

IA353 - Redes Neurais

* Técnicas de otimização não-linear suavizada

Neste material, vamos cobrir de maneira sucinta os aspectos básicos de diferentes técnicas de otimização não-linear que podem ser utilizadas no treinamento de redes neurais do tipo MLP.

④ Relação entre gradiente e Hessiana

- o método de Newton explora a aproximação local da função de erro $J(\theta)$ em série de Taylor até o termo de 2^a ordem.

$$J_{\text{quad}}(\theta) = J(\theta_i) + \nabla J(\theta_i)^T (\theta - \theta_i) + \frac{1}{2} [\theta - \theta_i]^T \nabla^2 J(\theta_i) [\theta - \theta_i]$$

Qual $\Delta\theta = \theta_{i+1} - \theta_i$ nos leva a um ponto ótimo da aproximação quadrática?

$$\frac{\partial J_{\text{quad}}(\theta_{i+1})}{\partial \theta_{i+1}} = 0 \quad (\text{condição de otimalidade do ponto } \theta_{i+1})$$

$$\underbrace{\frac{\partial J(\theta_i)}{\partial \theta_{i+1}}}^{=0} + \nabla J(\theta_i) + \nabla^2 J(\theta_i) [\theta_{i+1} - \theta_i] = 0$$

$$\Rightarrow \theta_{i+1} - \theta_i = \Delta\theta = -(\nabla^2 J(\theta_i))^{-1} \nabla J(\theta_i)$$

Se $J(\theta)$ não for uma função quadrática, $J(\theta_{i+1})$ não necessariamente será menor que $J(\theta_i)$. Por isso, a lei de ajuste passa a incluir um parâmetro α_i p/ controlar o tamanho do passo:

$$\theta_{i+1} = \theta_i - \alpha_i [\nabla^2 J(\theta_i)]^{-1} \nabla J(\theta_i)$$

PROBLEMAS : positividade da matriz Hessiana \Rightarrow p/ haver inversa e p/ buscarmos um ponto de mínimo

calculo exato da matriz Hessiana e de seu inverso

O vetor gradiente também pode ser aproximado por uma série de Taylor:

$$\nabla J(\theta_{i+1}) = \nabla J(\theta_i + \rho_i) = \nabla J(\theta_i + \alpha_i d_i)$$

↳ direção: $-\nabla^2 J(\theta) \nabla J(\theta)$ - NEWTON

$$\begin{aligned} &= \nabla J(\theta_i) + \nabla^2 J(\theta_i) [\theta_{i+1} - \theta_i] \\ &= \nabla J(\theta_i) + \nabla^2 J(\theta_i) \cdot \rho_i \end{aligned}$$

\Rightarrow Seja $g_{i+1} = \nabla J(\theta_{i+1})$ e $g_i = \nabla J(\theta_i)$. A diferença dos vetores gradientes acaba trazendo informações sobre a matriz Hessiana:

11

$$g_{i+1} - g_i = \nabla^2 J(\theta_i) \cdot p_i$$

Esta relação não é explorada nos métodos de otimização quasi-Newton, como os algoritmos DFP e BFGS, já mencionados anteriormente.

• Método de Levenberg-Maquardt

- pode ser visto como uma solução de compromisso entre o método de Newton - que converge rapidamente próximo de um mínimo local ou global - e o método de gradiente - que tem convergência asegurada com uma escolha apropriada do tamanho do passo, mas que pode ser lento.

$$\text{Idéia: } \Delta\theta = [H + \lambda I]^{-1} \cdot g \quad \rightarrow \text{vetor gradiente: } \nabla J(\theta)$$

$$\downarrow \text{matriz hermitiana: } \nabla^2 J(\theta)$$

λ é um parâmetro de regularização que força a matriz $H + \lambda I$ ser positiva definida.

