

# CS365: Deep Learning

## Deep Feedforward Network



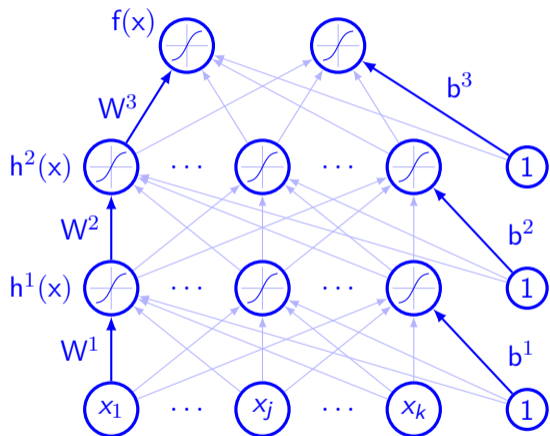
**Arijit Mondal**

Dept. of Computer Science & Engineering

Indian Institute of Technology Patna

arijit@iitp.ac.in

# Multilayer neural network



# Deep feedforward networks

- Also known as feedforward neural network or multilayer perceptron

# Deep feedforward networks

- Also known as feedforward neural network or multilayer perceptron
- Goal of such network is to approximate some function  $f^*$ 
  - For classifier,  $x$  is mapped to category  $y$  ie.  $y = f^*(x)$
  - A feedforward network maps  $y = f(x; \theta)$  and learns  $\theta$  for which the result is the best function approximation

# Deep feedforward networks

- Also known as feedforward neural network or multilayer perceptron
- Goal of such network is to approximate some function  $f^*$ 
  - For classifier,  $x$  is mapped to category  $y$  ie.  $y = f^*(x)$
  - A feedforward network maps  $y = f(x; \theta)$  and learns  $\theta$  for which the result is the best function approximation
- Information flows from input to intermediate to output
  - No feedback, directed acyclic graph
  - For general model, it can have feedback and known as recurrent neural network

# Deep feedforward networks

- Also known as feedforward neural network or multilayer perceptron
- Goal of such network is to approximate some function  $f^*$ 
  - For classifier,  $x$  is mapped to category  $y$  ie.  $y = f^*(x)$
  - A feedforward network maps  $y = f(x; \theta)$  and learns  $\theta$  for which the result is the best function approximation
- Information flows from input to intermediate to output
  - No feedback, directed acyclic graph
  - For general model, it can have feedback and known as recurrent neural network
- Typically it represents composition of functions
  - Three functions  $f^{(1)}, f^{(2)}, f^{(3)}$  are connected in chain
  - Overall function realized is  $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$
  - The number of layers provides the depth of the model

# Deep feedforward networks

- Also known as feedforward neural network or multilayer perceptron
- Goal of such network is to approximate some function  $f^*$ 
  - For classifier,  $x$  is mapped to category  $y$  ie.  $y = f^*(x)$
  - A feedforward network maps  $y = f(x; \theta)$  and learns  $\theta$  for which the result is the best function approximation
- Information flows from input to intermediate to output
  - No feedback, directed acyclic graph
  - For general model, it can have feedback and known as recurrent neural network
- Typically it represents composition of functions
  - Three functions  $f^{(1)}, f^{(2)}, f^{(3)}$  are connected in chain
  - Overall function realized is  $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$
  - The number of layers provides the depth of the model
- Goal of NN is not to model brain accurately!

# Issues with linear FFN

- Fit well for linear and logistic regression
- Convex optimization technique may be used
- Capacity of such function is limited
- Model cannot understand interaction between any two variables



# Overcome issues of linear FFN

- Transform  $x$  (input) into  $\phi(x)$  where  $\phi$  is nonlinear transformation

# Overcome issues of linear FFN

- Transform  $x$  (input) into  $\phi(x)$  where  $\phi$  is nonlinear transformation
- How to choose  $\phi$ ?

# Overcome issues of linear FFN

- Transform  $x$  (input) into  $\phi(x)$  where  $\phi$  is nonlinear transformation
- How to choose  $\phi$ ?
  - Use a very generic  $\phi$  of high dimension
    - Enough capacity but may result in poor generalization
    - Very generic feature mapping usually based on principle of local smoothness
    - Do not encode enough prior information

# Overcome issues of linear FFN

- Transform  $x$  (input) into  $\phi(x)$  where  $\phi$  is nonlinear transformation
- How to choose  $\phi$ ?
  - Use a very generic  $\phi$  of high dimension
    - Enough capacity but may result in poor generalization
    - Very generic feature mapping usually based on principle of local smoothness
    - Do not encode enough prior information
  - Manually design  $\phi$ 
    - Require domain knowledge

# Overcome issues of linear FFN

- Transform  $x$  (input) into  $\phi(x)$  where  $\phi$  is nonlinear transformation
- How to choose  $\phi$ ?
  - Use a very generic  $\phi$  of high dimension
    - Enough capacity but may result in poor generalization
    - Very generic feature mapping usually based on principle of local smoothness
    - Do not encode enough prior information
  - Manually design  $\phi$ 
    - Require domain knowledge
  - Strategy of deep learning is to learn  $\phi$

# Goal of deep learning

- We have a model  $y = f(x; \theta, w) = \phi(x; \theta)^T w$
- We use  $\theta$  to learn  $\phi$
- $w$  and  $\phi$  determines the output.  $\phi$  defines the hidden layer
- It loses the convexity of the training problem but benefits a lot
- Representation is parameterized as  $\phi(x, \theta)$ 
  - $\theta$  can be determined by solving optimization problem
- Advantages
  - $\phi$  can be very generic
  - Human practitioner can encode their knowledge to designing  $\phi(x; \theta)$

# Example

- Let us choose XOR function
- Target function is  $y = f^*(x)$  and our model provides  $y = f(x; \theta)$
- Learning algorithm will choose the parameters  $\theta$  to make  $f$  close to  $f^*$

# Example

- Let us choose XOR function
- Target function is  $y = f^*(x)$  and our model provides  $y = f(x; \theta)$
- Learning algorithm will choose the parameters  $\theta$  to make  $f$  close to  $f^*$
- Target is to fit output for  $X = \{[0, 0]^T, [0, 1]^T, [1, 0]^T, [1, 1]^T\}$
- This can be treated as regression problem and MSE error can be chosen as loss function  
$$(J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; \theta))^2)$$
- We need to choose  $f(x; \theta)$  where  $\theta$  depends on  $w$  and  $b$
- Let us consider a linear model  $f(x; w, b) = x^T w + b$



# Example

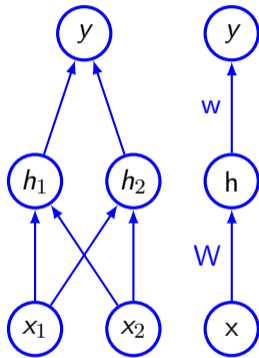
- Let us choose XOR function
- Target function is  $y = f^*(x)$  and our model provides  $y = f(x; \theta)$
- Learning algorithm will choose the parameters  $\theta$  to make  $f$  close to  $f^*$
- Target is to fit output for  $X = \{[0, 0]^T, [0, 1]^T, [1, 0]^T, [1, 1]^T\}$
- This can be treated as regression problem and MSE error can be chosen as loss function

$$(J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; \theta))^2)$$

- We need to choose  $f(x; \theta)$  where  $\theta$  depends on  $w$  and  $b$
- Let us consider a linear model  $f(x; w, b) = x^T w + b$
- Solving these, we get  $w = 0$  and  $b = \frac{1}{2}$

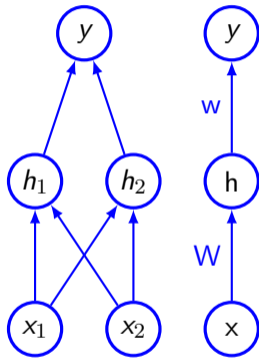
# Simple FFN with hidden layer

- Let us assume that the hidden unit  $h$  computes  $f^{(1)}(x; W, c)$



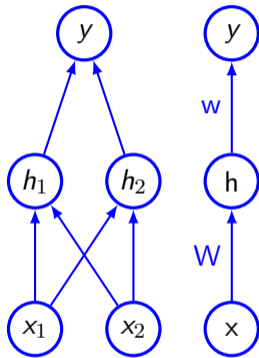
# Simple FFN with hidden layer

- Let us assume that the hidden unit  $h$  computes  $f^{(1)}(x; W, c)$
- In the next layer  $y = f^{(2)}(h; w, b)$  is computed



