

- Full gradient descent

$$\theta^t = \theta^{t-1} - lr_t \frac{1}{N} \sum_{i=1}^N \nabla F(x_i, \theta^{t-1})$$

- Highly Efficient
- but can not fully utilize the data

Gradient Descent

➤ Suppose we have N samples, at time step t

➤ Stochastic gradient descent

$$\theta^t = \theta^{t-1} - lr_t \nabla F(x_k, \theta^{t-1}), k \text{ in } \{1, 2, \dots, N\}$$

➤ Computation is slow

➤ Can fully utilize the training data

➤ Mini-batch with batch size s

$$\theta^t = \theta^{t-1} - lr_t \frac{1}{s} \sum_{k \in t_s} \nabla F(x_k, \theta^{t-1}), t_s \subseteq \{1, 2, \dots, N\}$$

➤ Trade-off between the computation speed and data use

Regularization

- Prevent overfitting
- L1 regularizer (Lasso regularizer)

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^N |w_i|$$

- E.g., Alignment in machine translation, graph for social network

- L2 regularizer (Ridge regularizer)

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

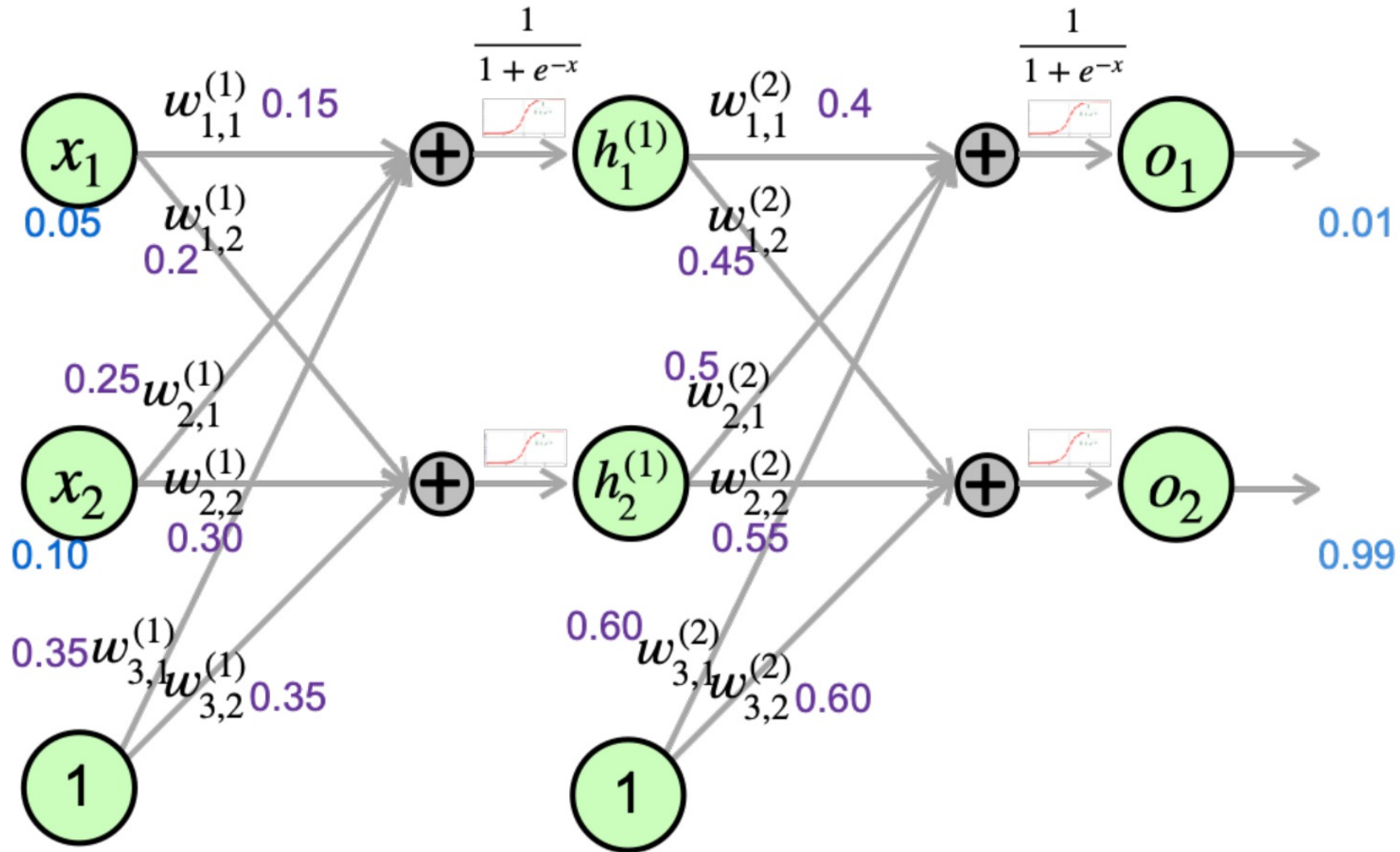
$\min_{\theta} J(\theta)$

- E.g. deep neural network training

- Dropout

- 0.1 ~ 0.5

Quiz



Why calculating gradient is necessary?

- **Sometimes we can't directly compute the gradient**
 - Incorporating a latent variable into the MLE objective

$$\log P(x; \theta) = \log \int P(x, z; \theta) dz$$

Any Question?