

# Machine Learning for Engineering and Science Applications

The NPTEL logo is a circular emblem with a stylized flower or star in the center. The petals or rays of the flower are colored in a gradient from light yellow to light pink. The outer ring of the emblem is composed of alternating yellow and pink segments.

Gradient Descent-2

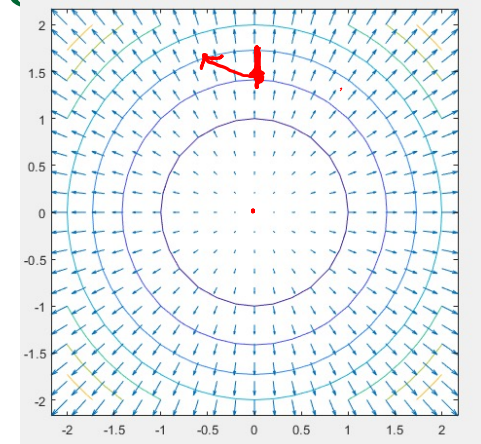
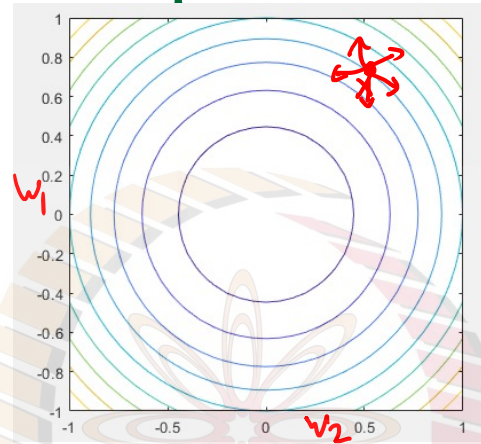
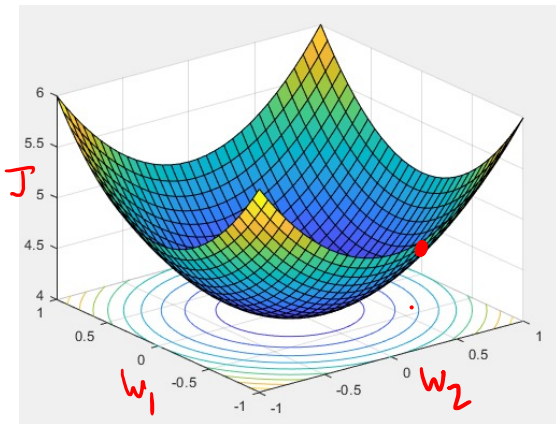
Proof of Steepest Descent

Numerical Gradient calculation

Stopping criteria

# Gradient and steepest change

$$\vec{w}, \vec{w} + \Delta \vec{w} \\ J, J + \Delta J \\ \text{Vectors} = \nabla J$$



**Claim** : The direction of maximum rate of change for a function  $J(w)$  is given by  $\nabla_w J$

**Proof**: Recall that rate of change in a given direction  $\mathbf{v}$  is given by  $\frac{\partial J}{\partial \mathbf{v}}$

That is, rate of change along  $\mathbf{v}$  is 
$$\frac{\partial J}{\partial v} = \nabla J \cdot \hat{\mathbf{v}} = |\nabla J| |\hat{\mathbf{v}}| \cos \theta$$

$\theta$  is  $\angle$  between  $\nabla J$  &  $\hat{\mathbf{v}}$   
 $\theta = 0 \Rightarrow \nabla J$  &  $\hat{\mathbf{v}}$  are  $\parallel$

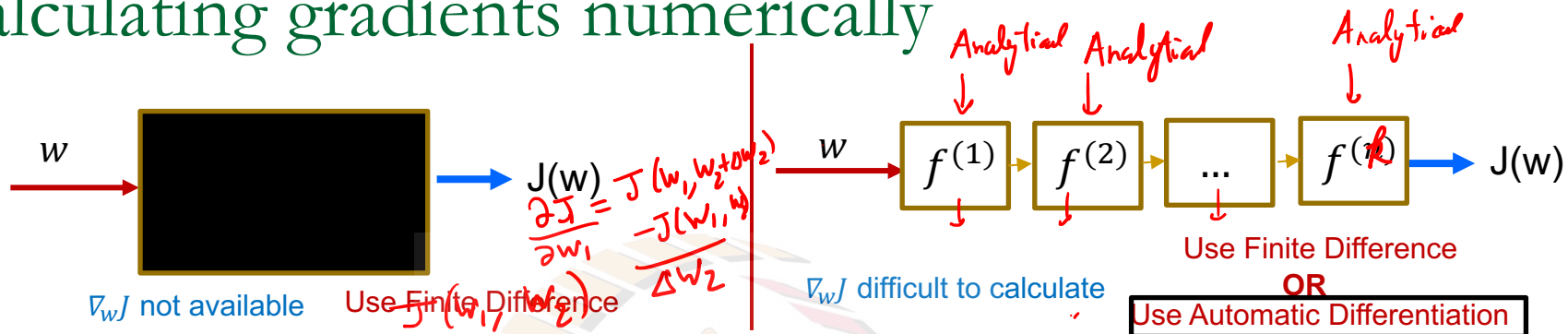
, where  $G = \nabla J$  and  $\theta$  is the angle between the gradient and  $v$

This is a maximum/minimum when  $\theta = 0$  and  $\pi$  respectively.

