# Computational Intelligence

Samaneh Hosseini

Isfahan University of Technology

# Outline

- Neural Networks in Practice

  - Regularization

    - Why Regularization Reduces Overfitting?

  - Regularization Techniques

    - L2 Regularization
    - Dropout

# Neural Networks in Practice: Regularization

# Regularization

## What is it?

Technique that constrains our optimization problem to discourage complex models

# Regularization

## What is it?

Technique that constrains our optimization problem to discourage complex models

## Why do we need it?

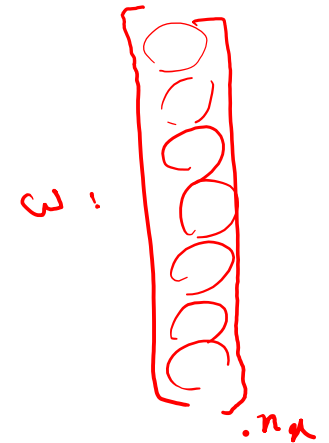Improve generalization of our model on unseen data

# Regularization in Logistic regression

$w \in R^n \qquad b \in R$

$$\min_{w,b} J(w,b)$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \boxed{\frac{\lambda}{2m} \|w\|_2^2} + \frac{\lambda}{2m} b^2$$

$L_2$ regularization: $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$

$\lambda$ = Regularization Parameter

$L_1$ reg... : $\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$

$w:$ $\quad .n_x$

# Regularization in Neural network

تعداد لایه‌ها

$$J(w^{[1]}, b^{[1]}, \ldots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{\ell=1}^{L} \| w^{[\ell]} \|_F^2$$

$$\| w^{[\ell]} \|_F^2 = \sum_{i=1}^{n^{[\ell]}} \sum_{j=1}^{n^{[\ell-1]}} (w_{ij}^{[\ell]})^2$$

$\| \cdot \|_2^2 \qquad \| \cdot \|_F^2$

"Frobenius norm"

$$dw^{[\ell]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[\ell]}$$

$$w^{[\ell]} = w^{[\ell]} - \alpha \, dw^{[\ell]} \qquad \frac{\partial J}{\partial w^{[\ell]}}$$

$$w^{[\ell]} = w^{[\ell]} - \alpha \left[ (\text{from back}) + \frac{\lambda}{m} w^{[\ell]} \right]$$

weight decay

$$= (1 - \frac{\alpha \lambda}{m}) w^{[\ell]} - \alpha (\text{from back} \cdots)$$
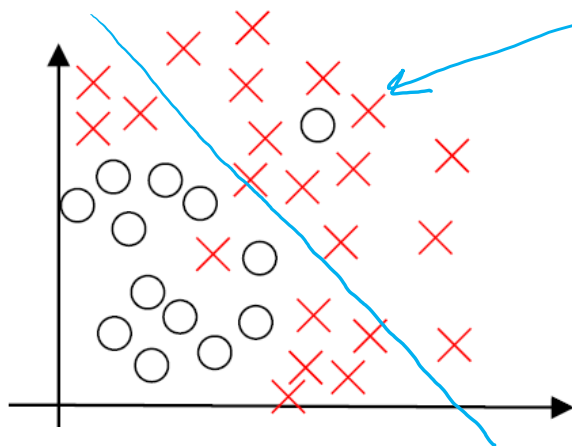
$n^{[\ell-1]} \qquad n^{[\ell]}$

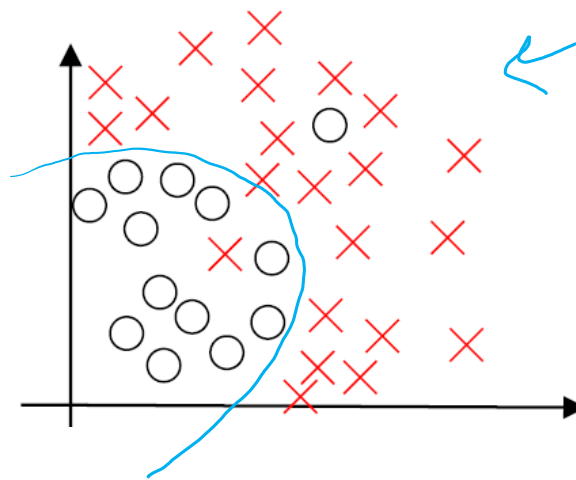# Regularization:
# Why Regularization Reduces Overfitting?

