

## 3-2 Regularization

Zhonglei Wang

WISE and SOE, XMU, 2025

# Contents

1. Penalties on parameters

2. Dropout

3. Early stopping

# Intuition

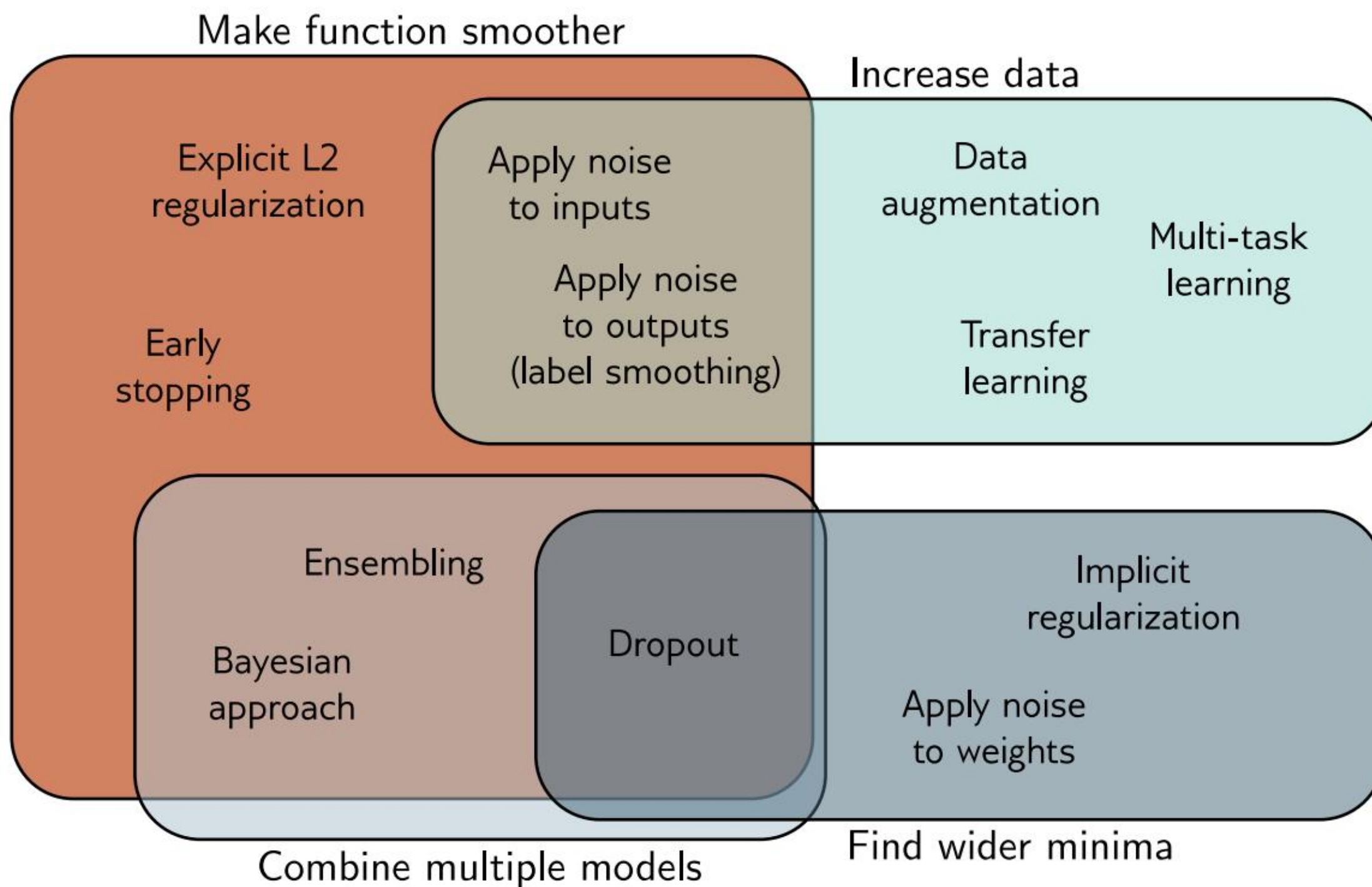
1. Fully connected neural network often involves too many model parameters
  - Complex models tend to **overfit** the training data
  - **Decrease the generality** of the trained model
2. Definitely, we can consider a simpler neural network
  - Simpler models may **underfit** the training data, however
3. Thus, we may would like to obtain a kind of less complex model **without changing** the architecture

