# L2 Regularization

03 January 2025    18:58

Applied to Model Weights:

- Regularization is applied to the weights of the model to penalize large values and encourage smaller, more generalizable weights.

Introduced via Loss Function or Optimizer:

- Adds a penalty term $\lambda \sum w_i^2$ to the loss function in L2 regularization.

$$\text{Loss}_{reg} = \text{Loss}_{original} + \lambda \sum w_i^2$$

$$\lambda \left( w_1^2 + w_2^2 + w_3^2 + w_4^2 \right)$$

- In weight decay, directly modifies the gradient update rule to include $\lambda w_i$, effectively shrinking weights during training.

$$w \leftarrow w - \eta(\nabla \text{Loss} + \lambda w)$$

weight_decay

Penalizes Large Weights:

- Encourages the network to distribute learning across multiple parameters, avoiding reliance on a few large weights.

Reduces Overfitting:

- Helps the model generalize better to unseen data by discouraging overly complex representations.

Controlled by a Hyperparameter:

- A regularization coefficient ($\lambda$ often set via weight_decay in optimizers) controls the strength of the penalty. Larger values lead to stronger regularization.

No Effect on Bias Terms:

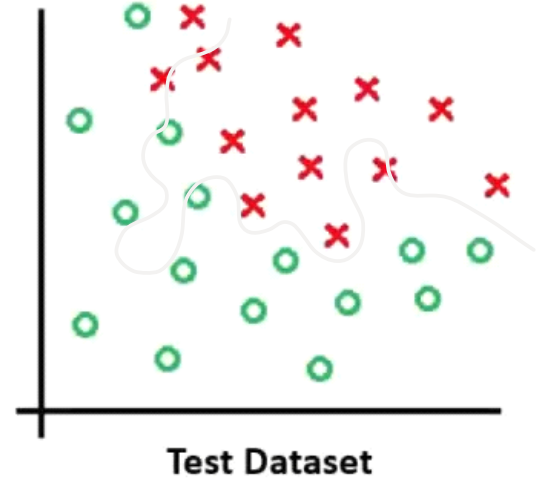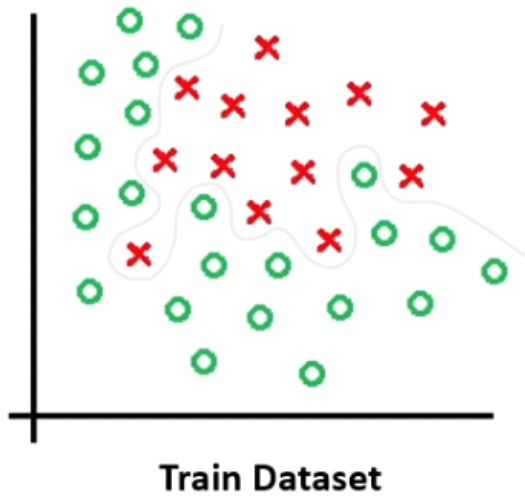- Regularization is typically applied only to weights, not biases, as biases don't directly affect model complexity.

Active During Training:

- Regularization affects weight updates only during training. It does not explicitly influence the model during inference.
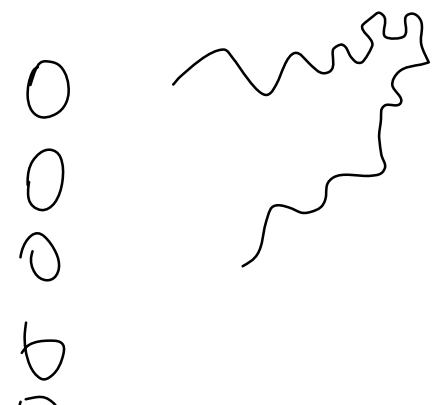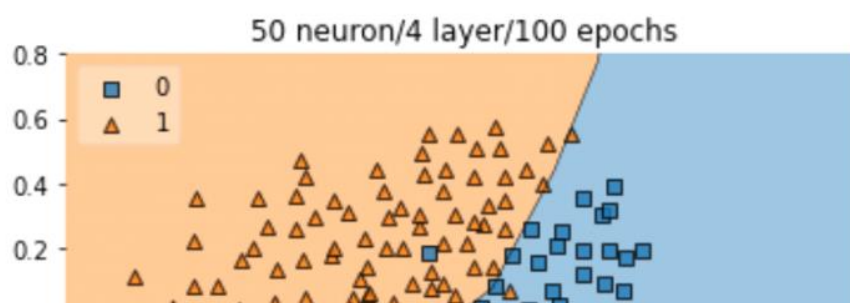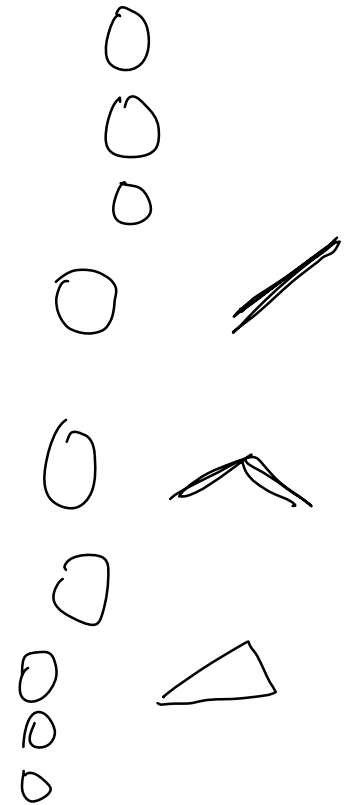
optimization