# How does regularization prevent overfitting?



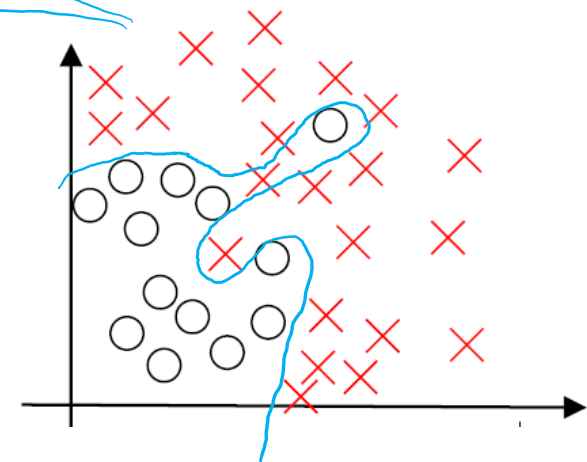$$J(w^{[1]}, b^{[1]}, \ldots w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{\ell=1}^{L} \|w^{[\ell]}\|_F^2$$

$$\lambda \rightarrow \text{بزرگ}$$

$$\|w^{[\ell]}\| \approx 0$$

high bias                    just right                    high variance

# How does regularization prevent overfitting?

$$g(z) = \tanh(z)$$

$\tanh$

