

# Iterative Reweighted Least Squares

Sargur N. Srihari

University at Buffalo, State University of New York  
USA

# Topics in Linear Classification using Probabilistic Discriminative Models

- Generative vs Discriminative
  1. Fixed basis functions in linear classification
  2. Logistic Regression (two-class)
  3. Iterative Reweighted Least Squares (IRLS)
  4. Multiclass Logistic Regression
  5. Probit Regression
  6. Canonical Link Functions

# Topics in IRLS

- What is IRLS
- Linear and Logistic Regression
- IRLS for Linear Regression
- IRLS for Logistic Regression

# What is IRLS?

- An iterative method to find solution  $w^*$ 
  - for linear regression and logistic regression
    - assuming least squares objective
- While simple gradient descent has the form

$$w^{(\text{new})} = w^{(\text{old})} - \eta \nabla E(w)$$

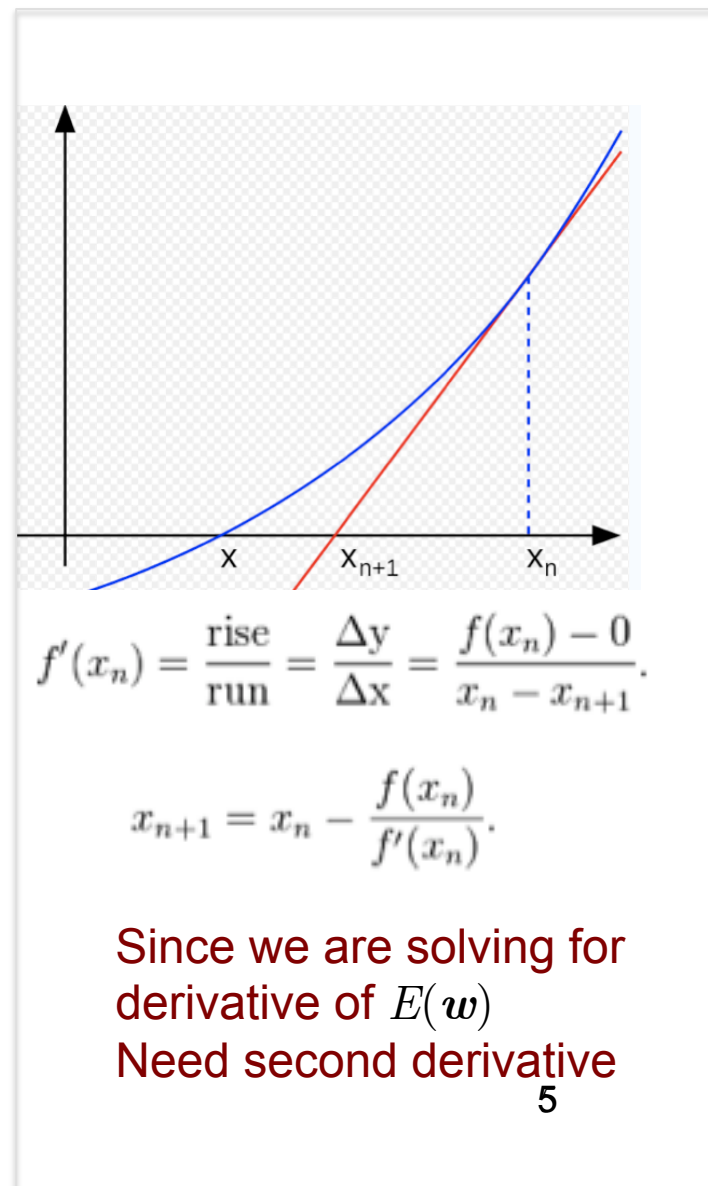
- IRLS uses the second derivative and has the form

$$w^{(\text{new})} = w^{(\text{old})} - H^{-1} \nabla E(w)$$

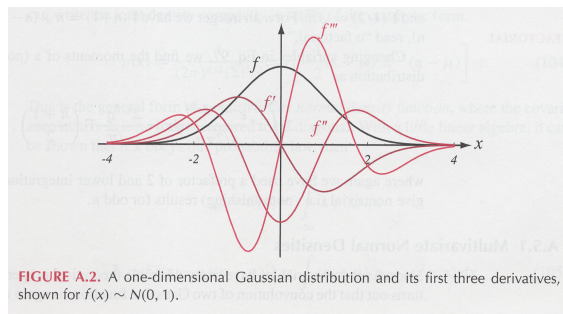
- It is derived from Newton-Raphson method
  - where  $H$  is the Hessian matrix whose elements are the second derivatives of  $E(w)$  wrt  $w$