Outra diferença deste método está associado à forma como a matriz hermitiana é calculada. Considera a função de erro $J(\theta) = \sum_{i=1}^N \sum_{j=1}^n (d_j(i) - y_j(i))^2 = \sum_{i=1}^N e_i^2$, onde e_i é o erro residual (errado 1 para cada saída da rede em cada instante de tempo ou p/ cada padrão de entrada).

$$(1) \frac{\partial J(\theta)}{\partial \theta} = \sum_{i=1}^N 2e_i \cdot \frac{\partial e_i}{\partial \theta}$$

Seja $J \in \mathbb{R}^{N \times N_{\text{pesos}}}$ o jacobiano do funcional $J(\theta)$

$$J = \begin{bmatrix} \nabla e_1^T \\ \vdots \\ \nabla e_N^T \end{bmatrix} = \begin{bmatrix} \frac{\partial e_1}{\partial \theta_1} & \dots & \frac{\partial e_1}{\partial \theta_{N_{\text{pesos}}}} \\ \vdots & \ddots & \vdots \\ \frac{\partial e_N}{\partial \theta_1} & \dots & \frac{\partial e_N}{\partial \theta_{N_{\text{pesos}}}} \end{bmatrix}$$

todas as derivadas do erro referente ao 1º-padrão c/ respeito a todos os N_{pesos} pesos da rede.

$$\begin{bmatrix} \nabla e_1^T \\ \vdots \\ \nabla e_N^T \end{bmatrix} = \begin{bmatrix} \frac{\partial e_1}{\partial \theta_1} & \dots & \frac{\partial e_1}{\partial \theta_{N_{\text{pesos}}}} \\ \vdots & \ddots & \vdots \\ \frac{\partial e_N}{\partial \theta_1} & \dots & \frac{\partial e_N}{\partial \theta_{N_{\text{pesos}}}} \end{bmatrix}$$

todas as derivadas do erro referente ao N -ésimo padrão, c/ respeito a todos os N_{pesos} pesos da rede (a mesma saída)

Então, podemos escrever que $\frac{\partial J(\theta)}{\partial \theta} = 2 J^T \cdot \underline{e}$, onde $\underline{e} = [e_1 \dots e_N]^T$.

$$(2) \frac{\partial^2 J(\theta)}{\partial \theta} = \sum_{i=1}^N 2 \cdot \frac{\partial e_i}{\partial \theta} \cdot \frac{\partial e_i}{\partial \theta}^T + 2 \sum_{i=1}^N e_i \cdot \frac{\partial^2 e_i}{\partial \theta^2}$$

$$= 2 \cdot J^T J + 2 \sum_{i=1}^N e_i \cdot \nabla^2 e_i$$

Quando os erros residuais são suficientemente pequenos, a matriz hermitiana pode ser aproximada pelo primeiro termo:

$$\frac{\partial^2 J(\theta)}{\partial \theta} = \nabla^2 J(\theta) \approx 2J^T J \Rightarrow \text{estimativa usando apenas as derivadas de } J^{\text{a}} \text{ ordem}$$

O parâmetro λ desempenha um papel bastante importante no funcionamento do algoritmo:

- λ muito pequeno \Rightarrow nos aproximações do método de Newton
- λ muito grande \Rightarrow a informação da Hessiana tem pouco efeito e o comportamento se assemelha ao do método de gradiente

PROCEDIMENTO SUGERIDO:

- (i) compute $J(\theta_i)$
- (ii) escolha um valor modesto para λ , e.g., $\lambda = 10^{-3}$
- (iii) obtenha $\Delta\theta = (H + \lambda I)^{-1}g$, calcule $J(\theta_{i+1} + \Delta\theta)$
- (iv) se $J(\theta_i) > J(\theta_{i+1})$ - ou seja, a atualização piorou o desempenho da rede - , aumente λ por um fator (e.g., 10) e retorne ao passo (iii).
- (v) se $J(\theta_i) < J(\theta_{i+1})$, reduza λ por um fator (e.g., 10), efetue a alteração nos pesos ($\theta_i = \theta_{i+1} + \Delta\theta$) e avance para a próxima iteração.

* Métodos Quasi-Newton : DFP e BFGS

- trabalham com uma aproximação iterativa da inversa da matriz Hessiana, de forma que: $\lim_{i \rightarrow \infty} H_i = \nabla^2 J(\theta)^{-1}$
- em cada iteração, a inversa da Hessiana é aprimorada pela soma de duas matrizes simétricas de rank 1, procedimento que é chamado de correção de parte 2.

$$\text{DFP: } H_{i+1} = H_i + \frac{p_i p_i^T}{p_i^T q_i} - \frac{H_i q_i q_i^T H_i}{q_i^T H_i q_i}, \quad i = 0, \dots, N_{\text{passos}} - 1,$$

onde $p_i = \alpha_i d_i$ e $q_i = g_{i+1} - g_i$.

