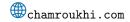


# **TC2: Optimization for Machine Learning**

## Master of Science in AI and Master of Science in Data Science @ UPSaclay 2024/2025.

FAÏCEL CHAMROUKHI







#### week 5 : November 28, 2024. mid-term exam.

#### **Continuous optimization**

(second order methods : Newton methods including Quasi-Newton, secant, IRLS)

# Newton's method I



#### Origin of its construction :

The Newton method is widely known as a technique for finding a root of a function of one variable. Let  $f(x) : \mathbb{R} \to \mathbb{R}$ . Consider the equation :

$$f(x) = 0.$$

The Newton method is based on linear approximation.

**Taylor Approximation for**  $f : \mathbb{R} \to \mathbb{R}$ : If f is differentiable, then for any  $x, x_0 \in \mathbb{R}$ , Provided that  $||x - x_0||$  is small (i.e., x is close to  $x_0$ ), by Taylor's theorem we have

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + R_2(x)$$

where  $R_2(x) = o(||x - x_0||^2)$  is the remainder which vanishes as  $x \to x_0$ . Assume that we get some  $x_0$  close enough to x.

Therefore, the equation f(x) = 0 can be approximated by the linear equation :

$$f(x_0) + f'(x_0)(x - x_0) = 0.$$

# Newton's method II



We then have :

$$x = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Converting this idea in an algorithmic form, we get the process :

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}.$$

In optimization, we are finding the roots of f'(x), i.e. solving

f'(x) = 0.

Hence, the Newton method for optimization problems appears to be in the form

$$x_{k+1} = x_k - [f''(x_k)]^{-1} f'(x_k).$$

assuming f is twice-differentiable.

F. Chamroukhi

## **Newton's Method**



#### Construction as a quadratic approximation

We can obtain the previous process using the idea of quadratic approximation. Consider the quadratic approximation of f(x), centered at the point  $x_k$ :

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2 + o(||x - x_k||^2)$$

Assume that  $f''(x_k) > 0$ . Then we choose  $x_{k+1}$  as the point that minimizes the quadratic approximation f(x).

This means we solve

$$f'(x) = \frac{d}{dx} \{ f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2} f''(x_k)(x - x_k)^2 \} = 0$$

which gives

$$f'(x_k) = f'(x_k) + f''(x_k)(x - x_k) = 0$$

Solving this gives us the Newton process

$$x_{k+1} = x_k - [f''(x_k)]^{-1} f'(x_k)$$

#### in $\mathbb{R}^n$ (Construction as a quadratic approximation) I

**Taylor** : If f is continuously twice differentiable, then for any  $x, x_0 \in \mathbb{R}^n$ , provided that  $||x - x_0||$  is small (i.e., x is close to  $x_0$ ), we can approximate f(x) by : (second-order approximation).

$$f(x) = f(x_0) + \nabla f(x_0)^T (x - x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0) + o \|x - x_0\|^3$$

To derive the Newton update form from the quadratic approximation, consider again the quadratic approximation of f(x) at the current point  $x_k$ :

$$f(x) = f(x_k) + \nabla f(x_k)^T (x - x_k) + \frac{1}{2} (x - x_k)^T \nabla^2 f(x_k) (x - x_k) + o||x - x_k||^3$$

The quadratic approximation,  $f(\boldsymbol{x}),$  captures the behavior of  $f(\boldsymbol{x})$  locally around  $\boldsymbol{x}_k$  :

The quadratic term incorporates the Hessian, which captures the curvature of f.

#### in $\mathbb{R}^n$ (Construction as a quadratic approximation) II



Assume that  $\nabla^2 f(x_k) \succ 0$ . Then we can choose  $x_{k+1}$  as a point of minimum of the quadratic function f(x).

Minimizing f(x) involves taking its gradient w.r.t x and setting it to zero :

$$\nabla f(x) = \nabla f(x_k) + \nabla^2 f(x_k)(x - x_k) = 0$$

This leads to :

$$\nabla^2 f(x_k)(x_{k+1} - x_k) = -\nabla f(x_k)$$

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

Step Size Adjustment :

- The term  $-[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$  uses both gradient and curvature information to determine the optimal direction and step size.

- This approach avoids issues in gradient descent, such as slow convergence in areas of different curvatures.

F. Chamroukhi



This iterative process is applied until convergence.