$\lambda \uparrow \quad \omega^{[\ell]} \downarrow$

$$Z^{[\ell]} = \omega^{[\ell]} a^{[\ell-1]} + b^{[\ell]}$$

$Z^{[\ell]} \downarrow$

$\lambda$ large $\approx$ linear

# Regularization Techniques:
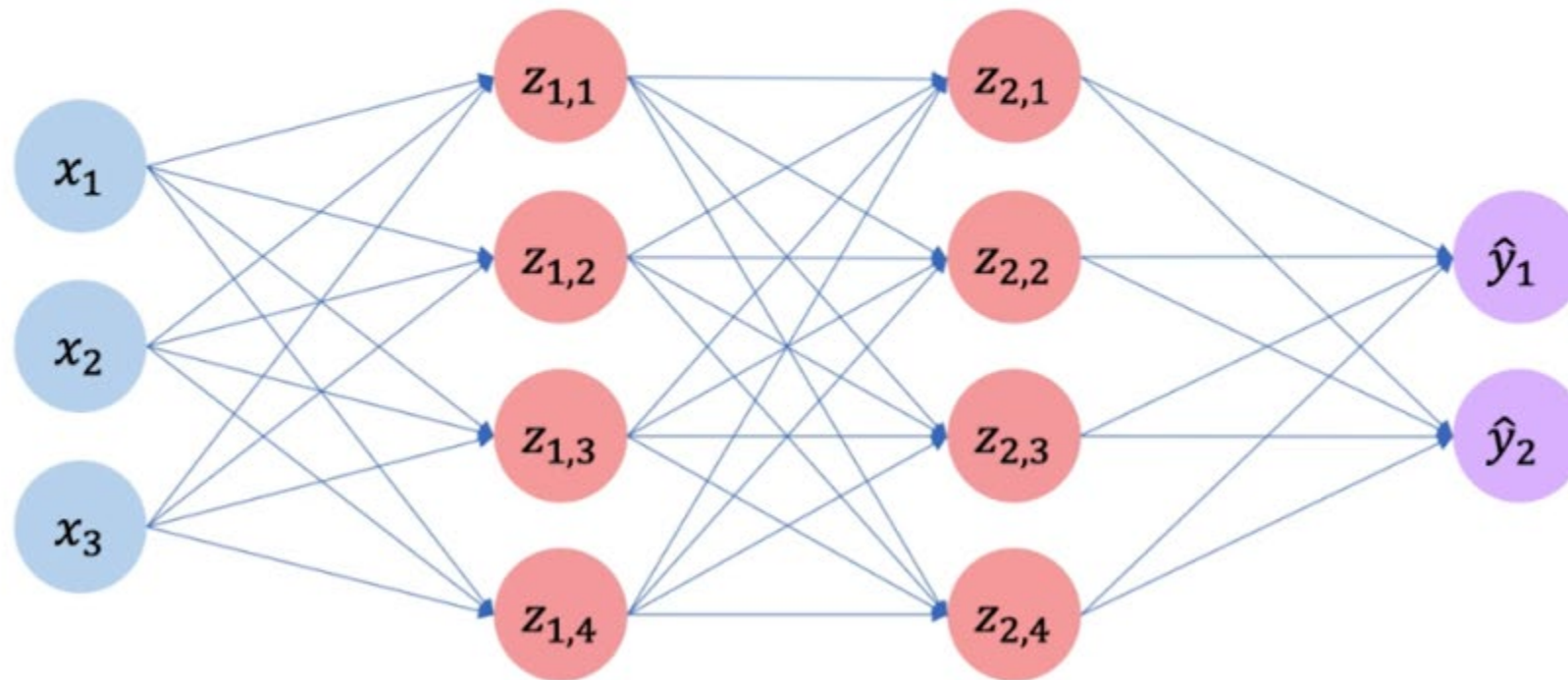## II: Dropout

I:      $L_2$ Regularization

II:     Dropout

$$\frac{\lambda}{2m} \sum_{i=1}^{\ell} \| w^{[\ell]} \|_F^2$$
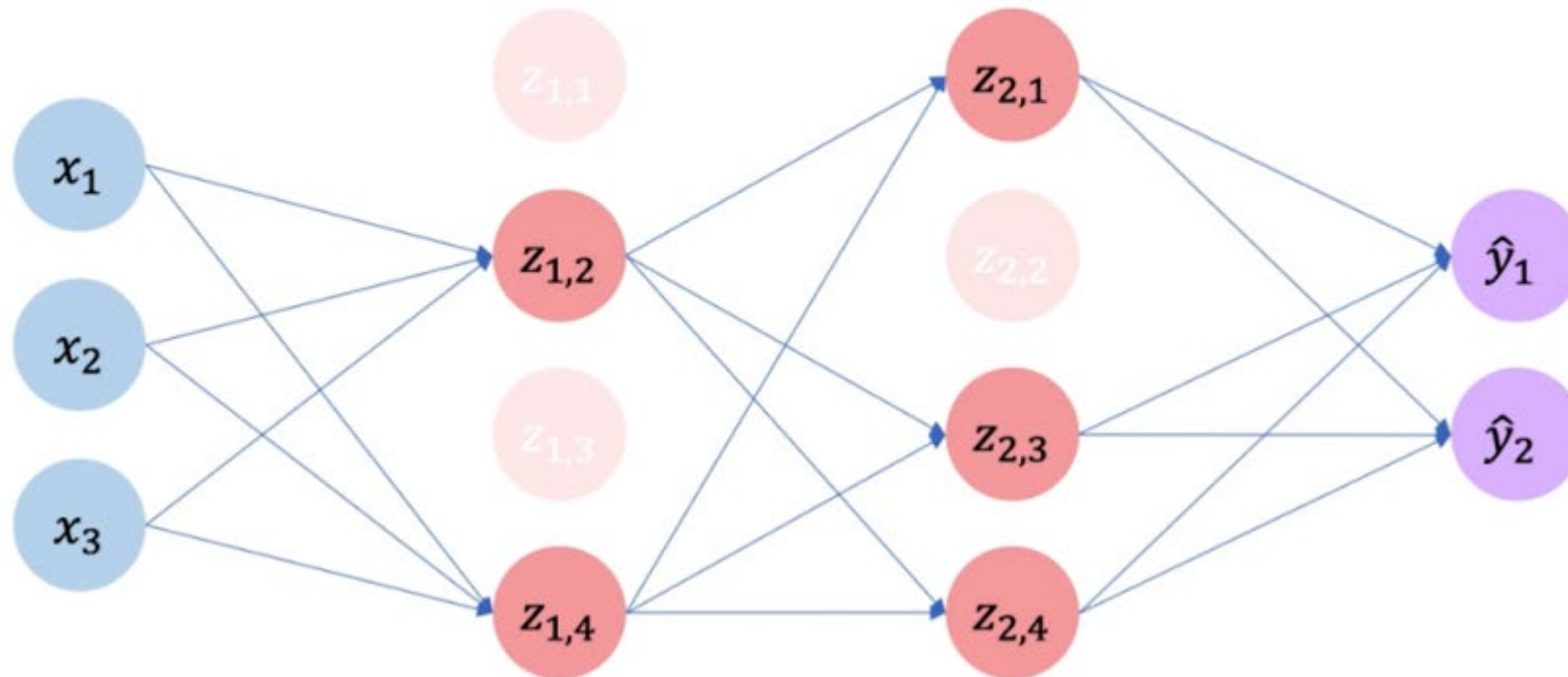
# Regularization II: Dropout