$L \rightarrow$ mse

$\rightarrow$ logloss

weights/bias

# Overfitting

**Train Dataset**



**Test Dataset**

# Why Neural Networks Overfit?

## 1 neuron/1 layer/100 epochs



perceptron

## 10 neuron/4 layer/100 epochs



## 50 neuron/4 layer/100 epochs

257 neuron/4 layer/100 epochs



1025 neuron/4 layer/100 epochs



+ weight ≃ 0

Simple

# Ways to solve overfitting

**Adding more data**

→ Add more rows

→ Data augmentation

**Reducing the compuxity model**

→ Dropout

→ Early stopping

→ Regularization

$\dfrac{L1}{}$

$\dfrac{L2}{} \to W$

$\dfrac{L1 + L2}{}$

# Regularization

$ANN \longrightarrow$ weights / bias

$\qquad \downarrow$ min Loss function

$L =$ mse

$\qquad \downarrow$ binary

$\boxed{L2}$

$\qquad L1$

$\boxed{W_1 \rightarrow W_{10}}$

$$C = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{y}_i) + \text{penalty term}$$

$$C = L + \frac{\lambda}{2n} \sum_{i=1}^{k} \|W_i\|^2 \qquad \text{weightage}$$

$$\frac{\lambda}{2n} \left[ W_1^2 + W_2^2 + \cdots + W_{10}^2 \right]$$

$\lambda =$ hyperparameter

$\boxed{\lambda = 0}$

$$C = \sum L(y_i, \hat{y}_i) + P$$

$\qquad\qquad\qquad\qquad \downarrow$
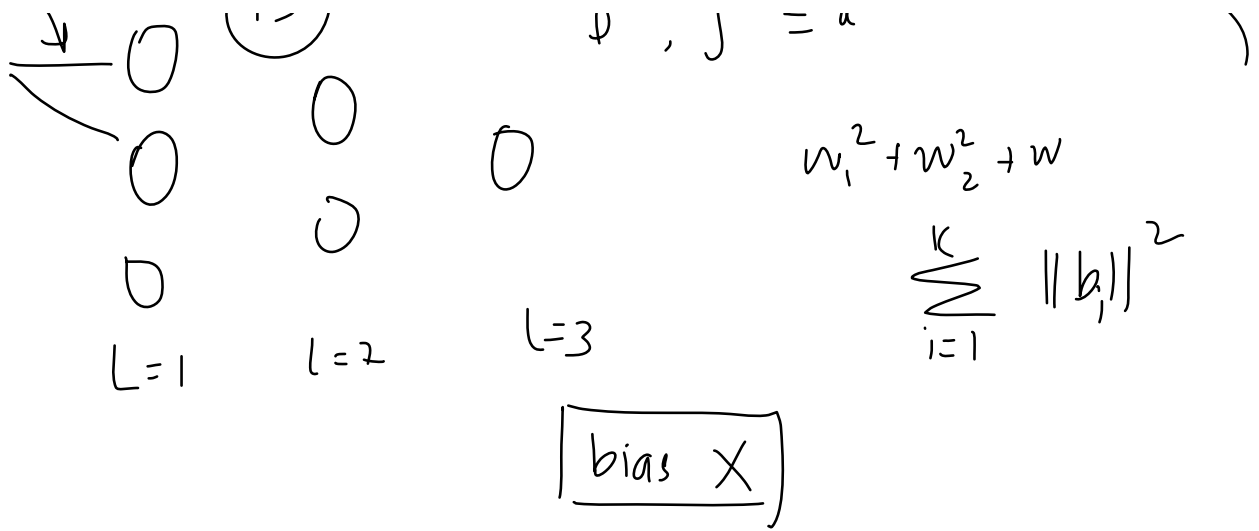
$\qquad\qquad\qquad\qquad W \simeq 0$

$L2 \rightarrow L1$

$\rightarrow L1$ norm

$$C = L + \frac{\lambda}{2n} \sum \|W_i\|$$

$$C = \sum_{i=1}^{n} L(y_i, \hat{y}_i) + \sum_{l=1}^{L} \sum_{i=1}^{} \sum_{j=1}^{} \|W_{ij}^{l}\|^2 \qquad \boxed{W \simeq 0}$$

$\rightarrow 0 \qquad \boxed{15} \qquad\qquad i , j = l$

$\psi$ , $J = a$ )

$w_1^2 + w_2^2 + w$

$$\sum_{i=1}^{K} \|b_i\|^2$$

L=1          l=2          L=3

bias X

# Intuition behind Regularization

$$W_\eta = W_0 - \eta \left( \frac{\partial L}{\partial W_0} \right) \qquad L$$

$$\boxed{1 - \eta \lambda}$$

positive

$$L' = L + \frac{\lambda}{2} \sum \|W_i\|^2$$

$$W_2^1 + W_2^2 + W_3^3 + \cdots \cdots$$

$$2W$$

$$\frac{\partial L'}{\partial W_0} = \frac{\partial L}{\partial W_0} + \frac{\lambda}{2} 2 W_0$$

$$W_\eta = W_0 - \eta \left( \frac{\partial L}{\partial W_0} + \lambda W_0 \right)$$

$$W_0 \qquad = \frac{\partial L}{\partial W_0} + \lambda W_0$$

$$W_\eta = W_0 - \eta \lambda W_0 - \eta \frac{\partial L}{\partial W_0}$$

$$W_0 << W_0$$

$$\boxed{W_\eta = (1 - \eta \lambda) W_b - \eta \frac{\partial L}{\partial W_0}}$$

L2 reg $\to$ weight decay     weight decay     $W_0$