# Intuition

1. In machine learning, regularization
  - usually adds additional terms to the loss function
  - improves generality of the deep learning model
2. There exist implicit regularization, including dropout, early stopping, ensembling

# Intuition

1. The following image is Figure 9.14 of Prince (2024)



# Penalties on parameters

## 1. $l_2$ penalty

$$\mathcal{J}_2 = \frac{\lambda}{2n} \sum_{l=1}^L \|\mathbf{W}^{[l]}\|_F^2$$

- $\lambda$  : hyperparameter
- $\|\mathbf{A}\|_F = (\sum_i \sum_j a_{ij}^2)^{1/2}$  : Frobenius norm of a real-valued matrix  $\mathbf{A} = (a_{ij})$

- Derivative

$$\frac{\partial \mathcal{J}_2}{\partial \mathbf{W}^{[l]}} = \frac{\lambda}{n} \mathbf{W}^{[l]}$$

- It is similar to ridge regression and is used to control the complexity of the model

# Penalties on parameters

## 1. $l_1$ penalty

$$\mathcal{J}_1 = \frac{\lambda}{n} \sum_{l=1}^L \|\mathbf{W}^{[l]}\|_1$$

- $\|\mathbf{A}\|_1 = \sum_i \sum_j |a_{ij}|$

- Derivative

$$\frac{\partial \mathcal{J}_1}{\partial \mathbf{W}^{[l]}} = \frac{\lambda}{n} \text{sign}(\mathbf{W}^{[l]})$$

- It is used to
  - ▷ control the complexity of the model
  - ▷ **induce sparsity**, which is similar to LASSO



# Remarks

1. We only implement  $l_2$  or  $l_1$  penalty on the **weight terms**  $\{\mathbf{W}^{[l]} : l = 1, \dots, L\}$ 
  - By implementing, we mean that the  $l_2$  or  $l_1$  penalty is directly added on the original cost function
  - We leave the bias terms  $\{\mathbf{b}^{[l]} : l = 1, \dots, L\}$  unpenalized
2. We “misuse”  $\|\mathbf{A}\|_1$  to denote the summation of absolute values of elements in  $\mathbf{A}$ 
  - More generally,  $\|\mathbf{A}\|_1 = \max_j \sum_i |a_{ij}|$  denotes the maximum **column summation** of absolute elements