Inicialmente, a direção de minimização é dada pelo vetor gradiente e a matriz Hessiana corresponde a uma matriz identidade. A cada N_{passos} iterações, o algoritmo é reinicializado.

$$\text{BFGS: } H_{i+1} = H_i + \frac{p_i p_i^T}{p_i^T q_i} \left[I + \frac{q_i^T H_i q_i}{p_i^T q_i} \right] - \frac{H_i q_i q_i^T H_i}{p_i^T q_i}$$

Em ambos os casos, o tamanho do passo α_i em cada iteração é determinado

com o auxílio de uma busca unidimensional.

④ Métodos baseados no gradiente conjugado

- são métodos projetados para exigir menos cálculos que o método de Newton e apresentar taxas de convergência maiores que as do método do gradiente.

A regra básica de atualização $\theta_{i+1} = \theta_i + \alpha_i d_i$ nos permite escrever que a solução ótima θ^* , uma vez obtida, é dada por:

$$\theta^* = \alpha_0 d_0 + \alpha_1 d_1 + \dots = \sum \alpha_i d_i$$

Como θ^* é um ponto no espaço $\mathbb{R}^{N_{\text{pesos}}}$, ele pode ser representado como uma combinação linear de N_{pesos} vetores que formam uma base.

$$\theta^* = \alpha_0 d_0 + \dots + \alpha_{N_{\text{pesos}}-1} d_{N_{\text{pesos}}-1} = \sum_{i=0}^{N_{\text{pesos}}-1} \alpha_i d_i \quad (\text{A})$$

desde que o conjunto $\{d_0, \dots, d_{N_{\text{pesos}}-1}\}$ forme uma base de $\mathbb{R}^{N_{\text{pesos}}}$ e $\alpha = [\alpha_0 \dots \alpha_{N_{\text{pesos}}-1}]^T$ seja a representação de θ^* neste base.

O princípio de operação do método do gradiente conjugado para obter θ^* em até N_{pesos} iterações consiste em buscar vetores $d_i \in \mathbb{R}^{N_{\text{pesos}}}$, $i=0, \dots, N_{\text{pesos}}-1$, linearmente independentes e encontrar a representação de θ^* neste base, ou seja, os coeficientes α_i .

Definição: dado uma matriz simétrica A ($N_{\text{pesos}} \times N_{\text{pesos}}$), as direções $d_i \in \mathbb{R}^{N_{\text{pesos}}}$ são A -conjugadas (ou A -ortogonais) se

$$d_j^T A d_i = 0, \quad i \neq j$$

- note que esta definição é uma generalização do conceito de ortogonalidade entre vetores (que ocorre p/ $A = I$).

- se A for simétrica e definida positiva, as direções A -conjugadas são necessariamente linearmente independentes. Consequentemente, tal conjunto $\{d_0, \dots, d_{N_{\text{pesos}}-1}\}$ forma uma base de $\mathbb{R}^{N_{\text{pesos}}}$.

- multiplicando (A) à esquerda por $d_j^T A$, obtemos

$$\begin{aligned} d_j^T A \theta^* &= \sum_{i=0}^{N_{\text{pesos}}-1} \alpha_i d_j^T A d_i \\ &= \alpha_j d_j^T A d_j \end{aligned}$$

Então, $\alpha_j = d_j^T A \theta^*$, $j = 0, \dots, N_{\text{pesos}}-1$

Esta expressão ainda não fornece uma solução viável, pois os coeficientes da representa-

glo de θ^* dependem do próprio vetor θ^* .

• para eliminar θ^* , duas hipóteses adicionais são feitas:

$$\ast \text{ o problema é quadrático} \Rightarrow J(\theta) = \frac{1}{2} \theta^T Q \theta - b^T \theta$$

No ponto de mínimo (θ^*), o gradiente deve ser nulo.

$$\nabla J(\theta^*) = 0 \Rightarrow Q\theta^* - b = 0 \Rightarrow Q\theta^* = b$$

• a matriz Q é a própria matriz A

Com isto:

$$\alpha_j = \frac{d_j^T b}{d_j^T Q d_j}, \quad j=0, \dots, N_{\text{passos}}-1$$

$$\theta^* = \sum_{i=0}^{N_{\text{passos}}-1} \frac{d_i^T b}{d_i^T Q d_i} \cdot d_i$$