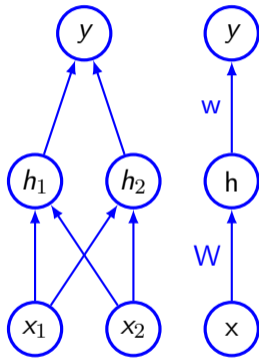
# Simple FFN with hidden layer

- Let us assume that the hidden unit  $h$  computes  $f^{(1)}(x; W, c)$
- In the next layer  $y = f^{(2)}(h; w, b)$  is computed
- Complete model  $f(x; W, c, w, b) = f^{(2)}(f^{(1)}(x))$



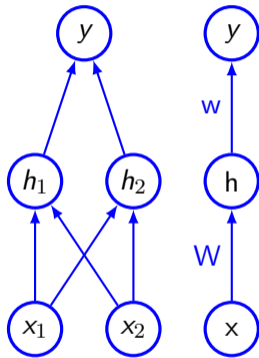
# Simple FFN with hidden layer

- Let us assume that the hidden unit  $h$  computes  $f^{(1)}(x; W, c)$
- In the next layer  $y = f^{(2)}(h; w, b)$  is computed
- Complete model  $f(x; W, c, w, b) = f^{(2)}(f^{(1)}(x))$
- Suppose  $f^{(1)}(x) = W^T x$  and  $f^2(h) = h^T w$



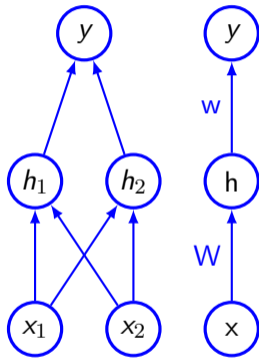
# Simple FFN with hidden layer

- Let us assume that the hidden unit  $h$  computes  $f^{(1)}(x; W, c)$
- In the next layer  $y = f^{(2)}(h; w, b)$  is computed
- Complete model  $f(x; W, c, w, b) = f^{(2)}(f^{(1)}(x))$
- Suppose  $f^{(1)}(x) = W^T x$  and  $f^{(2)}(h) = h^T w$  then  $f(x) = w^T W^T x$



# Simple FFN with hidden layer (contd.)

- We need to have nonlinear function to describe the features
- Usually NN have affine transformation of learned parameters followed by nonlinear activation function
- Let us use  $h = g(W^T x + c)$
- Let us use ReLU as activation function  $g(z) = \max\{0, z\}$
- $g$  is chosen element wise  $h_i = g(x^T W_{:,i} + c_i)$



# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$



# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$
- A solution for XOR problem can be as follows
  - $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- X

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW$

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ ,

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ , add bias  $c$

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ , add bias  $c \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ ,

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ , add bias  $c \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ , compute  $h$



# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ , add bias  $c \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ , compute  $h \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ ,

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$

- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ , add bias  $c \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ , compute  $h \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ , multiply

with  $w$

# Simple FFN with hidden layer (contd.)

- Complete network is  $f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$
- A solution for XOR problem can be as follows

- $W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$

- Now we have

- $X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$ , add bias  $c \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ , compute  $h \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ , multiply

with  $w \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$

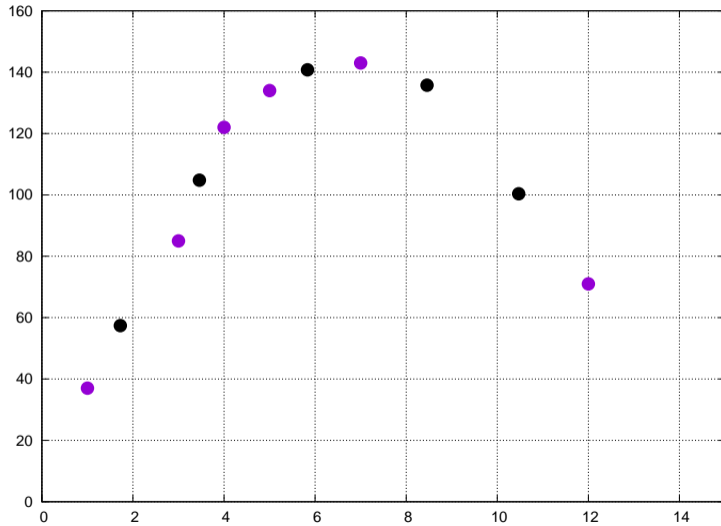
# Gradient based learning

- Similar to machine learning tasks, gradient descent based learning is used
  - Need to specify optimization procedure, cost function and model family
- For NN, model is nonlinear and function becomes nonconvex
  - Usually trained by iterative, gradient based optimizer
- Solved by using gradient descent or stochastic gradient descent (SGD)

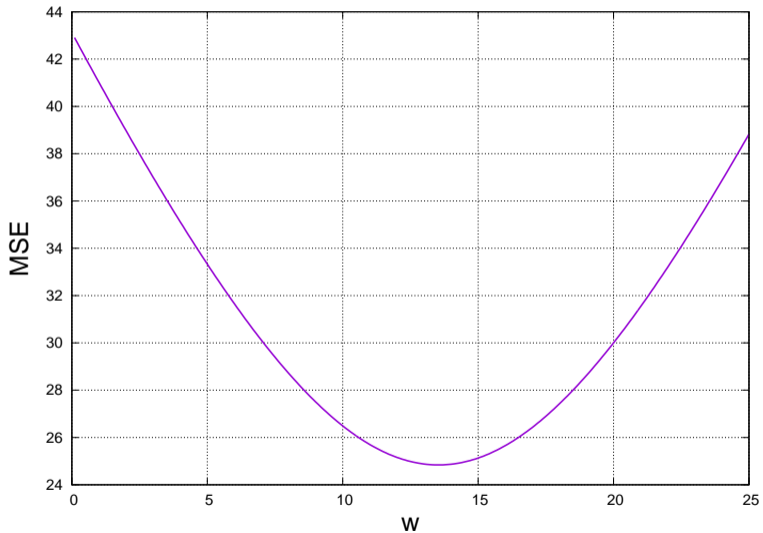
# Gradient descent

- For a function  $y = f(x)$ , derivative (slope at point  $x$ ) of it is  $f'(x) = \frac{dy}{dx}$
- A small change in the input can cause output to move to a value given by  $f(x+\epsilon) \approx f(x) + \epsilon f'(x)$
- We need to take a jump so that  $y$  reduces (assuming minimization problem)
- We can say that  $f(x - \epsilon \text{sign}(f'(x)))$  is less than  $f(x)$
- For multiple inputs partial derivatives are used ie.  $\frac{\partial}{\partial x_i} f(x)$
- Gradient vector is represented as  $\nabla_x f(x)$
- Gradient descent proposes a new point as  $x' = x - \epsilon \nabla_x f(x)$  where  $\epsilon$  is the learning rate