- During training, randomly set some activations to 0

# Regularization II: Dropout

- During training, randomly set some activations to 0
  - Typically 'drop' 50% of activations in layer
  - Forces network to not rely on any 1 node

`tf.keras.layers.Dropout(p=0.5)`

# Regularization II: Dropout

- During training, randomly set some activations to 0
  - Typically 'drop' 50% of activations in layer
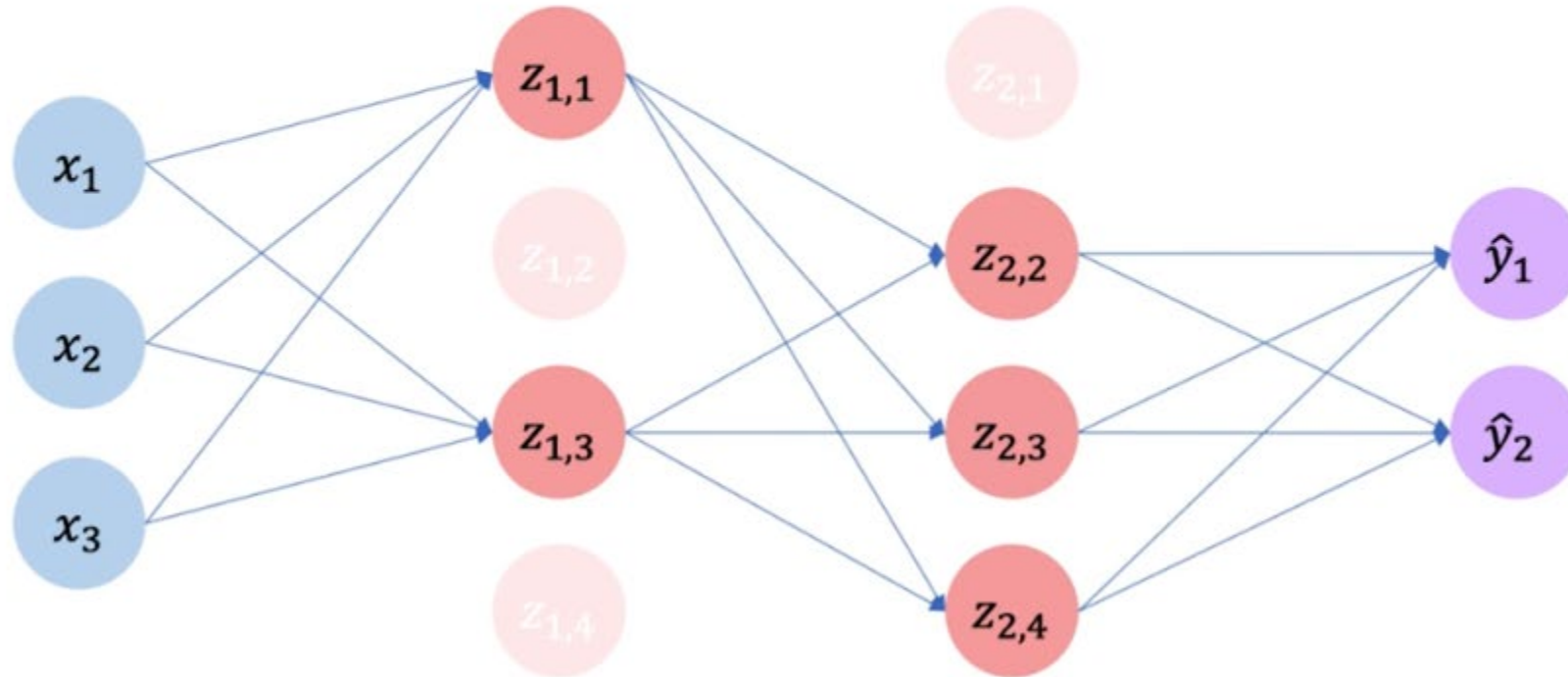  - Forces network to not rely on any 1 node

```
tf.keras.layers.Dropout(p=0.5)
```