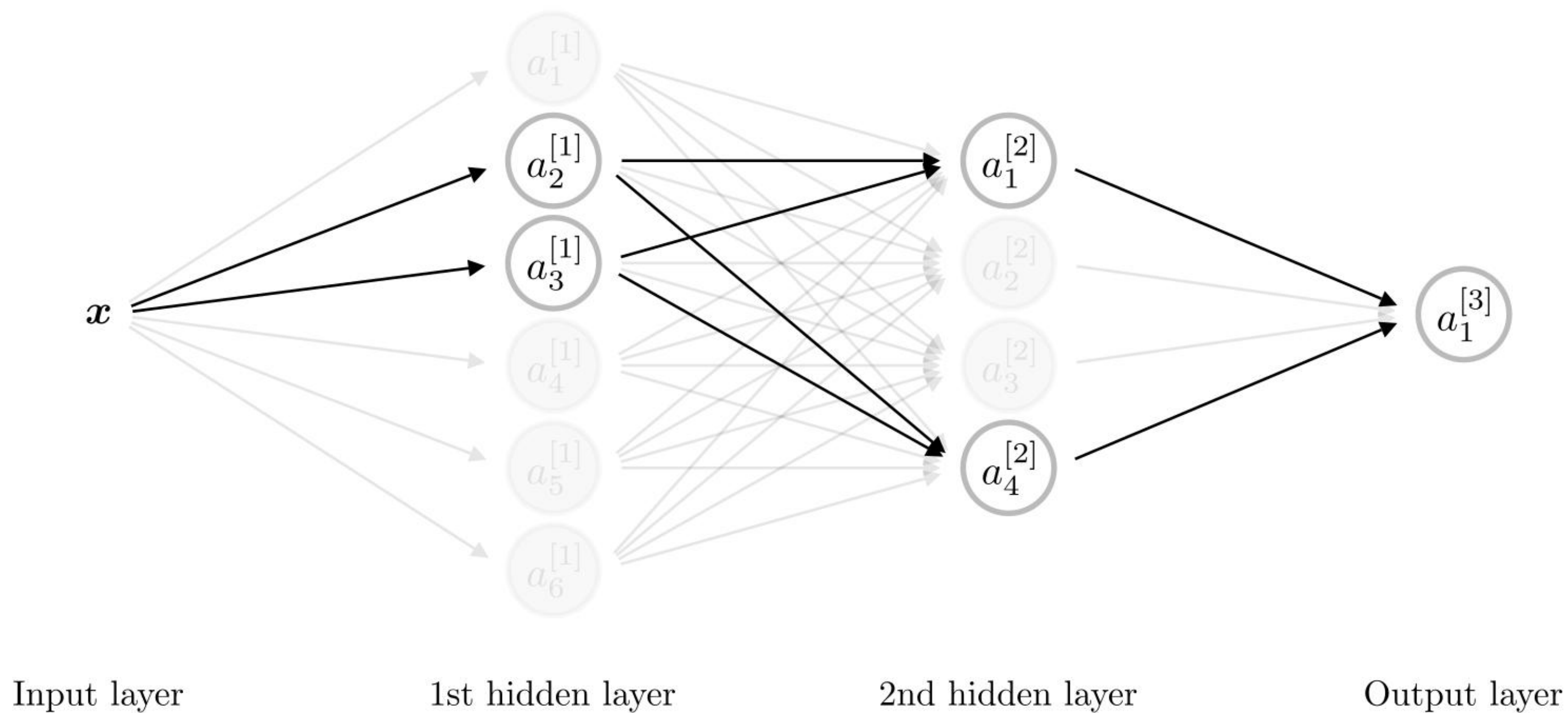
# Dropout

1. Commonly used regularization **during the training step**
  - Randomly remove neurons in each **hidden** layer
  - If a neuron is removed, its activated value is set to be 0
  - Removal is done for each training example individually
2. Improves the robustness

# Dropout

1. For activated values of each training example in the  $l$ th layer ( $l = 1, \dots, L - 1$ )
  - Set the dropout rate to be  $p^{[l]}$
  - Randomly assign them to be 0 using probability  $p^{[l]}$
  - Need to rescale the remainings by  $(1 - p^{[l]})^{-1}$  to compensate information loss
2. Dropout is implemented using a **mask matrix** for each layer

# Dropout



# Dropout -- training procedure

1. Let  $\mathbf{M}^{[l]} \in \{0, 1\}^{n \times d^{[l]}}$  be a mask matrix
  - Each element of  $\mathbf{M}^{[l]}$  is a Bernoulli random variable with success probability  $1 - p^{[l]}$
2. Forward propagation
$$\mathbf{A}_d^{[l]} = \frac{\mathbf{A}^{[l]} \circ \mathbf{M}^{[l]}}{1 - p^{[l]}}.$$
3. Backpropagation: as usual
  - No additional model parameters are introduced for dropout

# Dropout -- training procedure

1. Dropout is multiplying Bernoulli random noise to the neural network
2. Generally, we can add or multiply other noise to the neural network
  - Add noise to the training data
  - Add noise to the weights
  - Perturb the labels

# Early stopping

1. The performance on the training dataset improves as training procedure precedes
2. It may overfit the training dataset
  - The performance on the validation/test dataset may not always improve
3. Early stopping monitors the performance on the validation set
  - Stop the training if the performance does not improve on the validation set



# Ensembling

1. To improve generalization of the model, we may
  - Build several models and take their average as the final model
    - ▷ For example, using different initializations
  - Use resampling methods to generate different training sets
    - ▷ This method is commonly used to get uncertainty of the deep learning model