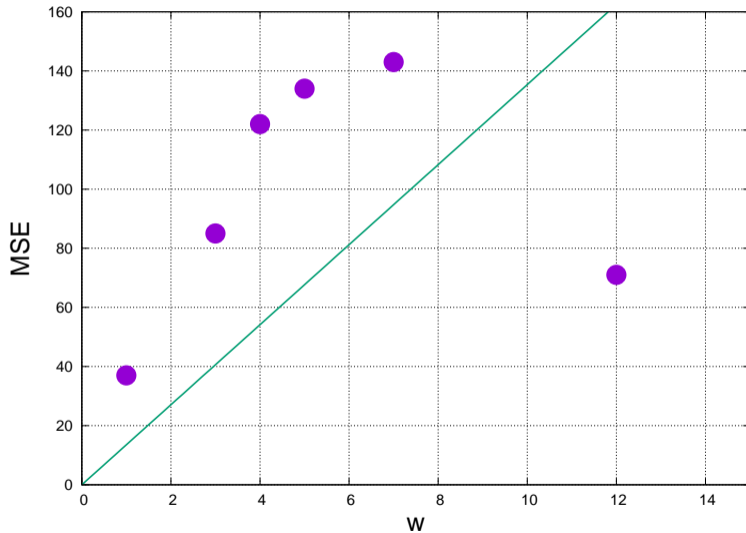
# Example



# Example: Variation of MSE wrt $w$



# Example: Best fit





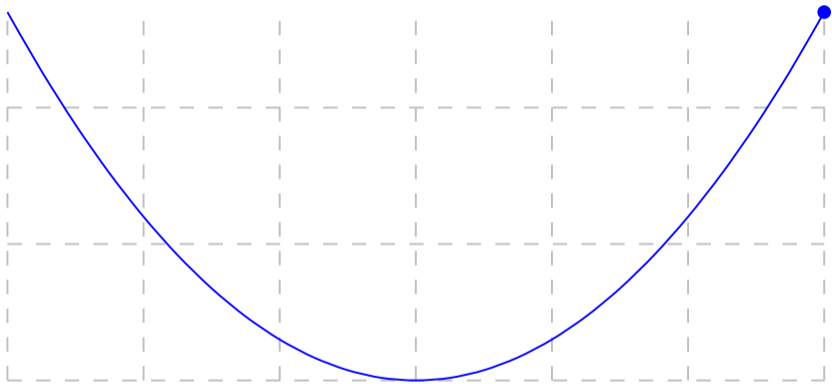
# Minimization of MSE: Gradient descent

- Assuming  $\text{MSE}_{(\text{train})} = J(w_1, w_2)$
- Target is to  $\min_{w_1, w_2} J(w_1, w_2)$
- Approach
  - Start with some  $w_1, w_2$
  - Keep modifying  $w_1, w_2$  so that  $J(w_1, w_2)$  reduces till the desired accuracy is achieved

# Minimization of MSE: Gradient descent

- Assuming  $\text{MSE}_{(\text{train})} = J(w_1, w_2)$
- Target is to  $\min_{w_1, w_2} J(w_1, w_2)$
- Approach
  - Start with some  $w_1, w_2$
  - Keep modifying  $w_1, w_2$  so that  $J(w_1, w_2)$  reduces till the desired accuracy is achieved
- Algorithm
  - Repeat the following until convergence  $w_j = w_j - \frac{\partial}{\partial w_j} J(w_1, w_2)$
- Gradient descent proposes a new point as  $w' = w - \epsilon \nabla_w f(w)$  where  $\epsilon$  is the learning rate

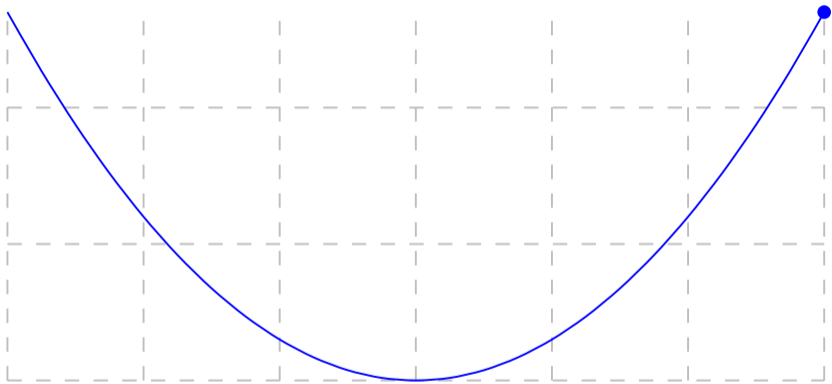
# Gradient descent



$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 1.80002$$

# Gradient descent

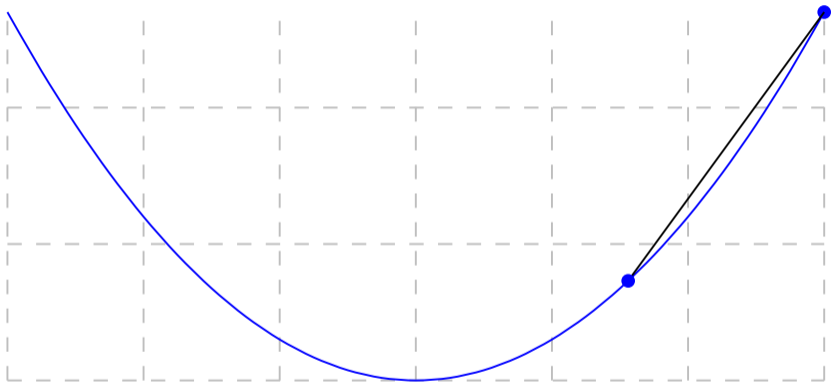


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 1.80002$$

$$x_{\text{new}} = 1.56001$$

# Gradient descent

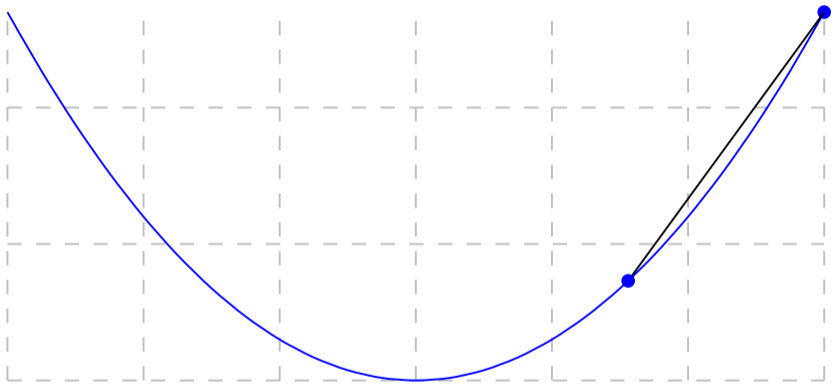


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 1.80002$$

$$x_{\text{new}} = 1.56001$$

# Gradient descent

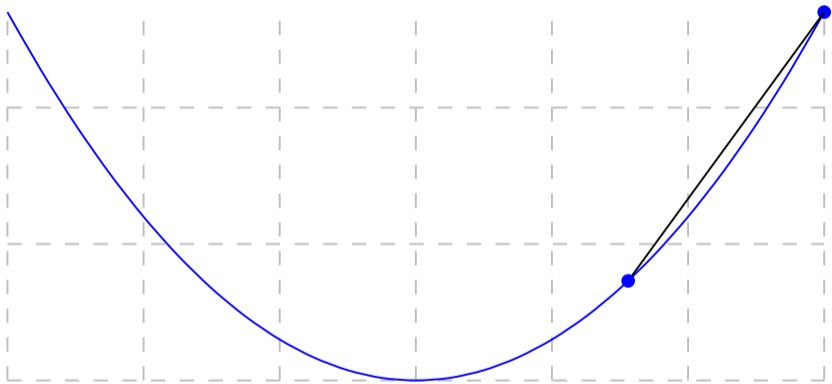


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.936$$

$$x_{\text{new}} = 1.56001$$

# Gradient descent

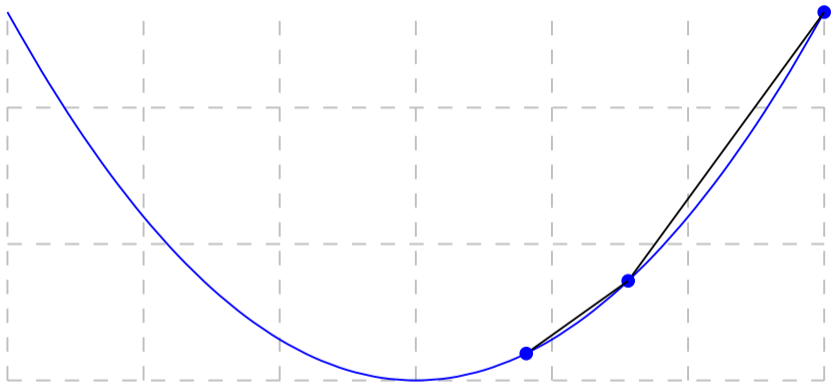


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.936$$

$$x_{\text{new}} = 0.81122$$

# Gradient descent



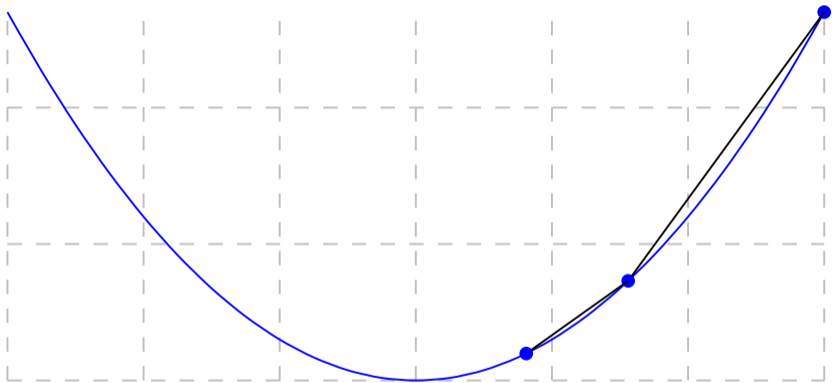
$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.936$$