Each iteration involves recalculating the quadratic approximation at the new point  $x_{k+1}$ , then minimizing this approximation to determine the next point :

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

- The gradient, as usual, tells us the direction of steepest ascent (the anti-gradient determines the steepest descent)

- The Hessian allows the method to **adapt the step size** depending on the curvature.

- In regions of **high curvature**, the inverse Hessian shrinks the step size, preventing overshooting.

- In flat regions, it permits larger steps, speeding up convergence.
- The Newton method can converge faster than standard gradient descent, especially near well-behaved minima, due to its use of second-order information.



# **Quasi-Newton Methods**

Efficient second-order approximation without explicitly computing the Hessian.



# Secant Method

Approximates derivatives using prior information, avoiding full differentiation.



#### Issue : Lack of Positive Definiteness

The matrix  $\nabla^2 f(x)$  may not be positive definite. As a result :

- The problem might lack solutions.
- The descent direction  $-[\nabla^2 f(x)]^{-1} \nabla f(x)$  may not always be effective.

#### Solution : Add Diagonal Matrix E

To resolve this, we add a diagonal matrix **E** such that  $\nabla^2 f(x) + \mathbf{E}$  becomes positive definite.

- **Example** : Let  $\mathbf{E} = \gamma I^n$ , with  $-\gamma$  chosen to be smaller than all non-positive eigenvalues of  $\nabla^2 f(x)$ .
- This modification shifts the original eigenvalues by  $\gamma > 0$ .

The required value of  $\gamma$  is found when solving the "Newton equation" :

$$\nabla^2 f(x)\mathbf{p} = -\nabla f(x)$$

This approach is known as the *Levenberg-Marguardt* modification. Note : As  $\gamma$  becomes larger, **p** increasingly resembles the steepest descent direction F. Chamboukhi



**Issue : Lack of Enough Differentiability** The function f may not belong to  $C^2$ , or computing the second derivatives  $(\nabla^2 f)$  might be too costly.

#### Solution : Use Quasi-Newton Methods

In quasi-Newton methods, we approximate the Newton equation by replacing the Hessian  $\nabla^2 f(x_k)$  with a more computationally efficient matrix  $\mathbf{B}_k$ .

- **B**<sub>k</sub> is computed using gradient values at the current and previous points.
- Using a first-order Taylor expansion for  $\nabla f(x_k)$  :

$$\nabla f(x_k) \approx \nabla f(x_{k-1}) + \nabla^2 f(x_k)(x_k - x_{k-1})$$

then

$$\nabla^2 f(x_k)(x_k - x_{k-1}) \approx \nabla f(x_k) - \nabla f(x_{k-1})$$



#### Updating the Approximation

The matrix  $\mathbf{B}_k$  is chosen to satisfy the following system :

$$\mathbf{B}_k(x_k - x_{k-1}) = \nabla f(x_k) - \nabla f(x_{k-1})$$

#### Secant Method : Special Case for n = 1

For n = 1, this process reduces to the *secant method*, which approximates the second derivative as :

$$f''(x_k) \approx \frac{f'(x_k) - f'(x_{k-1})}{x_k - x_{k-1}}$$

#### universite PARIS-SACLAY

#### Matrix $\mathbf{B}_k$ Properties

The matrix  $\mathbf{B}_k$  has  $n^2$  elements, which makes it *under-determined* by the n equations available. Additional requirements, such as making sure that  $\mathbf{B}_k$  is *symmetric* and *positive definite*, result in specific quasi-Newton methods. **Initialization and Update of \mathbf{B}\_k**. Typically :

• Start with 
$$\mathbf{B}_0 = I^n$$

- Update  $\mathbf{B}_k$  to  $\mathbf{B}_{k+1}$  using a rank-one or rank-two matrix update .
- allows efficient updating of the factorization of  $\mathbf{B}_k$ , utilizing linear algebra

