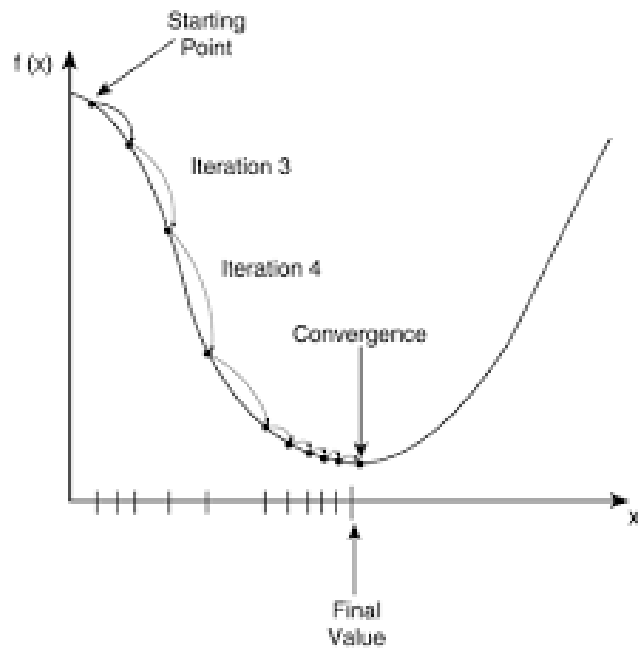


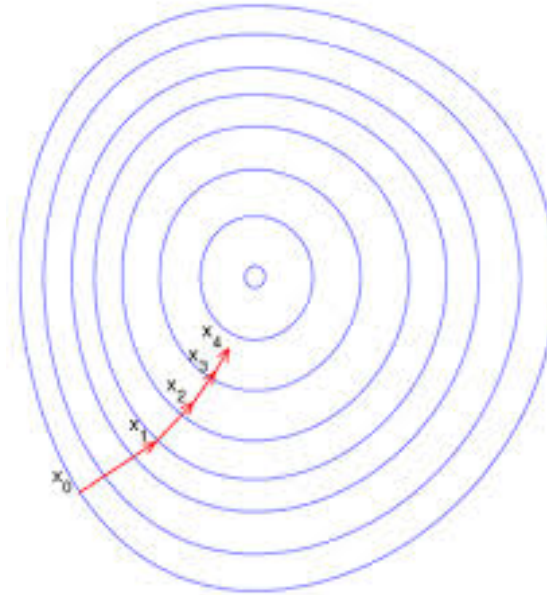
Gradient Descent

- ❑ Most simple method for unconstrained optimization
- ❑ Key property of gradient, $\nabla_w f(\mathbf{w})$
 - $-\nabla_w f(\mathbf{w})$ = Points in the direction of steepest decrease
- ❑ Gradient descent algorithm:
 - Start with initial w^0
 - $w^{k+1} = w^k - \alpha_k \nabla f(w^k)$
 - Repeat until some stopping criteria
- ❑ α_k is called the **step size**
 - In machine learning, this is called the **learning rate**

Gradient Descent Illustrated



□ $M = 1$



• $M = 2$

Gradient Descent Analysis

□ Using gradient update rule

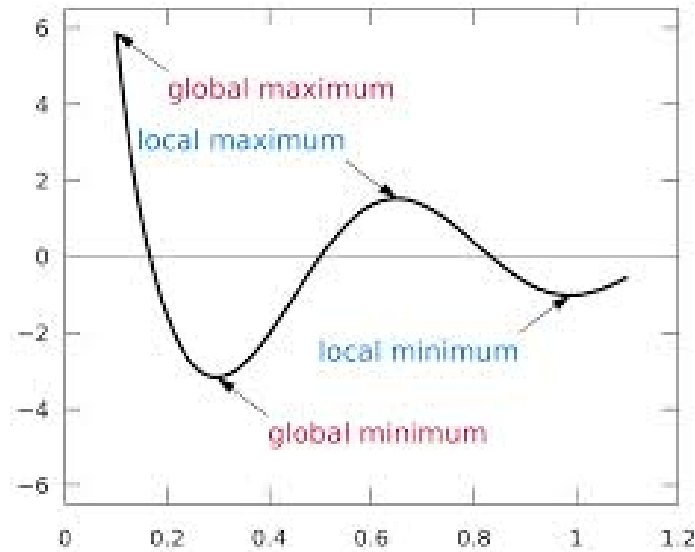
$$\begin{aligned} f(w^{k+1}) &= f(w^k) + \nabla f(w^k) \cdot (w^{k+1} - w^k) + O\|w^{k+1} - w^k\|^2 \\ &= f(w^k) - \alpha \nabla f(w^k) \cdot \nabla f(w^k) + O(\alpha^2) \\ &= f(w^k) - \alpha \|\nabla f(w^k)\|^2 + O(\alpha^2) \end{aligned}$$