$$x_{\text{new}} = 0.81122$$



# Gradient descent

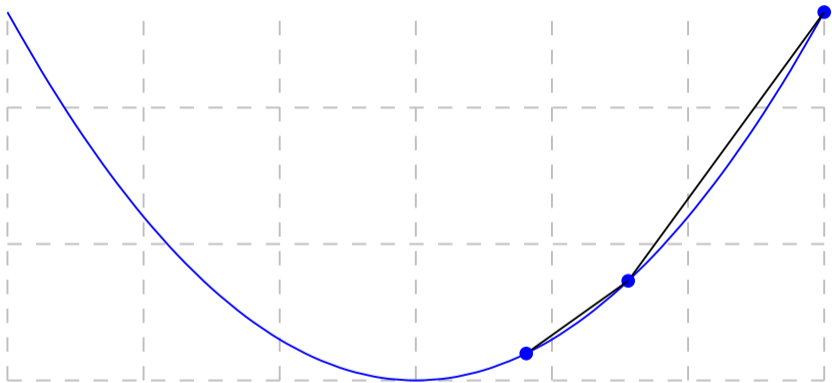


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.48672$$

$$x_{\text{new}} = 0.81122$$

# Gradient descent

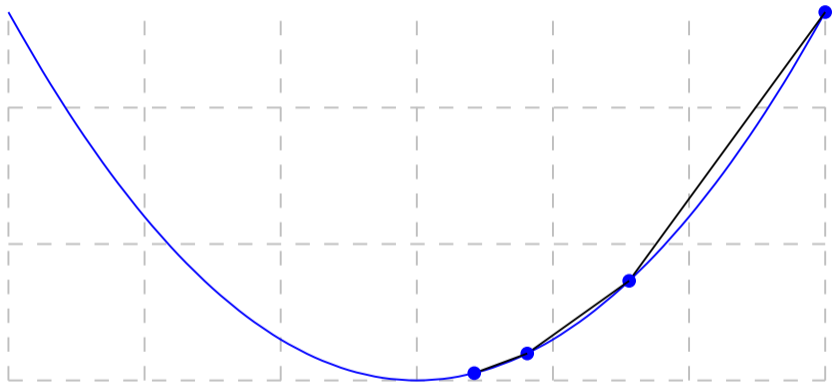


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.48672$$

$$x_{\text{new}} = 0.42184$$

# Gradient descent

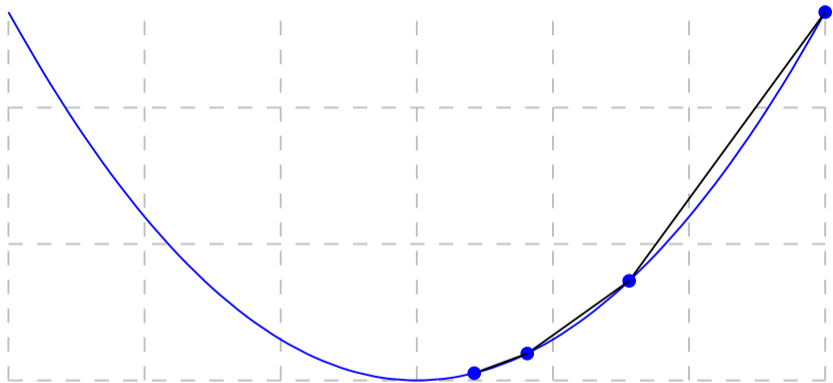


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.48672$$

$$x_{new} = 0.42184$$

# Gradient descent

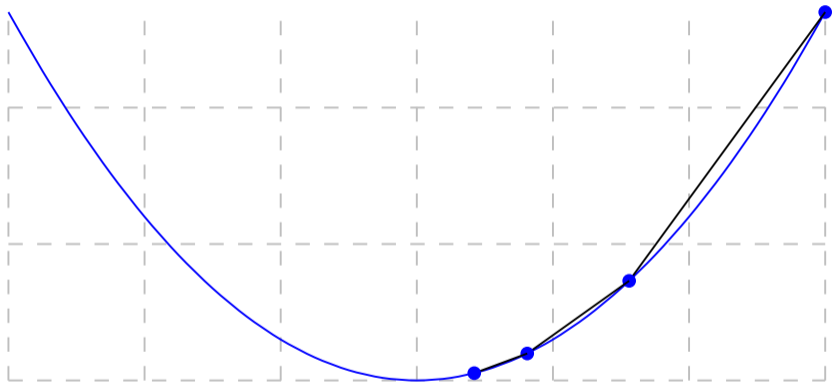


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.2531$$

$$x_{new} = 0.42184$$

# Gradient descent

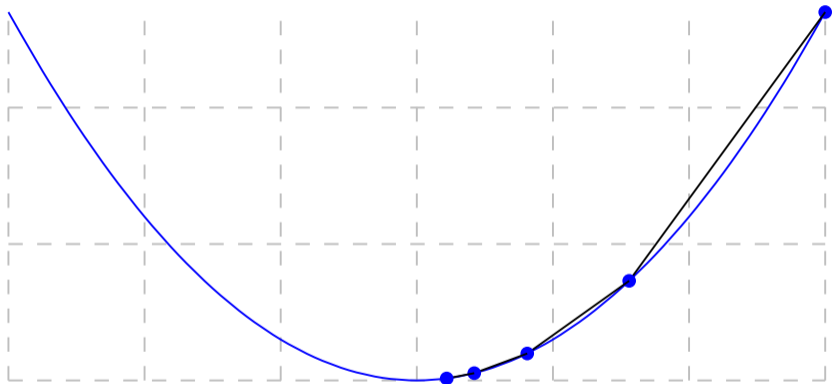


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.2531$$

$$x_{new} = 0.21938$$

# Gradient descent

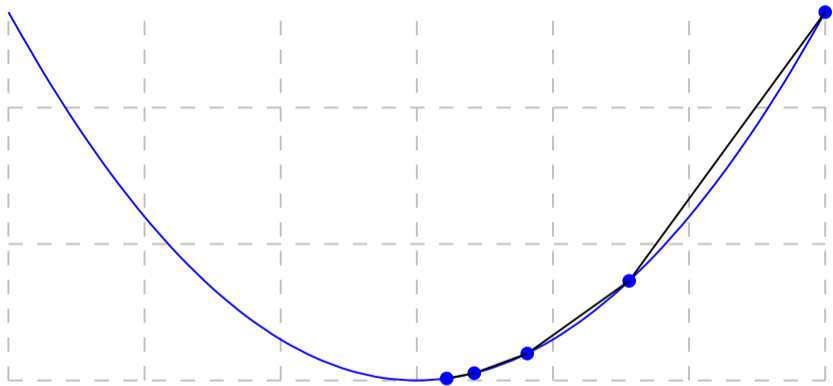


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.2531$$

$$x_{new} = 0.21938$$

# Gradient descent

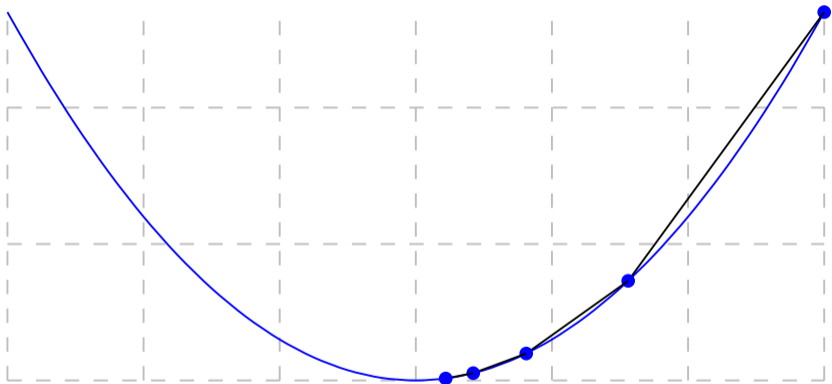


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.13162$$

$$x_{new} = 0.21938$$

# Gradient descent



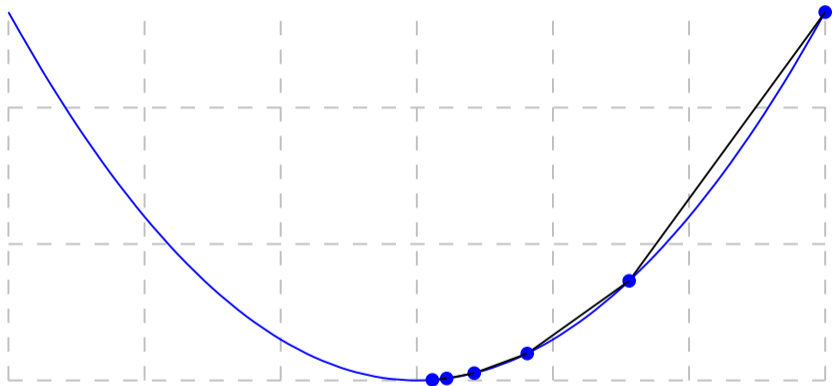
$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.13162$$

$$x_{new} = 0.11409$$



# Gradient descent

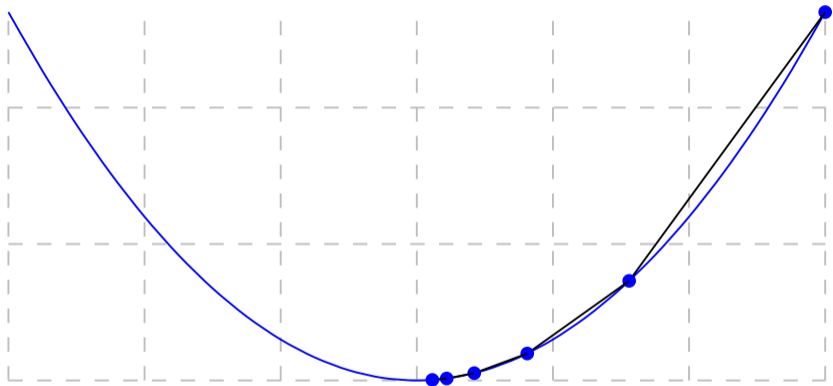


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.13162$$