**BFGS Update Rule** : There are many ways to update  $\mathbf{B}_k$ , but the **Broyden-Fletcher-Goldfarb-Shanno** method is among the most effective :  $\mathbf{B}_{k+1} = \mathbf{B}_k - \frac{(\mathbf{B}_k \mathbf{s}_k)(\mathbf{B}_k \mathbf{s}_k)^{\top}}{\mathbf{s}_k^{\top} \mathbf{B}_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^{\top}}{\mathbf{y}_k^{\top} \mathbf{s}_k}$  where :

$$\bullet \mathbf{s}_k = x_{k+1} - x_k$$

• 
$$\mathbf{y}_k = \nabla f(x_{k+1}) - \nabla f(x_k)$$

**Remark** If f is quadratic,  $\mathbf{B}_k$  will converge to the Hessian after n steps.



## IRLS

# Iteratively Reweighted Least Squares Newto descent for Logistic Regression



• **Problem :** Given  $(x_i, y_i) \in \mathbb{R}^p \times \{0, 1\}, i = 1, \dots, n$ , Let  $p_i(\theta) = \frac{\exp(x_i^T \theta)}{1 + \exp(x_i^T \theta)}$  be the logistic function. In logistic regression we minimize the negative log-likelihood :  $\min_{\theta} f(\theta)$ 

$$f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left[ -y_i x_i^T \theta + \log \left( 1 + \exp(x_i^T \theta) \right) \right].$$

• Newton algorithm :  $\theta^{(k+1)} = \theta^{(k)} - \left[\nabla^2 f(\theta^{(k)})\right]^{-1} \nabla f(\theta^{(k)})$ 

Gradient vector and Hessian Matrix : (details in the next slide)

- Gradient vector  $\nabla f(\theta) = \frac{\partial L(\theta)}{\partial \theta} = -\sum_{i=1}^{n} x_i (y_i p_i(\theta))$ .
- Hessian Matrix  $\nabla^2 f(\theta) = \sum_{i=1}^n x_i x_i^\top p_i(\theta) (1 p_i(\theta))$
- The Newton iterative update of  $\theta$  has therefore the following expression :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i (y_i - p_i(x_i; \theta^{(k)}))$$



• **Problem :** Given  $(x_i, y_i) \in \mathbb{R}^p \times \{0, 1\}, i = 1, ..., n$ , Let  $p_i(\theta) = \frac{\exp(x_i^T \theta)}{1 + \exp(x_i^T \theta)}$  be the logistic function.

In logistic regression we minimize the negative log-likelihood :  $\min_{ heta} f( heta)$ 

$$f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left[ -y_i x_i^T \theta + \log \left( 1 + \exp(x_i^T \theta) \right) \right].$$

• Newton algorithm :  $\theta^{(k+1)} = \theta^{(k)} - \left[\nabla^2 f(\theta^{(k)})\right]^{-1} \nabla f(\theta^{(k)})$ 

Gradient vector and Hessian Matrix : (details in the next slide)

- Gradient vector  $\nabla f(\theta) = \frac{\partial L(\theta)}{\partial \theta} = -\sum_{i=1}^{n} x_i(y_i p_i(\theta))$ .
- Hessian Matrix  $\nabla^2 f(\theta) = \sum_{i=1}^n x_i x_i^{\top} p_i(\theta) (1 p_i(\theta))$
- The Newton iterative update of  $\theta$  has therefore the following expression :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i(y_i - p_i(x_i; \theta^{(k)}))$$



**Problem :** Given  $(x_i, y_i) \in \mathbb{R}^p \times \{0, 1\}, i = 1, ..., n$ , Let  $p_i(\theta) = \frac{\exp(x_i^T \theta)}{1 + \exp(x_i^T \theta)}$  be the logistic function.

In logistic regression we minimize the negative log-likelihood :  $\min_{ heta} f( heta)$ 

$$f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left[ -y_i x_i^T \theta + \log \left( 1 + \exp(x_i^T \theta) \right) \right].$$

• Newton algorithm : 
$$\theta^{(k+1)} = \theta^{(k)} - \left[\nabla^2 f(\theta^{(k)})\right]^{-1} \nabla f(\theta^{(k)})$$

Gradient vector and Hessian Matrix : (details in the next slide)

- Gradient vector  $\nabla f(\theta) = \frac{\partial L(\theta)}{\partial \theta} = -\sum_{i=1}^{n} x_i (y_i p_i(\theta))$ .
- Hessian Matrix  $\nabla^2 f(\theta) = \sum_{i=1}^n x_i x_i^\top p_i(\theta) (1 p_i(\theta))$

The Newton iterative update of heta has therefore the following expression :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^\top p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i (y_i - p_i(x_i; \theta^{(k)}))$$



**Problem :** Given  $(x_i, y_i) \in \mathbb{R}^p \times \{0, 1\}, i = 1, ..., n$ , Let  $p_i(\theta) = \frac{\exp(x_i^T \theta)}{1 + \exp(x_i^T \theta)}$  be the logistic function.

