

# Foundations of Machine Learning

## AI2000 and AI5000

FoML-21

Logistic Regression - Newton Raphson optimization

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad

July-Nov 2025



भारतीय नॉर्केटिक विज्ञान संस्था  
भारतीय प्रौद्योगिकी संस्थान हैदराबाद  
Indian Institute of Technology Hyderabad

**DIL**

Data-driven Intelligence  
& Learning Lab

# So far in FoML

- Intro to ML and Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Supervised learning
  - a. Linear Regression with basis functions (regularization, model selection)
  - b. Bias-Variance Decomposition (Bayesian Regression)
  - c. Decision Theory - three broad classification strategies
    - Probabilistic Generative Models - Continuous & discrete data
    - (Linear) Discriminant Functions - least squares solution, Perceptron
    - Probabilistic Discriminative Models - Logistic Regression



# Logistic Regression - IRLS



ભારતીય નોંકેટિક વિજ્ઞાન સંસ્કૃત પ્રેરણાખાડ  
ભારતીય પ્રૌદ્યોગિકી સંસ્થાન હૈદરાબાદ  
Indian Institute of Technology Hyderabad



# Loss function approximation

- Taylor series expansion (univariate case)

$$E(w + \Delta w) = E(w) + \Delta w \cdot E'(w) + \frac{1}{2} \Delta w^2 E''(w) + \dots$$



# Loss function approximation

- Taylor series expansion (multivariate case)

$$E(\mathbf{w} + \Delta\mathbf{w}) = E(\underline{\mathbf{w}}) + \underline{\Delta\mathbf{w}}^T \cdot \nabla E + \frac{1}{2} \underline{\Delta\mathbf{w}}^T \cdot \underline{H} \cdot \underline{\Delta\mathbf{w}} + \dots$$

$$\underline{H}_{ij} = \frac{\partial^2 E(\underline{\mathbf{w}})}{\partial w_i \partial w_j}$$

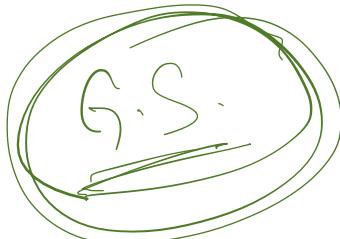
Hessian of  $E$  at  $\underline{\mathbf{w}}$



# Linear (first-order) approximation

$$E(\mathbf{w} + \Delta\mathbf{w}) \approx E(\mathbf{w}) + \Delta\mathbf{w}^T \nabla E(\mathbf{w})$$

for  $\Delta\mathbf{w} \ll E(\mathbf{w})$



$$\Rightarrow \Delta\mathbf{w} = -\eta \nabla E(\mathbf{w})$$

i.e., take a small step  
in the opposite  
direction of gradient



# Quadratic (second-order) approximation

$$E(\mathbf{w} + \Delta\mathbf{w}) \approx E(\mathbf{w}) + \Delta\mathbf{w}^T \nabla E(\mathbf{w}) + \frac{1}{2} \Delta\mathbf{w}^T \nabla^2 E(\mathbf{w}) \Delta\mathbf{w}$$

what is the best  $\Delta\mathbf{w}$  for minimizing  $E(\mathbf{w})$

Hessian

$$\Rightarrow \frac{\partial E(\mathbf{w} + \Delta\mathbf{w})}{\partial \Delta\mathbf{w}} = 0$$

$\Rightarrow$

$$\nabla E(\mathbf{w}) + \nabla^2 E(\mathbf{w}) \cdot \Delta\mathbf{w} = 0$$

$$\mathbf{w}^{k+1} = \mathbf{w}^k - H^{-1} \nabla E(\mathbf{w})$$

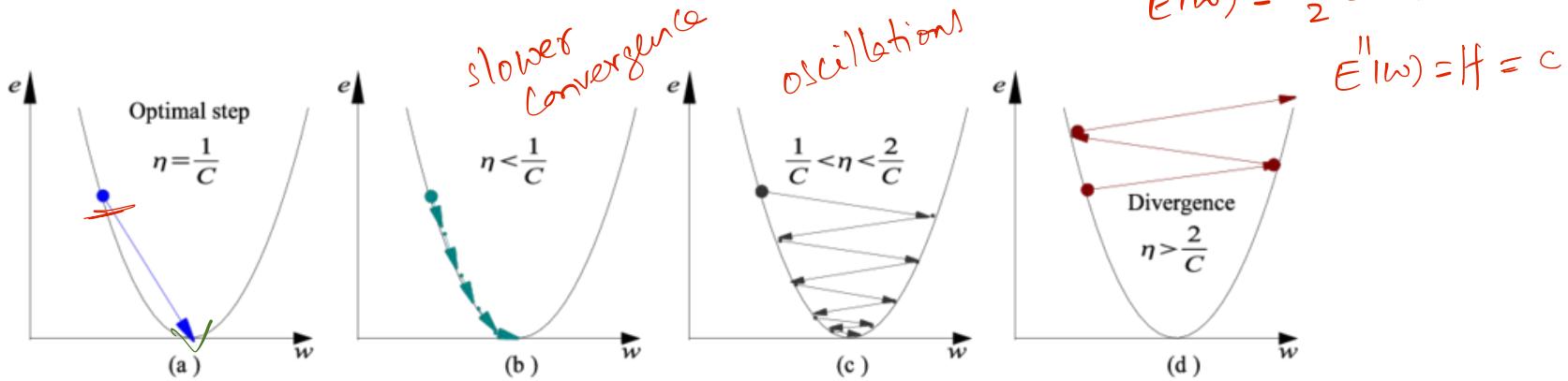
$\Leftarrow$

$$\Delta\mathbf{w} = - \left[ \nabla^2 E(\mathbf{w}) \right]^{-1} \nabla E(\mathbf{w})$$

