

1. **(Math)** Nonlinear least-squares. Suppose that $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})) : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbf{x} \in \mathbb{R}^n, \mathbf{f} \in \mathbb{R}^m$ and some $f_i(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is a (are) non-linear function(s). Then, the problem,

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{f}(\mathbf{x})\|_2^2 = \arg \min_{\mathbf{x}} \frac{1}{2} (\mathbf{f}(\mathbf{x}))^T \mathbf{f}(\mathbf{x})$$

is a nonlinear least-squares problem. In our lecture, we mentioned that Levenberg-Marquardt algorithm is a typical method to solve this problem. In L-M algorithm, for each updating step, at the current \mathbf{x} , a local approximation model is constructed as,

$$\begin{aligned} L(\mathbf{h}) &= \frac{1}{2} (\mathbf{f}(\mathbf{x} + \mathbf{h}))^T \mathbf{f}(\mathbf{x} + \mathbf{h}) + \frac{1}{2} \mu \mathbf{h}^T \mathbf{h} \\ &= \frac{1}{2} (\mathbf{f}(\mathbf{x}))^T \mathbf{f}(\mathbf{x}) + \mathbf{h}^T (\mathbf{J}(\mathbf{x}))^T \mathbf{f}(\mathbf{x}) + \frac{1}{2} \mathbf{h}^T (\mathbf{J}(\mathbf{x}))^T \mathbf{J}(\mathbf{x}) \mathbf{h} + \frac{1}{2} \mu \mathbf{h}^T \mathbf{h} \end{aligned}$$

where $\mathbf{J}(\mathbf{x})$ is $\mathbf{f}(\mathbf{x})$'s Jacobian matrix, and $\mu > 0$ is the damped coefficient. Please prove that $L(\mathbf{h})$ is a strictly convex function. (Hint: If a function $L(\mathbf{h})$ is differentiable up to at least second order, L is strictly convex if its Hessian matrix is positive definite.)

we know $\mathbf{h} \in \mathbb{R}^n$, and $L(\mathbf{h})$ is differentiable at second order, $\text{dom } L \in \mathbb{R}^n$, then

$$\begin{aligned} \nabla L(\mathbf{h}) &= (\mathbf{J}(\mathbf{x}))^T \mathbf{f}(\mathbf{x}) + \frac{1}{2} ((\mathbf{J}(\mathbf{x}))^T \mathbf{J}(\mathbf{x}) + (\mathbf{J}(\mathbf{x}))^T \mathbf{J}(\mathbf{x})) \mathbf{h} + \mu \mathbf{h} \\ &= (\mathbf{J}(\mathbf{x}))^T \mathbf{f}(\mathbf{x}) + (\mathbf{J}(\mathbf{x}))^T \mathbf{J}(\mathbf{x}) \mathbf{h} + \mu \mathbf{h} \end{aligned}$$

$$\text{so } \nabla^2 L(\mathbf{h}) = \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \mu \mathbf{I} \quad (\mu > 0)$$

note $\mathbf{J}(\mathbf{x})$ as \mathbf{J} , $\mathbf{J}^T \in \mathbb{S}^n$, $\mathbf{I} \in \mathbb{S}^n$, so $\mathbf{J}^T \mathbf{J} + \mu \mathbf{I} \in \mathbb{S}^n$

we need to prove $\mathbf{J}^T \mathbf{J} + \mu \mathbf{I}$ is positive definite

$$\forall \mathbf{x} \neq 0, \text{ we have } \mathbf{x}^T \mathbf{J}^T \mathbf{J} \mathbf{x} = (\mathbf{J} \mathbf{x})^T \mathbf{J}(\mathbf{x}) = \|\mathbf{J} \mathbf{x}\|_2^2 \geq 0$$

then $\mathbf{J}^T \mathbf{J}$ is positive semi-positive

for $\mathbf{J}^T \mathbf{J}$'s eigenvalues $\lambda_i, i=1, \dots, n$, we have $\lambda_i \geq 0$
for positive semi-positive matrix

$$\mathbf{J}^T \mathbf{J} \mathbf{x}_i = \lambda_i \mathbf{x}_i$$

$$(\mathbf{J}^T \mathbf{J} + \mu \mathbf{I}) \mathbf{x}_i = (\lambda_i + \mu) \mathbf{x}_i, \lambda_i + \mu \text{ is the eigenvalue for}$$

$$\mathbf{J}^T \mathbf{J} + \mu \mathbf{I}, \text{ and } \mu > 0 \Rightarrow \lambda_i + \mu > 0$$

thus $\mathbf{J}^T \mathbf{J} + \mu \mathbf{I}$, the Hessian matrix of $L(\mathbf{h})$

is positive definite

then $L(\mathbf{h})$ is strictly convex

2. **(Math)** In our lecture, we mentioned that for logistic regression, the cost function

is,

$$J(\theta) = -\sum_{i=1}^m y_i \log(h_{\theta}(x_i)) + (1-y_i) \log(1-h_{\theta}(x_i))$$

Please verify that the gradient of this cost function is

$$\nabla_{\theta} J(\theta) = \sum_{i=1}^m x_i (h_{\theta}(x_i) - y_i)$$

$$h_{\theta}(x) = \sigma(\theta^T x + b) = \frac{1}{1 + e^{-(\theta^T x + b)}}$$

$$\text{Let } t = \theta^T x + b, \quad \frac{\partial t}{\partial \theta} = x$$

$$\frac{\partial h_{\theta}(x)}{\partial t} = \frac{\partial \sigma(t)}{\partial t} = \sigma(t) (1 - \sigma(t)) = h_{\theta}(x) (1 - h_{\theta}(x))$$

$$\begin{aligned} \text{so: } \nabla_{\theta} J(\theta) &= -\sum_{i=1}^m \left(y_i \frac{\partial \log(h_{\theta}(x_i))}{\partial h_{\theta}(x_i)} \cdot \frac{\partial h_{\theta}(x_i)}{\partial t} \cdot \frac{\partial t}{\partial \theta} \right. \\ &\quad \left. + (1-y_i) \frac{\partial \log(1-h_{\theta}(x_i))}{\partial h_{\theta}(x_i)} \cdot \frac{\partial h_{\theta}(x_i)}{\partial t} \cdot \frac{\partial t}{\partial \theta} \right) \\ &= -\sum_{i=1}^m \left(y_i \frac{1}{h_{\theta}(x_i)} \cdot h_{\theta}(x_i) (1-h_{\theta}(x_i)) x_i \right. \\ &\quad \left. + (1-y_i) \frac{-1}{(1-h_{\theta}(x_i))} h_{\theta}(x_i) (1-h_{\theta}(x_i)) x_i \right) \\ &= -\sum_{i=1}^m y_i (1-h_{\theta}(x_i)) x_i - (1-y_i) h_{\theta}(x_i) x_i \\ &= -\sum_{i=1}^m x_i (y_i - h_{\theta}(x_i)) \\ &= \sum_{i=1}^m x_i (h_{\theta}(x_i) - y_i) \end{aligned}$$