$$x_{new} = 0.11409$$

# Gradient descent

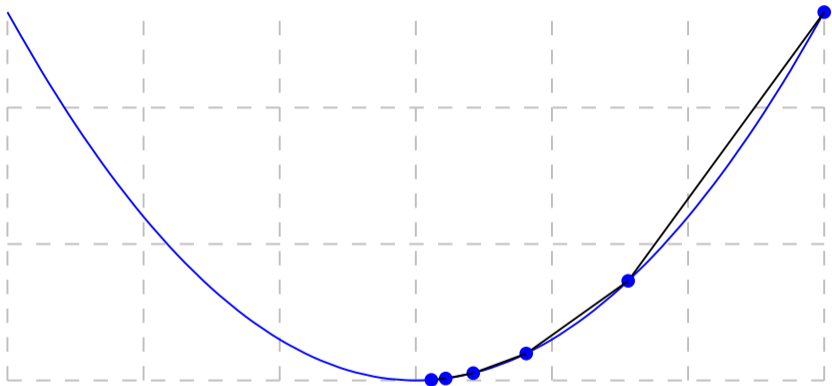


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.06845$$

$$x_{new} = 0.11409$$

# Gradient descent

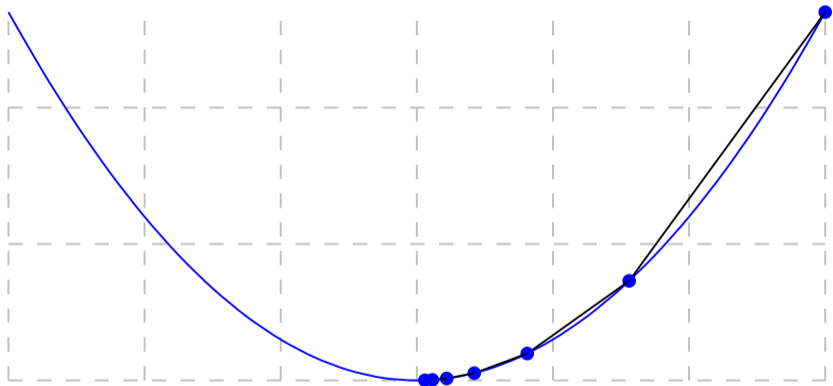


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.06845$$

$$x_{new} = 0.05934$$

# Gradient descent

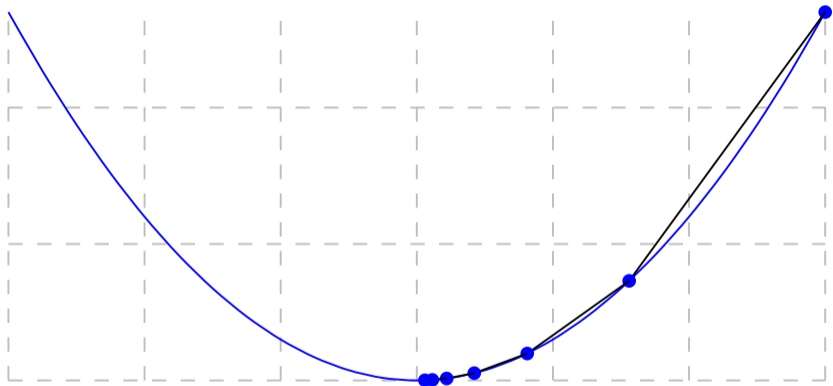


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.06845$$

$$x_{new} = 0.05934$$

# Gradient descent

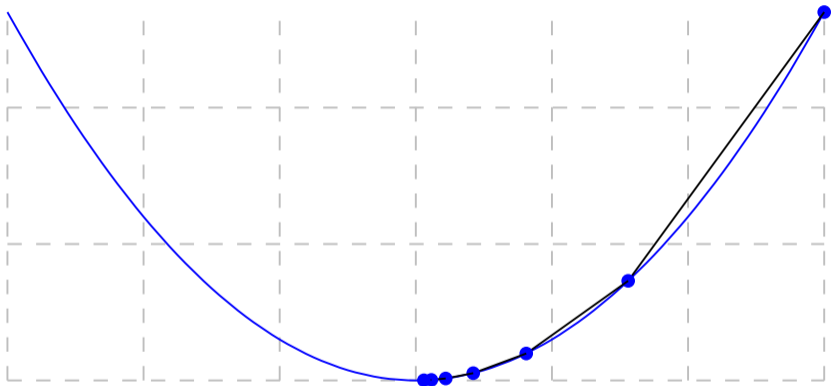


$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.0356$$

$$x_{new} = 0.05934$$

# Gradient descent



$$y = 0.3x^2, x_0 = 3, \alpha = 0.8$$

$$\text{gradient} = 0.0356$$

$$x_{new} = 0.03087$$

# Stochastic gradient descent

- Large training set are necessary for good generalization
- Cost function used for optimization is  $J(\theta) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}, \theta)$
- Gradient descent requires  $\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta)$

# Stochastic gradient descent

- Large training set are necessary for good generalization
- Cost function used for optimization is  $J(\theta) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}, \theta)$
- Gradient descent requires  $\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta)$ 
  - Computation cost is  $O(m)$