Best  $\Delta\mathbf{w}$  to minimize  
 $E(\mathbf{w})$



# Convergence for Quadratic functions



quadratic loss function is convex

$\Rightarrow$  2nd order approximation will be exact for such functions  
 $\Rightarrow$  single step convergence (irrespective of  $w^\circ$ )

$$\hat{w} = w^\circ - \hat{H}^{-1} \nabla E(w)$$



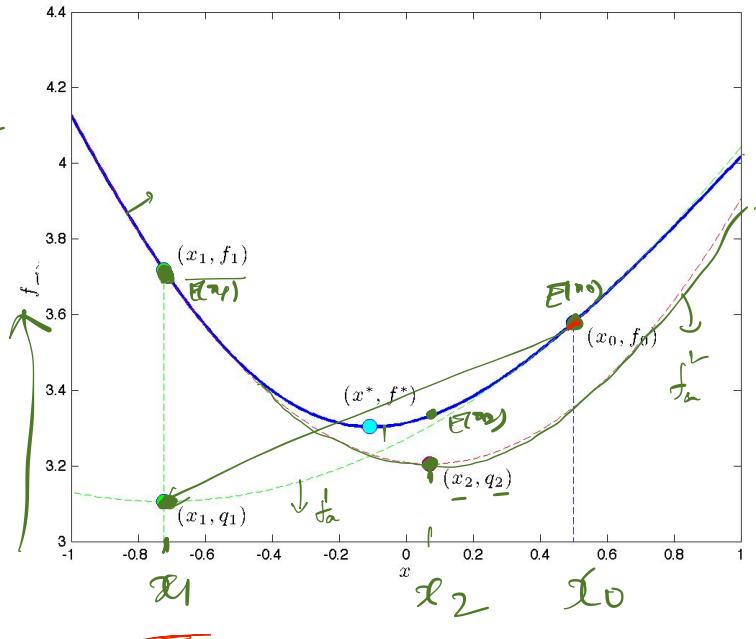
optimizing

# Generic Convex functions

with

Second-order  
approx.

example for  
generic (degree > 2) but  
convex loss function



$\rightarrow x_0 = \text{initial guess}$

$$\underline{x}_1 = x_0 - \frac{1}{f'(x_0)} f''(x_0)$$

$$\underline{x}_2 = x_1 - \frac{1}{f'(x_1)} f''(x_1)$$



भारतीय नॉर्केटिक विज्ञान संस्था पूर्वभार्द

भारतीय प्रौद्योगिकी संस्थान हैदराबाद

Indian Institute of Technology Hyderabad

DIL

Data-driven Intelligence  
& Learning Lab

# Newton-Raphson Iterative optimization

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{H}^{-1} \nabla E(\mathbf{w}^t)$$

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

$$\phi_n = \underline{\phi(x_n)}$$

$$\Phi = \begin{bmatrix} -\phi(x_1) & \vdots & x(y_1 - t_1) \\ -\phi(x_2) & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ -\phi(x_N) & \vdots & x(y_N - t_N) \end{bmatrix}_{N \times M}$$

$$\bar{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

$$\bar{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$$



# Newton-Raphson Iterative optimization

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{H}^{-1} \nabla E(\mathbf{w}^t)$$

$$\Phi = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_N \end{bmatrix} \quad N \times M$$

$$H_{ij} = \frac{\partial^2 E(\mathbf{w}^t)}{\partial \mathbf{w}_i \partial \mathbf{w}_j} = \frac{\partial}{\partial \mathbf{w}_i} \sum_{n=1}^N (y_n - t_n) \phi_j(\mathbf{x}_n) =$$

$$\sum_{i=1}^N \bar{y}_n (1 - \bar{y}_n) \phi_i(x_n) \phi_j(x_n)$$

$$H = \begin{bmatrix} * & \cdots & \cdots & \cdots \\ \vdots & & & \end{bmatrix} = \underbrace{\Phi}_{M \times M} \underbrace{R}_{M \times N} \underbrace{\Phi^T}_{N \times M}$$

$$R = \underbrace{\Phi^T R \Phi}_{N \times N}$$

$$R_{nn} = \bar{y}_n (1 - \bar{y}_n)$$



# Newton-Raphson Iterative optimization

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{H}^{-1} \nabla E(\mathbf{w}^t)$$

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

(y-t)

$$R_{nn} = y_n(1-y_n) f(\mathbf{w}^t, \Phi_n) = y_n = P(x)$$

Iterative Reweighted Least Squares



$$E(w) = \sum_{i=1}^n (f_i - y_i)^2$$

# Next Neural Networks