In logistic regression we minimize the negative log-likelihood :  $\min_{ heta} f( heta)$ 

$$f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left[ -y_i x_i^T \theta + \log \left( 1 + \exp(x_i^T \theta) \right) \right].$$

• Newton algorithm : 
$$\theta^{(k+1)} = \theta^{(k)} - \left[\nabla^2 f(\theta^{(k)})\right]^{-1} \nabla f(\theta^{(k)})$$

Gradient vector and Hessian Matrix : (details in the next slide)

- Gradient vector  $\nabla f(\theta) = \frac{\partial L(\theta)}{\partial \theta} = -\sum_{i=1}^{n} x_i (y_i p_i(\theta))$ .
- Hessian Matrix  $\nabla^2 f(\theta) = \sum_{i=1}^n x_i x_i^{\top} p_i(\theta) (1 p_i(\theta))$
- The Newton iterative update of  $\theta$  has therefore the following expression :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i (y_i - p_i(x_i; \theta^{(k)}))$$

#### gradient and hessian for the logistic regression's objective



#### Gradient vector

$$\nabla f(\theta) = \frac{\partial f(\theta)}{\partial \theta} = \sum_{i=1}^{n} \left[ \frac{\partial}{\partial \theta} - y_i x_i^{\top} \theta + \frac{\partial}{\partial \theta} \log(1 + \exp(x_i^{\top} \theta)) \right]$$
$$= \sum_{i=1}^{n} -y_i x_i + x_i p_i(\theta) p_i(\theta)$$
$$= -\sum_{i=1}^{n} x_i (y_i - p_i(\theta)) \cdot$$

Hessian matrix :

$$\nabla^2 f(\theta) = \frac{\partial^2 f(\theta)}{\partial \theta \partial \theta^{\top}} = \sum_{i=1}^n x_i \frac{\partial}{\partial \theta^{\top}} \{ \frac{\exp(x_i^{\top} \theta)}{1 + \exp(x_i^{\top} \theta)} \}$$
$$= \sum_{i=1}^n x_i \frac{x_i^{\top} \exp(x_i^{\top} \theta)}{\left(1 + \exp(x_i^{\top} \theta)\right)^2}$$
$$= \sum_{i=1}^n x_i x_i^{\top} p_i(\theta) (1 - p_i(\theta))$$

Iteratively Reweighted Least Squares (IRLS) Let's write the previous quantities in a matrix form (also very useful for coding)<sup>2</sup> Definition of the state of the

- $X = (x_1, \ldots, x_n)^\top$  matrix whose rows are the input vectors  $x_i$
- $y = (y_1, \ldots, y_n)^{\top}$  the vector on binary labels  $y_i$
- $p = (p_1(\theta), \dots, p_n(\theta))^{\top}$  the vector of logistic probabilities
- $W = \operatorname{diag}(p \odot (\mathbf{1}_n p))$  diagonal matrix with  $(W)_{ii} = p_i(\theta) (1 p_i(\theta))$
- z = Xθ<sup>(k)</sup> + (W<sup>(k)</sup>)<sup>-1</sup>(y − p<sup>(k)</sup>) the current approximate response Then we can write the vectorized forms of the gradient and the hesssian :

$$\hookrightarrow \nabla f(\theta) = -\sum_{i=1}^{n} x_i (y_i - p_i(x_i; \theta^{(k)})) = -X^\top (y - p^{(k)})$$

Then we get the vectorized form of the Newton iteration :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i (y_i - p_i(x_i; \theta^{(k)}))$$
$$= \theta^{(k)} + \left(X^{\top} W X\right)^{-1} X^{\top} (y - p^{(k)})$$
$$= (X^{\top} W^{(k)} X)^{-1} X^{\top} W^{(k)} z$$
(1)

The NR update has the form of Least Squares, with weights W, and this is calculated iteratively  $\hookrightarrow$  we call it iteratively re-weighted least squares (IRLS).