# Stochastic gradient descent

- Large training set are necessary for good generalization
- Cost function used for optimization is  $J(\theta) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}, \theta)$
- Gradient descent requires  $\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta)$ 
  - Computation cost is  $O(m)$
- For SGD, gradient is an expectation estimated from a small sample known as minibatch ( $\mathbb{B} = \{x^{(1)}, \dots, x^{(m')}\}$ )
- Estimated gradient is  $g = \frac{1}{m'} \sum_{i=1}^{m'} \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta)$
- New point will be  $\theta = \theta - \epsilon g$

# SGD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2)$ ,  $(2, 4)$ ,  $(3, 6)$ ,  $(4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $x(w \times x - y)$

Step	Point	Derivative	New w
------	-------	------------	-------

# SGD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2)$ ,  $(2, 4)$ ,  $(3, 6)$ ,  $(4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $x(w \times x - y)$

Step	Point	Derivative	New w
1	(1,2)	$1*(4.0*1-2)=2.0$	3.80

# SGD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2)$ ,  $(2, 4)$ ,  $(3, 6)$ ,  $(4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $x(w \times x - y)$

Step	Point	Derivative	New w
1	(1,2)	$1*(4.0*1-2)=2.0$	3.80
2	(2,4)	$2*(3.8*2-4)=7.2$	3.08

# SGD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2)$ ,  $(2, 4)$ ,  $(3, 6)$ ,  $(4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $x(w \times x - y)$

Step	Point	Derivative	New w
1	(1,2)	$1*(4.0*1-2)=2.0$	3.80
2	(2,4)	$2*(3.8*2-4)=7.2$	3.08
3	(3,6)	$3*(3.1*3-6)=9.7$	2.11
4	(4,8)	$4*(2.1*4-8)=1.7$	1.94
5	(1,2)	$1*(1.9*1-2)=-0.1$	1.94
6	(2,4)	$2*(1.9*2-4)=-0.2$	1.97
7	(3,6)	$3*(2.0*3-6)=-0.3$	1.99
8	(4,8)	$4*(2.0*4-8)=-0.1$	2.00
9	(1,2)	$1*(2.0*1-2)=0.0$	2.00

# GD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2)$ ,  $(2, 4)$ ,  $(3, 6)$ ,  $(4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $\frac{1}{4} \sum_i x_i (w \times x_i - y_i)$

Step	Derivative	New w
------	------------	-------

# GD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2), (2, 4), (3, 6), (4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $\frac{1}{4} \sum_i x_i (w \times x_i - y_i)$

Step	Derivative	New w
1	15	2.5

# GD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2), (2, 4), (3, 6), (4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $\frac{1}{4} \sum_i x_i (w \times x_i - y_i)$

Step	Derivative	New w
1	15	2.5
2	3.75	2.13



# GD example

- Consider the following pair  $(x, y)$  of points -  $(1, 2), (2, 4), (3, 6), (4, 8)$
- Let us try to fit a curve as follows  $y = w \times x$  where  $w$  is initialized with 4, learning rate as 0.1
- MSE as cost function. Derivative will be  $\frac{1}{4} \sum_i x_i (w \times x_i - y_i)$

Step	Derivative	New w
1	15	2.5
2	3.75	2.13
3	0.94	2.03
4	0.23	2.01
5	0.06	2.00

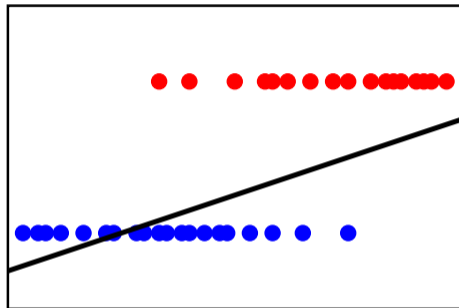
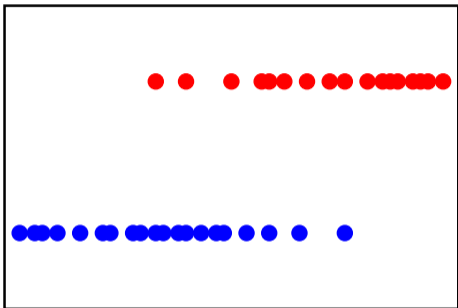
# Cost function

- Similar to other parametric model like linear models
- Parametric model defines distribution  $p(y|x; \theta)$
- Principle of maximum likelihood is used (cross entropy between training data and model prediction)
- Instead of predicting the whole distribution of  $y$ , some statistic of  $y$  conditioned on  $x$  is predicted
- It can also contain regularization term

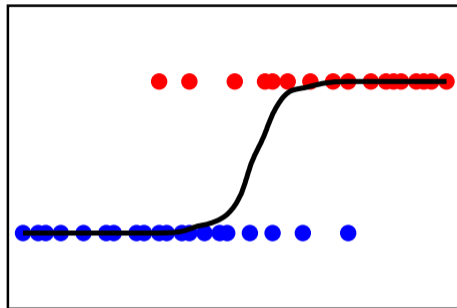
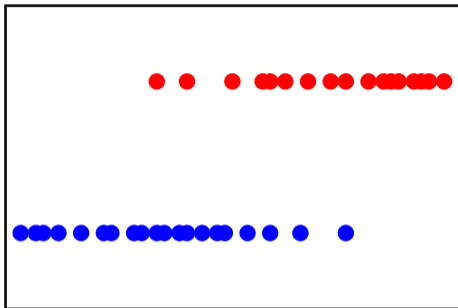
# Logistic regression

- Responses may be qualitative (categorical)
  - Example: ⟨Hours of study, pass/fail⟩, ⟨MRI scan, benign/malignant⟩
  - Output should be 0 or 1
- Predicting qualitative response is known as classification
- Linear regression does not help

# Issues with linear regression



# Logistic regression



# Logistic model

- Linear regression model to represent non-normalized probability  $p'(x) = w_0 + w_1x$
- To avoid problem, we use function  $p(x) = \frac{e^{w_0+w_1x}}{1 + e^{w_0+w_1x}}$
- Quantity  $\frac{p(x)}{1-p(x)} = e^{w_0+w_1x}$  is known as odds
- Taking log on both the sides, we get  $\log\left(\frac{p(x)}{1-p(x)}\right) = w_0 + w_1x$
- Coefficient can be determined using maximum likelihood
  - $l(w_0, w_1) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} p(x_j)$
- Similar to linear regression except the output is mapped between 0 and 1 ie.

$$p(y|x, \theta) = \sigma(\theta^T x)$$

where  $\sigma(x) = \frac{1}{1 + \exp(-x)}$  (Sigmoid function)

# Maximum likelihood estimation

- Consider a set of  $m$  examples  $\mathbb{X} = \{x^{(1)}, \dots, x^{(m)}\}$  drawn independently from the true but unknown data generating distribution  $p_{data}(x)$
- Let  $p_{model}(x; \theta)$  be a parametric family of probability distribution

# Maximum likelihood estimation

- Consider a set of  $m$  examples  $\mathbb{X} = \{x^{(1)}, \dots, x^{(m)}\}$  drawn independently from the true but unknown data generating distribution  $p_{data}(x)$
- Let  $p_{model}(x; \theta)$  be a parametric family of probability distribution
- Maximum likelihood estimator for  $\theta$  is defined as

$$\theta_{ML} = \arg \max_{\theta} p_{model}(\mathbb{X}; \theta) = \arg \max_{\theta} \prod_{i=1}^m p_{model}(x^{(i)}; \theta)$$



# Maximum likelihood estimation

- Consider a set of  $m$  examples  $\mathbb{X} = \{x^{(1)}, \dots, x^{(m)}\}$  drawn independently from the true but unknown data generating distribution  $p_{data}(x)$
- Let  $p_{model}(x; \theta)$  be a parametric family of probability distribution
- Maximum likelihood estimator for  $\theta$  is defined as

$$\theta_{ML} = \arg \max_{\theta} p_{model}(\mathbb{X}; \theta) = \arg \max_{\theta} \prod_{i=1}^m p_{model}(x^{(i)}; \theta)$$