# Newton-Raphson Method (1-D)

- Based on second derivatives
  - Derivative at point  $x$  of a function is the slope of its tangent at that point
- Illustration of second derivative
  - Derivatives of Gaussian  $p(x) \sim N(0, \sigma)$



$$\begin{aligned} \frac{\partial}{\partial x} \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{-x}{\sqrt{2\pi}\sigma^3} e^{-x^2/(2\sigma^2)} = \frac{-x}{\sigma^2} p(x) \\ \frac{\partial^2}{\partial x^2} \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{1}{\sqrt{2\pi}\sigma^5} (-\sigma^2 + x^2) e^{-x^2/(2\sigma^2)} = \frac{-\sigma^2 + x^2}{\sigma^4} p(x) \\ \frac{\partial^3}{\partial x^3} \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{1}{\sqrt{2\pi}\sigma^7} (3x\sigma^2 - x^3) e^{-x^2/(2\sigma^2)} = \frac{-3x\sigma^2 - x^3}{\sigma^6} p(x), \end{aligned}$$

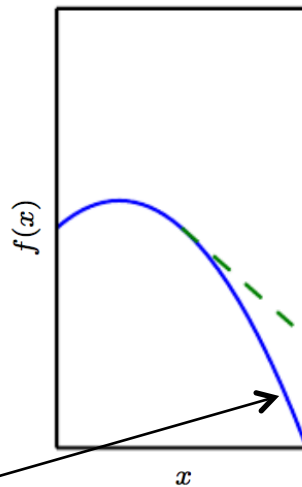


# Second derivative measures curvature

- Derivative of a derivative
- Quadratic functions with different curvatures

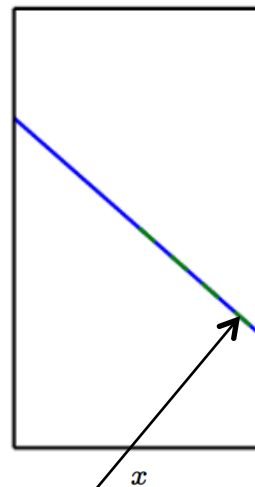
Dashed line is value of cost function predicted by gradient alone

Negative curvature



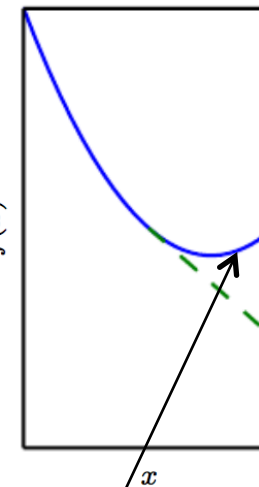
Decrease is faster than predicted by Gradient Descent

No curvature

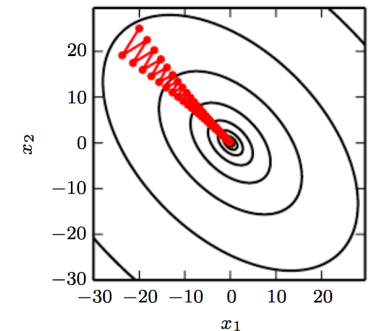


Gradient Predicts decrease correctly

Positive curvature



Decrease is slower than expected  
Actually increases



# Learning rate from Hessian

- Taylor's series of  $f(\mathbf{x})$  around current point  $\mathbf{x}^{(0)}$

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^T \mathbf{g} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(0)})^T H(\mathbf{x} - \mathbf{x}^{(0)})$$

- where  $\mathbf{g}$  is the gradient and  $H$  is the Hessian at  $\mathbf{x}^{(0)}$

– If we use learning rate  $\varepsilon$  the new point  $\mathbf{x}$  is given by  $\mathbf{x}^{(0)} - \varepsilon \mathbf{g}$ . Thus we get

$$f(\mathbf{x}^{(0)} - \varepsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \varepsilon \mathbf{g}^T \mathbf{g} + \frac{1}{2} \varepsilon^2 \mathbf{g}^T H \mathbf{g}$$

- There are three terms:
  - original value of  $f$ ,
  - expected improvement due to slope, and
  - correction to be applied due to curvature
- Solving for step size when correction is least gives

$$\varepsilon^* \approx \frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T H \mathbf{g}}$$

# Another 2<sup>nd</sup> Derivative Method

- Using Taylor's series of  $f(\mathbf{x})$  around current  $\mathbf{x}^{(0)}$

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^T \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^T H(f)(\mathbf{x} - \mathbf{x}^{(0)}) (\mathbf{x} - \mathbf{x}^{(0)})$$