F. Chamroukh

TC2: Optimization for Machine Learning:(In progress)

18/18

Iteratively Reweighted Least Squares (IRLS) Let's write the previous quantities in a matrix form (also very useful for coding)<sup>2</sup> Definition of the second s

- $X = (x_1, \ldots, x_n)^{\top}$  matrix whose rows are the input vectors  $x_i$
- $y = (y_1, \ldots, y_n)^{\top}$  the vector on binary labels  $y_i$
- $p = (p_1(\theta), \dots, p_n(\theta))^{\top}$  the vector of logistic probabilities
- $W = \operatorname{diag}(p \odot (\mathbf{1}_n p))$  diagonal matrix with  $(W)_{ii} = p_i(\theta) (1 p_i(\theta))$
- z = Xθ<sup>(k)</sup> + (W<sup>(k)</sup>)<sup>-1</sup>(y − p<sup>(k)</sup>) the current approximate response Then we can write the vectorized forms of the gradient and the hesssian :

Then we get the vectorized form of the Newton iteration :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i (y_i - p_i(x_i; \theta^{(k)}))$$
$$= \theta^{(k)} + \left(X^{\top} W X\right)^{-1} X^{\top} (y - p^{(k)})$$
$$= (X^{\top} W^{(k)} X)^{-1} X^{\top} W^{(k)} z \qquad (1$$

The NR update has the form of Least Squares, with weights W, and this is calculated iteratively  $\hookrightarrow$  we call it iteratively re-weighted least squares (IRLS).

F. Chamroukh

TC2: Optimization for Machine Learning:(In progress

Iteratively Reweighted Least Squares (IRLS) University Let's write the previous quantities in a matrix form (also very useful for coding)<sup>PALIS SACLA</sup>

- $X = (x_1, \ldots, x_n)^\top$  matrix whose rows are the input vectors  $x_i$
- $y = (y_1, \ldots, y_n)^{\top}$  the vector on binary labels  $y_i$
- $p = (p_1(\theta), \dots, p_n(\theta))^{\top}$  the vector of logistic probabilities
- $W = \operatorname{diag}(p \odot (\mathbf{1}_n p))$  diagonal matrix with  $(W)_{ii} = p_i(\theta) (1 p_i(\theta))$
- z = Xθ<sup>(k)</sup> + (W<sup>(k)</sup>)<sup>-1</sup>(y − p<sup>(k)</sup>) the current approximate response Then we can write the vectorized forms of the gradient and the hesssian :

Then we get the vectorized form of the Newton iteration :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i(y_i - p_i(x_i; \theta^{(k)}))$$
$$= \theta^{(k)} + \left(X^{\top} W X\right)^{-1} X^{\top} (y - p^{(k)})$$
$$= (X^{\top} W^{(k)} X)^{-1} X^{\top} W^{(k)} z$$
(1)

The NR update has the form of Least Squares, with weights W, and this is calculated iteratively  $\hookrightarrow$  we call it iteratively re-weighted least squares (IRLS).

F. Chamroukhi

TC2: Optimization for Machine Learning:(In progress

Iteratively Reweighted Least Squares (IRLS) Let's write the previous quantities in a matrix form (also very useful for coding)<sup>PABLS-SACLA</sup>

- $X = (x_1, \ldots, x_n)^\top$  matrix whose rows are the input vectors  $x_i$
- $y = (y_1, \ldots, y_n)^{\top}$  the vector on binary labels  $y_i$
- $p = (p_1(\theta), \dots, p_n(\theta))^{\top}$  the vector of logistic probabilities
- $W = \operatorname{diag}(p \odot (\mathbf{1}_n p))$  diagonal matrix with  $(W)_{ii} = p_i(\theta) (1 p_i(\theta))$
- z = Xθ<sup>(k)</sup> + (W<sup>(k)</sup>)<sup>-1</sup>(y − p<sup>(k)</sup>) the current approximate response Then we can write the vectorized forms of the gradient and the hesssian :

Then we get the vectorized form of the Newton iteration :

$$\theta^{(k+1)} = \theta^{(k)} + \left[\sum_{i=1}^{n} x_i x_i^{\top} p_i(x_i; \theta^{(k)}) (1 - p_i(x_i; \theta^{(k)}))\right]^{-1} \sum_{i=1}^{n} x_i(y_i - p_i(x_i; \theta^{(k)}))$$
  
$$= \theta^{(k)} + \left(X^{\top} W X\right)^{-1} X^{\top} (y - p^{(k)})$$
  
$$= (X^{\top} W^{(k)} X)^{-1} X^{\top} W^{(k)} z$$
(1)

The NR update has the form of Least Squares, with weights W, and this is calculated iteratively  $\hookrightarrow$  we call it iteratively re-weighted least squares (IRLS).

F. Chamroukh

C2: Optimization for Machine Learning:(In progress

18/18