□ Consequence: If step size α is small, then $f(w^k)$ decreases

□ Theorem:

If $f''(w)$ is bounded above, $f(w)$ is bounded below, and α is chosen sufficiently small,
Then gradient descent converges to **local** minima

Local vs. Global Minima



□ Definitions:

- w^* is a **global minima** if $f(w) \geq f(w^*)$ for all w
- w^* is a **local minima** if $f(w) \geq f(w^*)$ for all w in some open neighborhood of w^*

□ Most numerical methods:

- Generally only guarantee convergence to **local minima**

□ Convex functions: Have only global minima (more later)

Gradients for Logistic Regression

□ Logistic regression

- Linear function: $z_i = w_0 + \sum_{j=1}^d X_{ij}w_j$
- Output probability: $P(y = 1|x) = \frac{1}{1+e^{-z_i}}$
- Binary cross-entropy loss: $J(\mathbf{w}) = \sum_{i=1}^n \{\ln[1 + e^{z_i}] - y_i z_i\}$

□ Compute gradients:

- Define $A = [1 \ X]$, matrix with ones on the first column
- Then, $z_i = w_0 + \sum_{j=1}^d X_{ij}w_j = \sum_{j=0}^d A_{ij}w_j$
- Let $p_i = \frac{1}{1+e^{-z_i}}$
- Observe $\frac{\partial J}{\partial z_i} = \frac{e^{z_i}}{1+e^{z_i}} - y_i = p_i - y_i$
- Use multi-variable chain rule:

$$\frac{\partial J}{\partial w_j} = \sum_{i=1}^N \frac{\partial J}{\partial z_i} \frac{\partial z_i}{\partial w_j} = \sum_{i=1}^N (p_i - y_i) A_{ij}$$

Matrix Form

□ Logistic regression

- Linear function: $z_i = \sum_{j=0}^d A_{ij} w_j$
- Output probability: $P(y = 1|x) = \frac{1}{1+e^{-z_i}}$
- BCE: $J = \sum_{i=1}^n \{\ln[1 + e^{z_i}] - y_i z_i\}$
- $\frac{\partial J}{\partial z_i} = p_i - y_i$
- $\frac{\partial J}{\partial w_j} = \sum_{i=1}^N \frac{\partial J}{\partial z_i} \frac{\partial z_i}{\partial w_j} = \sum_{i=1}^N (p_i - y_i) A_{ij}$

□ Matrix form:

- $z = Aw$
- Let $p = \frac{1}{1+e^{-z}}$
- $\frac{\partial J}{\partial z} = p - y$
- $\frac{\partial J}{\partial w} = A^T \frac{\partial J}{\partial z}$

```
def feval(w,X,y):  
    """  
    Compute the loss and gradient given w,X,y  
    """  
  
    # Construct transform matrix  
    n = X.shape[0]  
    A = np.column_stack((np.ones(n,), X))  
  
    # The loss is the binary cross entropy  
    z = A.dot(w)  
    py = 1/(1+np.exp(-z))  
    f = np.sum((1-y)*z - np.log(py))  
  
    # Gradient  
    df_dz = py-y  
    fgrad = A.T.dot(df_dz)  
    return f, fgrad
```

Implementation in Python

❑ Optimizer requires a python method to compute:

- Objective function $f(\mathbf{w})$, and
- Gradient $\nabla f(\mathbf{w})$

❑ For logistic loss:

$$f(\mathbf{w}) := \sum_{i=1}^N -y_i z_i + \ln[1 + e^{z_i}], \quad z = A\mathbf{w}$$

❑ Thus, $f(\mathbf{w})$ and $\nabla f(\mathbf{w})$ depends on training data (\mathbf{x}_i, y_i)

- How do we pass these?

❑ Two methods to pass data to the function:

- Method 1: Use a class
- Method 2: Use lambda calculus

Training data

```
def feval(w,X,y):  
    """  
    Compute the loss and gradient given w,X,y  
    """  
    # Construct transform matrix  
    n = X.shape[0]  
    A = np.column_stack((np.ones(n,), X))  
  
    # The loss is the binary cross entropy  
    z = A.dot(w)  
    py = 1/(1+np.exp(-z))  
    f = np.sum((1-y)*z - np.log(py))  
  
    # Gradient  
    df_dz = py-y  
    fgrad = A.T.dot(df_dz)  
    return f, fgrad
```