

Утверждение (Pedregosa, 2016)

Пусть L — дифференцируемая функция, такая что все стационарные точки L являются локальными минимумами. Пусть также гессиан H^{-1} функции потерь L является обратимым в каждой стационарной точке. Тогда

$$\nabla_{\mathbf{h}} Q(T(\theta_0, \mathbf{h}), \mathbf{h}) = \nabla_{\mathbf{h}} Q(\theta^\eta, \mathbf{h}) - \nabla_{\mathbf{h}} \nabla_{\theta} L(\theta^\eta, \mathbf{h})^\top \mathbf{H}^{-1} \nabla_{\theta} Q(\theta^\eta, \mathbf{h})$$

Доказательство

Для стационарной точки выполняется, что $\nabla_{\theta} L(T(\theta_0, h)) = 0$

Следовательно $\nabla_h (\nabla_{\theta} L(T(\theta_0, h))) = \nabla_{\theta, h} L(\theta^\eta, h) + \nabla_{\theta}^2 L(\theta^\eta, h) \nabla_h \theta = 0$

Из предыдущего выражаем $\nabla_h \theta = -(\nabla_{\theta}^2 L(\theta^\eta, h))^{-1} \nabla_{\theta, h} L(\theta^\eta, h)$

Тогда

$$\begin{aligned} \nabla_h Q(T(\theta_0, h)) &= \nabla_h Q(\theta^\eta, h) + \nabla_{\theta} Q(\theta^\eta, h)^\top \nabla_h \theta \\ &= \nabla_h Q(\theta^\eta, h) - \nabla_{\theta} Q(\theta^\eta, h)^\top (\nabla_{\theta}^2 L(\theta^\eta, h))^{-1} \nabla_{\theta, h} L(\theta^\eta, h) \\ &= \nabla_h Q(\theta^\eta, h) - \nabla_{\theta, h} L(\theta^\eta, h)^\top (\nabla_{\theta}^2 L(\theta^\eta, h))^{-1} \nabla_{\theta} Q(\theta^\eta, h) \end{aligned}$$