- solve for the critical point of this function to give

$$\mathbf{x}^* = \mathbf{x}^{(0)} - H(f)(\mathbf{x}^{(0)})^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$$

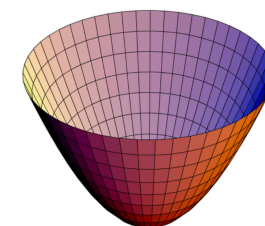
- When  $f$  is a quadratic (positive definite) function use solution to jump to the minimum function directly
- When not quadratic apply solution iteratively
- Can reach critical point much faster than gradient descent
  - But useful only when nearby point is a minimum



# Linear and Logistic Regression

- In linear regression there is a closed-form max likelihood solution for determining  $\mathbf{w}$ 
  - on the assumption of Gaussian noise model
  - Due to quadratic dependence of log-likelihood on  $\mathbf{w}$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2$$



- For logistic regression: No closed-form maximum likelihood solution
  - Due to nonlinearity of logistic sigmoid

$$E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

$$y_n = \sigma(\mathbf{w}^T \phi_n)$$

- But departure from quadratic is not substantial
  - Error function is concave, i.e., unique minimum<sup>9</sup>

# Two applications of IRLS

- IRLS is applicable to both Linear Regression and Logistic Regression
- We discuss both, for each we need

1. Model

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

$$p(C_1 | \boldsymbol{\phi}) = y(\boldsymbol{\phi}) = \sigma(\mathbf{w}^T \boldsymbol{\phi})$$

2. Objective function  $E(\mathbf{w})$

- Linear Regression: Sum of Squared Errors
- Logistic Regression: Bernoulli Likelihood Function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

3. Gradient

$$\nabla E(\mathbf{w})$$

4. Hessian

$$H = \nabla \nabla E(\mathbf{w})$$

5. Newton-Raphson update

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$$

# IRLS for Linear Regression

1. Model:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

2. Error Function: Sum of Squares:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \right\}^2$$

for data set  $\mathbf{X} = \{ \mathbf{x}_n, t_n \} \quad n=1, \dots, N$

3. Gradient of Error Function is:

$$\begin{aligned} \nabla E(\mathbf{w}) &= \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}_n - t_n) \boldsymbol{\phi}_n \\ &= \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{t} \end{aligned}$$

where  $\boldsymbol{\Phi}$  is the  $N \times M$  design matrix whose  $n^{\text{th}}$  row is given by  $\boldsymbol{\phi}_n^T$

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

4. Hessian is:

$$H = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

5. Newton-Raphson Update:

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$$

Substituting:  $H = \boldsymbol{\Phi}^T \boldsymbol{\Phi}$  and  $\nabla E(\mathbf{w}) = \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{t}$

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \{ \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w}^{(\text{old})} - \boldsymbol{\Phi}^T \mathbf{t} \} \\ &= (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t} \end{aligned}$$

which is the standard least squares solution

Since it is independent of  $\mathbf{w}$ , Newton-Raphson gives exact solution in one step

# IRLS for Logistic Regression

1. **Model:**  $p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$  **Likelihood:**

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

for data set  $\{\phi_n, t_n\}$ ,  $t_n \in \{0, 1\}$ ,  $y_n = \phi(\mathbf{x}_n)$

2. **Objective Function:** **Negative log-likelihood:**

$$E(\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

Cross-entropy error function

3. **Gradient of Error Function is:**

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t})$$

where  $\Phi$  is the  $N \times M$  design matrix whose  $n^{\text{th}}$  row is given by  $\phi_n^T$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

4. **Hessian is:**

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi$$

$\mathbf{R}$  is  $N \times N$  diagonal matrix with elements

$$R_{nn} = y_n(1 - y_n) = \mathbf{w}^T \phi_n (1 - \mathbf{w}^T \phi_n)$$

Hessian is not constant and depends on  $\mathbf{w}$  through  $\mathbf{R}$ . Since  $\mathbf{H}$  is positive-definite (i.e., for arbitrary  $\mathbf{u}$ ,  $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$ ) error function is a concave function of  $\mathbf{w}$  and so has a unique minimum

5. **Newton-Raphson Update:**

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

**Substituting**  $\mathbf{H} = \Phi^T \mathbf{R} \Phi$  **and**  $\nabla E(\mathbf{w}) = \Phi^T (\mathbf{y} - \mathbf{t})$

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi \Phi^T \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} \end{aligned}$$

where  $\mathbf{z}$  is a  $N$ -dimensional vector with elements  $\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$

Update formula is a set of normal equations.

Since Hessian depends on  $\mathbf{w}$  apply them iteratively each time using the new weight vector