# Implementing dropout ("Inverted dropout")

Illustrate with layer $l = 3$

keep-prob = 0.8     0.2

$d_3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep-prob$

$a3 = np.multiply(a3, d3)$

$a3 /= keep-prob$

50 unit $\longrightarrow$ 10 unit shut off

$Z^{[4]} = \omega^{[4]} a^{[3]} + b^{[4]}$

$\uparrow$ reduced 20%

$/= 0.8$

# Making predictions at test time

$$a^{[0]} = X$$

No drop out

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

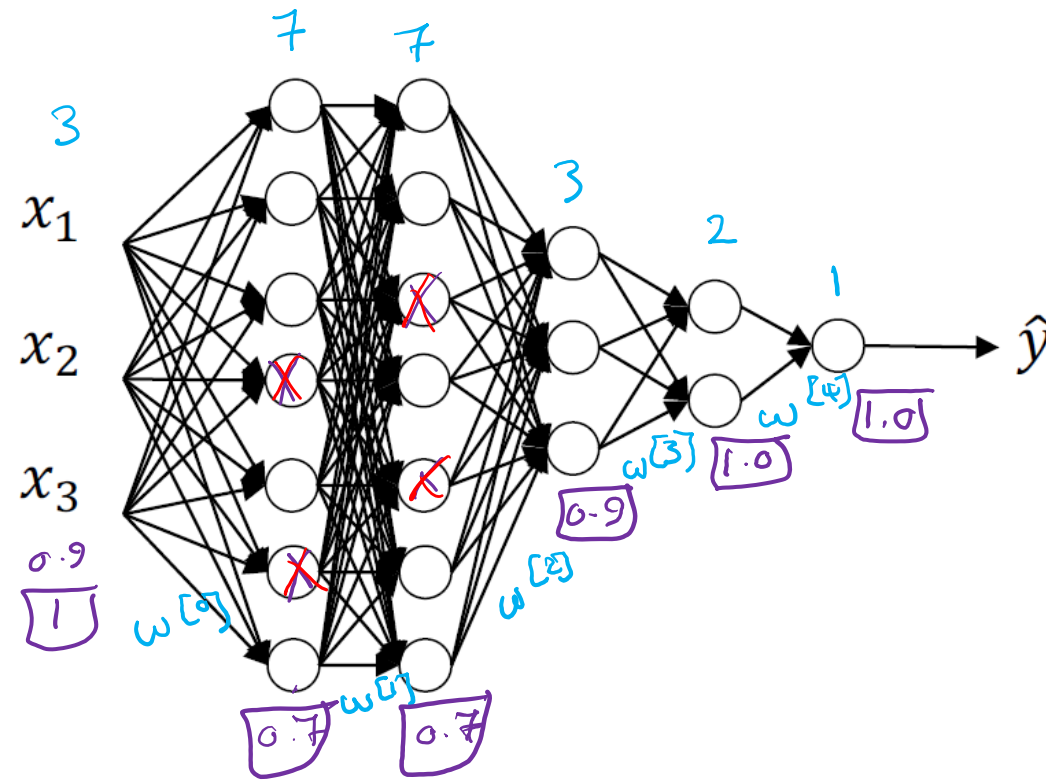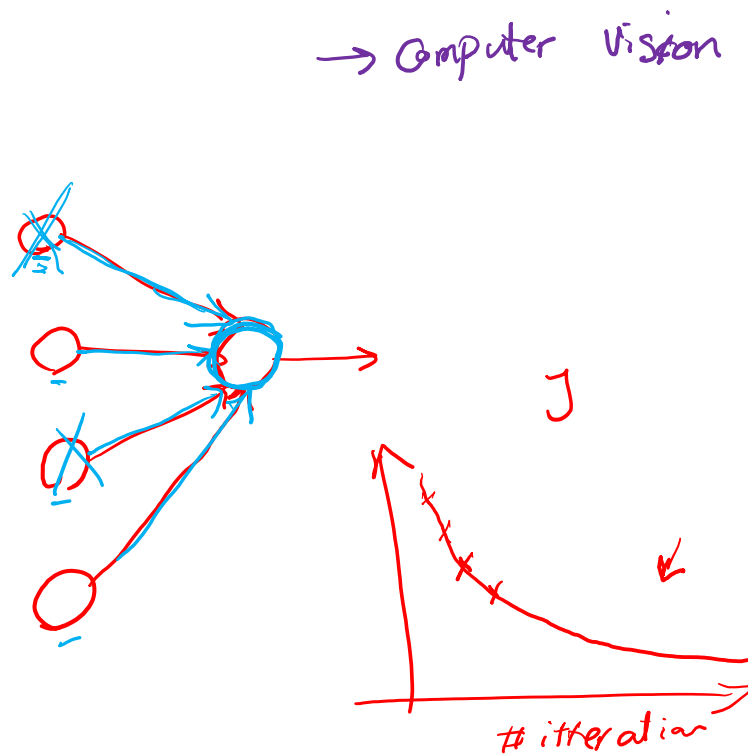$$\vdots$$

$$\hat{y}$$

$$/ = keep\_prob$$

# Understanding Dropout

# Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights

Shrink    weight    $\longrightarrow$  L2

→ Computer Vision

$x_1$   $x_2$   $x_3$

3    7    7    3    2    1

J

$\hat{y}$

dw

$7 - \alpha \, dw$

6

0.9
$\boxed{1}$    $\omega^{[1]}$    $\boxed{0.7}$ $\omega^{[1]}$  $\boxed{0.7}$    $\omega^{[2]}$    $\boxed{0.9}$ $\omega^{[3]}$  $\boxed{1.0}$    $\omega^{[4]}$ $\boxed{1.0}$

# iteration

Keep _ Prob

# Core Foundation Review

Regularization

How does regularization prevent overfitting

Regularization techniques: L2, Dropout

Technique that constrains our optimization problem to discourage complex models