That is, the maximum increase is along the gradient and maximum decrease is opposite to the gradient.

**Proved**

# Calculating gradients numerically



- In most cases, we do not have explicit expressions for the gradient.
- This can usually happen because of two reasons
  - There is no analytical expression for  $J$  as it is available only as a black box
  - The expression for  $J$  is available as a composition of functions and is too complicated.
- In both cases, the simplest solution is to use the **Finite Difference Method**

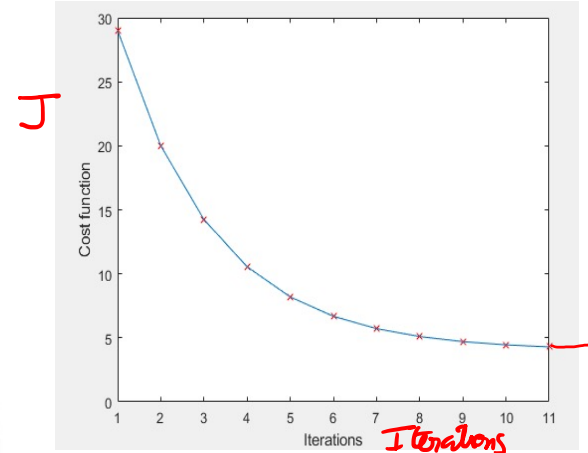
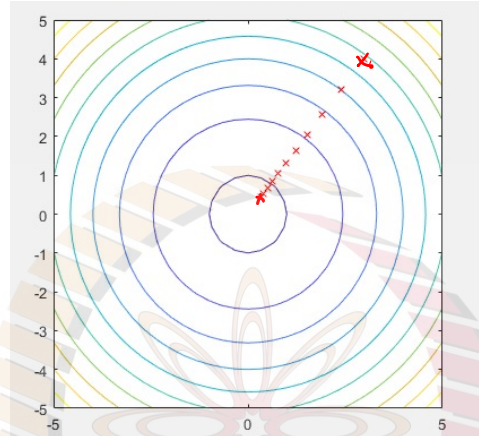
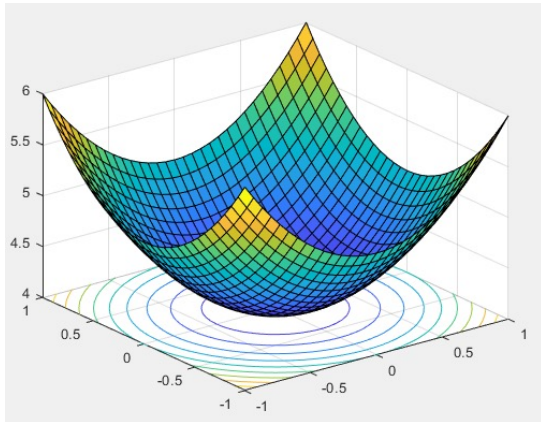
$$\frac{\partial J}{\partial w_j} \approx \frac{J(w_1, w_2, \dots, w_j + \delta w_j, \dots, w_n) - J(w_1, w_2, \dots, w_j - \delta w_j, \dots, w_n)}{2\delta w_j}$$

- This can, however, become very expensive if the number of features is large
- In case  $J$  is analytical but result of a chain  $J = f(g(h(\dots)))$  one can use **Automatic Differentiation**
  - This is a cheaper, numerical way to calculate the derivative via chain rule
  - For Neural Networks the equivalent method is called **Backpropagation**

# Stopping criteria

$$\vec{w}^{(3)} = \begin{bmatrix} 1.012 & 2.013 \end{bmatrix}$$

$$\vec{w}^{(4)} = \begin{bmatrix} 1.009 & 1.010 \end{bmatrix}$$



- Ideally, we should stop when  $\vec{\nabla}_{\vec{w}} J = \vec{0}$ . This almost never happens in practice as
  - The number of iterations required could be infinite
  - Because of finite precision
- In practice, we decide on some precision  $\epsilon$  (say,  $\epsilon \approx 10^{-5}$ ).
- Multiple options for stopping criteria. Stop when
  - $\|\vec{w}^{k+1} - \vec{w}^k\| \leq \epsilon$  ①
  - $\|\vec{\nabla}_{\vec{w}} J(\vec{w}^k)\| \leq \epsilon$  ②
  - $|J(\vec{w}^{k+1}) - J(\vec{w}^k)| \leq \epsilon$

# Summary of the Gradient Descent procedure

1. Decide on  $\alpha, \epsilon$  and stopping criterion

2. Make an initial guess for  $w = w^0$

3. Calculate  $w^{k+1} = w^k - \alpha \nabla_w J$

1. Calculate gradient numerically, if required

4. Calculate stopping criterion

1. If satisfied, stop

2. If not satisfied, go to Step 3