- It can be written as  $\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^m \log p_{model}(x^{(i)}; \theta)$

# Maximum likelihood estimation

- Consider a set of  $m$  examples  $\mathbb{X} = \{x^{(1)}, \dots, x^{(m)}\}$  drawn independently from the true but unknown data generating distribution  $p_{data}(x)$
- Let  $p_{model}(x; \theta)$  be a parametric family of probability distribution
- Maximum likelihood estimator for  $\theta$  is defined as

$$\theta_{ML} = \arg \max_{\theta} p_{model}(\mathbb{X}; \theta) = \arg \max_{\theta} \prod_{i=1}^m p_{model}(x^{(i)}; \theta)$$

- It can be written as  $\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^m \log p_{model}(x^{(i)}; \theta)$
- By dividing  $m$  we get  $\theta_{ML} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{data}} \log p_{model}(x; \theta)$

# Maximum likelihood estimation (cont.)

- Minimizing dissimilarity between the empirical  $\hat{p}_{data}$  and model distribution  $p_{model}$  and it is measured by KL divergence

$$D_{KL}(\hat{p}_{data} || p_{model}) = \arg \min_{\theta} \mathbb{E}_{x \sim \hat{p}_{data}} [\log \hat{p}_{data}(x) - \log p_{model}(x; \theta)]$$

# Maximum likelihood estimation (cont.)

- Minimizing dissimilarity between the empirical  $\hat{p}_{data}$  and model distribution  $p_{model}$  and it is measured by KL divergence

$$D_{KL}(\hat{p}_{data} || p_{model}) = \arg \min_{\theta} \mathbb{E}_{x \sim \hat{p}_{data}} [\log \hat{p}_{data}(x) - \log p_{model}(x; \theta)]$$

- We need to minimize  $-\arg \min_{\theta} \mathbb{E}_{x \sim \hat{p}_{data}} \log p_{model}(x; \theta)$

# Conditional log-likelihood

- In most of the supervised learning we estimate  $P(y|x; \theta)$
- If  $X$  be the all inputs and  $Y$  be observed targets then conditional maximum likelihood estimator is  $\theta_{ML} = \arg \max_{\theta} P(Y|X; \theta)$
- If the examples are assumed to be i.i.d then we can say

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^m \log P(y^{(i)}|x^{(i)}; \theta)$$

# Linear regression as maximum likelihood

- Instead of producing single prediction  $\hat{y}$  for a given  $x$ , we assume the model produces conditional distribution  $p(y|x)$
- For infinitely large training set, we can observe multiple examples having the same  $x$  but different values of  $y$
- Goal is to fit the distribution  $p(y|x)$

# Linear regression as maximum likelihood

- Instead of producing single prediction  $\hat{y}$  for a given  $\mathbf{x}$ , we assume the model produces conditional distribution  $p(y|\mathbf{x})$
- For infinitely large training set, we can observe multiple examples having the same  $\mathbf{x}$  but different values of  $y$
- Goal is to fit the distribution  $p(y|\mathbf{x})$
- Let us assume,  $p(y|\mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$

# Linear regression as maximum likelihood

- Instead of producing single prediction  $\hat{y}$  for a given  $\mathbf{x}$ , we assume the model produces conditional distribution  $p(y|\mathbf{x})$
- For infinitely large training set, we can observe multiple examples having the same  $\mathbf{x}$  but different values of  $y$
- Goal is to fit the distribution  $p(y|\mathbf{x})$
- Let us assume,  $p(y|\mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$
- Since the examples are assumed to be i.i.d, conditional log-likelihood is given by

$$\sum_{i=1}^m \log p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$$



# Linear regression as maximum likelihood

- Instead of producing single prediction  $\hat{y}$  for a given  $\mathbf{x}$ , we assume the model produces conditional distribution  $p(y|\mathbf{x})$
- For infinitely large training set, we can observe multiple examples having the same  $\mathbf{x}$  but different values of  $y$
- Goal is to fit the distribution  $p(y|\mathbf{x})$
- Let us assume,  $p(y|\mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$
- Since the examples are assumed to be i.i.d, conditional log-likelihood is given by

$$\sum_{i=1}^m \log p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta}) = -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|\hat{\mathbf{y}}^{(i)} - y^{(i)}\|^2}{2\sigma^2}$$

# Learning conditional distributions

- Usually neural networks are trained using maximum likelihood. Therefore the cost function is negative log-likelihood. Also known as cross entropy between training data and model distribution
- Cost function  $J(\theta) = -\mathbb{E}_{X, Y \sim \hat{p}_{data}} \log p_{model}(y|x, \theta)$
- Uniform across different models
- Gradient of cost function is very much crucial
  - Large and predictable gradient can serve good guide for learning process
  - Function that saturates will have small gradient
    - Activation function usually produces values in a bounded zone (saturates)
  - Negative log-likelihood can overcome some of the problems
    - Output unit having exp function can saturate for high negative value
    - Log-likelihood cost function undoes the exp of some output functions

# Learning conditional statistics

- Instead of learning the whole distribution  $p(y|x; \theta)$ , we want to learn one conditional statistics of  $y$  given  $x$ 
  - For a predicting function  $f(x; \theta)$ , we would like to predict the mean of  $y$

# Learning conditional statistics

- Instead of learning the whole distribution  $p(y|x; \theta)$ , we want to learn one conditional statistics of  $y$  given  $x$ 
  - For a predicting function  $f(x; \theta)$ , we would like to predict the mean of  $y$
- Neural network can represent any function  $f$  from a very wide range of functions
- Range of function is limited by features like continuity, boundedness, etc.

# Learning conditional statistics

- Instead of learning the whole distribution  $p(y|x; \theta)$ , we want to learn one conditional statistics of  $y$  given  $x$ 
  - For a predicting function  $f(x; \theta)$ , we would like to predict the mean of  $y$
- Neural network can represent any function  $f$  from a very wide range of functions
- Range of function is limited by features like continuity, boundedness, etc.
- Cost function becomes functional rather than a function

# Learning conditional statistics

- Need to solve the optimization problem

$$f^* = \arg \min_f \mathbb{E}_{X, Y \sim p_{data}} \|y - f(x)\|^2$$

# Learning conditional statistics

- Need to solve the optimization problem

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim p_{data}} \|\mathbf{y} - f(\mathbf{x})\|^2$$

- Using calculus of variation, it gives  $f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{Y} \sim p_{data}(\mathbf{y}|\mathbf{x})}[\mathbf{y}]$ 
  - Mean of  $\mathbf{y}$  for each value of  $\mathbf{x}$

# Learning conditional statistics

- Need to solve the optimization problem

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim p_{data}} \|\mathbf{y} - f(\mathbf{x})\|^2$$

- Using calculus of variation, it gives  $f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{Y} \sim p_{data}(\mathbf{y}|\mathbf{x})}[\mathbf{y}]$

- Mean of  $\mathbf{y}$  for each value of  $\mathbf{x}$

- Using a different cost function  $f^* = \arg \min_f \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim p_{data}} \|\mathbf{y} - f(\mathbf{x})\|_1$



# Learning conditional statistics

- Need to solve the optimization problem

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim p_{data}} \|\mathbf{y} - f(\mathbf{x})\|^2$$

- Using calculus of variation, it gives  $f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{Y} \sim p_{data}(\mathbf{y}|\mathbf{x})}[\mathbf{y}]$

- Mean of  $\mathbf{y}$  for each value of  $\mathbf{x}$

- Using a different cost function  $f^* = \arg \min_f \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim p_{data}} \|\mathbf{y} - f(\mathbf{x})\|_1$

- Median of  $\mathbf{y}$  for each value of  $\mathbf{x}$

*Thank you!*

# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$

# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$
- Let  $J[y]$  has local minima at  $f$ . Therefore, we can say  $J[f] \leq J[f + \varepsilon\eta]$ 
  - $\eta$  is an arbitrary function of  $x$  such that  $\eta(x_1) = \eta(x_2) = 0$  and differentiable

# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$
- Let  $J[y]$  has local minima at  $f$ . Therefore, we can say  $J[f] \leq J[f + \varepsilon\eta]$ 
  - $\eta$  is an arbitrary function of  $x$  such that  $\eta(x_1) = \eta(x_2) = 0$  and differentiable
- Let us assume  $\Phi(\varepsilon) = J[f + \varepsilon\eta]$ . Therefore,  $\Phi'(0) \equiv \left. \frac{d\Phi}{d\varepsilon} \right|_{\varepsilon=0} = \int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx$

# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$
- Let  $J[y]$  has local minima at  $f$ . Therefore, we can say  $J[f] \leq J[f + \varepsilon\eta]$ 
  - $\eta$  is an arbitrary function of  $x$  such that  $\eta(x_1) = \eta(x_2) = 0$  and differentiable
- Let us assume  $\Phi(\varepsilon) = J[f + \varepsilon\eta]$ . Therefore,  $\Phi'(0) \equiv \left. \frac{d\Phi}{d\varepsilon} \right|_{\varepsilon=0} = \int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = 0$

# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$
- Let  $J[y]$  has local minima at  $f$ . Therefore, we can say  $J[f] \leq J[f + \varepsilon\eta]$ 
  - $\eta$  is an arbitrary function of  $x$  such that  $\eta(x_1) = \eta(x_2) = 0$  and differentiable
- Let us assume  $\Phi(\varepsilon) = J[f + \varepsilon\eta]$ . Therefore,  $\Phi'(0) \equiv \left. \frac{d\Phi}{d\varepsilon} \right|_{\varepsilon=0} = \int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = 0$
- Now we can say,  $\frac{dL}{d\varepsilon} = \frac{\partial L}{\partial y} \frac{dy}{d\varepsilon} + \frac{\partial L}{\partial y'} \frac{dy'}{d\varepsilon}$

# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$
- Let  $J[y]$  has local minima at  $f$ . Therefore, we can say  $J[f] \leq J[f + \varepsilon\eta]$ 
  - $\eta$  is an arbitrary function of  $x$  such that  $\eta(x_1) = \eta(x_2) = 0$  and differentiable
- Let us assume  $\Phi(\varepsilon) = J[f + \varepsilon\eta]$ . Therefore,  $\Phi'(0) \equiv \left. \frac{d\Phi}{d\varepsilon} \right|_{\varepsilon=0} = \int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = 0$
- Now we can say,  $\frac{dL}{d\varepsilon} = \frac{\partial L}{\partial y} \frac{dy}{d\varepsilon} + \frac{\partial L}{\partial y'} \frac{dy'}{d\varepsilon}$
- As we have  $y = f + \varepsilon\eta$  and  $y' = f' + \varepsilon\eta'$ , therefore,  $\frac{dL}{d\varepsilon}$



# Calculus of variations

- Let us consider functional  $J[y] = \int_{x_1}^{x_2} L(x, y(x), y'(x)) dx$
- Let  $J[y]$  has local minima at  $f$ . Therefore, we can say  $J[f] \leq J[f + \varepsilon\eta]$ 
  - $\eta$  is an arbitrary function of  $x$  such that  $\eta(x_1) = \eta(x_2) = 0$  and differentiable
- Let us assume  $\Phi(\varepsilon) = J[f + \varepsilon\eta]$ . Therefore,  $\Phi'(0) \equiv \left. \frac{d\Phi}{d\varepsilon} \right|_{\varepsilon=0} = \int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = 0$
- Now we can say,  $\frac{dL}{d\varepsilon} = \frac{\partial L}{\partial y} \frac{dy}{d\varepsilon} + \frac{\partial L}{\partial y'} \frac{dy'}{d\varepsilon}$
- As we have  $y = f + \varepsilon\eta$  and  $y' = f' + \varepsilon\eta'$ , therefore,  $\frac{dL}{d\varepsilon} = \frac{\partial L}{\partial y}\eta + \frac{\partial L}{\partial y'}\eta'$

# Calculus of variations (contd.)

- Now we have

$$\int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta + \frac{\partial L}{\partial f'} \eta' \right) dx$$

# Calculus of variations (contd.)

- Now we have

$$\int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta + \frac{\partial L}{\partial f'} \eta' \right) dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta - \eta \frac{d}{dx} \frac{\partial L}{\partial f'} \right) dx + \left. \frac{\partial L}{\partial f'} \eta \right|_{x_1}^{x_2}$$

# Calculus of variations (contd.)

- Now we have

$$\int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta + \frac{\partial L}{\partial f'} \eta' \right) dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta - \eta \frac{d}{dx} \frac{\partial L}{\partial f'} \right) dx + \left. \frac{\partial L}{\partial f'} \eta \right|_{x_1}^{x_2}$$

- Hence  $\int_{x_1}^{x_2} \eta \left( \frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} \right) dx = 0$

# Calculus of variations (contd.)

- Now we have

$$\int_{x_1}^{x_2} \left. \frac{dL}{d\varepsilon} \right|_{\varepsilon=0} dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta + \frac{\partial L}{\partial f'} \eta' \right) dx = \int_{x_1}^{x_2} \left( \frac{\partial L}{\partial f} \eta - \eta \frac{d}{dx} \frac{\partial L}{\partial f'} \right) dx + \left. \frac{\partial L}{\partial f'} \eta \right|_{x_1}^{x_2}$$

- Hence  $\int_{x_1}^{x_2} \eta \left( \frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} \right) dx = 0$

- Euler-Lagrange equation  $\frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} = 0$

# Example

- Let us consider distance between two points  $A[y] = \int_{x_1}^{x_2} \sqrt{1 + [y'(x)]^2} dx$ 
  - $y'(x) = \frac{dy}{dx}$ ,  $y_1 = f(x_1)$ ,  $y_2 = f(x_2)$

# Example

- Let us consider distance between two points  $A[y] = \int_{x_1}^{x_2} \sqrt{1 + [y'(x)]^2} dx$ 
  - $y'(x) = \frac{dy}{dx}$ ,  $y_1 = f(x_1)$ ,  $y_2 = f(x_2)$
- We have,  $\frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} = 0$  where  $L = \sqrt{1 + [f'(x)]^2}$

# Example

- Let us consider distance between two points  $A[y] = \int_{x_1}^{x_2} \sqrt{1 + [y'(x)]^2} dx$ 
  - $y'(x) = \frac{dy}{dx}$ ,  $y_1 = f(x_1)$ ,  $y_2 = f(x_2)$
- We have,  $\frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} = 0$  where  $L = \sqrt{1 + [f'(x)]^2}$
- As  $f$  does not appear explicitly in  $L$ , hence  $\frac{d}{dx} \frac{\partial L}{\partial f'} = 0$



# Example

- Let us consider distance between two points  $A[y] = \int_{x_1}^{x_2} \sqrt{1 + [y'(x)]^2} dx$ 
  - $y'(x) = \frac{dy}{dx}$ ,  $y_1 = f(x_1)$ ,  $y_2 = f(x_2)$
- We have,  $\frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} = 0$  where  $L = \sqrt{1 + [f'(x)]^2}$
- As  $f$  does not appear explicitly in  $L$ , hence  $\frac{d}{dx} \frac{\partial L}{\partial f'} = 0$
- Now we have,  $\frac{d}{dx} \frac{f'(x)}{\sqrt{1 + [f'(x)]^2}} = 0$

# Example

- Taking derivative we get  $\frac{d^2 f}{dx^2} \cdot \frac{1}{\left[\sqrt{1 + [f'(x)]^2}\right]^3} = 0$

# Example

- Taking derivative we get  $\frac{d^2 f}{dx^2} \cdot \frac{1}{\left[\sqrt{1 + [f'(x)]^2}\right]^3} = 0$
- Therefore we have,  $\frac{d^2 f}{dx^2} = 0$

# Example

- Taking derivative we get  $\frac{d^2 f}{dx^2} \cdot \frac{1}{\left[\sqrt{1 + [f'(x)]^2}\right]^3} = 0$
- Therefore we have,  $\frac{d^2 f}{dx^2} = 0$
- Hence we have  $f(x) = mx + b$  with  $m = \frac{y_2 - y_1}{x_2 - x_1}$  and  $b = \frac{x_2 y_1 - x_1 y_2}{x_2 - x_1}$