Como a solução θ^* deve ser obtida iterativamente a partir de uma condicão inicial θ_0 , podemos considerar que

$$\theta^* = \theta_0 + \alpha_0^* d_0 + \alpha_1^* d_1 + \dots + \alpha_{N_{\text{passos}}-1}^* d_{N_{\text{passos}}-1}$$

sendo os coeficientes dados por

$$\alpha_j^* = \frac{d_j^T Q(\theta^* - \theta_0)}{d_j^T Q d_j}, \quad j=0, \dots, N_{\text{passos}}-1$$

No j -ésima iteração ($0 < j < N_{\text{passos}}-1$), θ_j assume a forma

$$\theta_j = \theta_0 + \alpha_0^* d_0 + \dots + \alpha_{j-1}^* d_{j-1}$$

Põe-multiplicando à esquerda por $d_j^T Q$,

$$d_j^T Q \theta_j = d_j^T Q \theta_0 + \underbrace{\alpha_0^* d_j^T Q d_0}_{=0 \text{ pelo } Q\text{-ortogonalidade}} + \dots$$

Usando esta igualdade na equação de α_j^* :

$$\alpha_j^* = \frac{d_j^T Q(\theta^* - \theta_j)}{d_j^T Q d_j}, \quad j=0, \dots, N_{\text{passos}}-1$$

Como $\nabla J(\theta_j) = Q\theta_j - b$ e $Q\theta^* = b$, é possível perceber que

$$\alpha_j^* = - \frac{d_j^T \nabla J(\theta_j)}{d_j^T Q d_j}, \quad j=0, \dots, N_{\text{passos}}-1$$

Assim, a lei de ajuste do método dos direcções conjugadas é dada por:

$$\theta_{i+1} = \theta_i - \frac{d_i^T \nabla J(\theta_i)}{d_i^T Q d_i} \cdot d_i \quad (B)$$

alcançando a solução ótima em N_{passos} iterações no problema for quadrático.

⇒ Antes de aplicar a regra de atualização em (B), é necessário obter as direções conjugadas $d_i \in \mathbb{R}^{N_{\text{passos}}}$, $i = 0, \dots, N_{\text{passos}} - 1$. Uma maneira de obtê-las é dada a seguir:

$$\begin{cases} d_0 = -\nabla J(\theta_0) \\ d_{i+1} = -\nabla J(\theta_{i+1}) + \beta_i d_i, \end{cases}$$

com $\beta_i = \frac{\nabla J(\theta_{i+1})^T Q d_i}{d_i^T Q d_i}$

$$d_i^T Q d_i$$

Quando temos um problema não-quadrático, a matriz Q será aproximada pela Hessiana calculada sobre o ponto θ_i (lembre da expansão em série de Taylor), ou seja, $Q = \nabla^2 J(\theta_i)$.

Agora, o metodo do gradiente conjugado pode ser resumido da seguinte forma:

- calcule $\nabla J(\theta_0)$ e $d_0 = -\nabla J(\theta_0)$

- para $i = 0, \dots, N_{\text{passos}} - 1$, faça:

(a) calcule $\nabla^2 J(\theta_i)$

(b) $\theta_{i+1} = \theta_i + \alpha_i d_i$, com $\alpha_i = \frac{-d_i^T \nabla J(\theta_i)}{d_i^T \nabla^2 J(\theta_i) d_i}$

(c) calcule $\nabla^2 J(\theta_{i+1})$

(d) se $i < N_{\text{passos}} - 1$, $d_{i+1} = -\nabla J(\theta_{i+1}) + \beta_i d_i$, onde

$$\beta_i = \frac{\nabla J(\theta_{i+1})^T \nabla^2 J(\theta_i) d_i}{d_i^T \nabla^2 J(\theta_i) d_i}$$

Vantagens: (i) não há busca unidimensional de passo

(ii) não há inversão da matriz Hessiana

Desvantagens: (i) cálculo e armazenamento de $\nabla^2 J(\cdot)$ a cada iteração

(ii) como o ajuste de α considera informações da função até a 2ª ordem (e a mesma não necessariamente é quadrática), não há garantia de que o passo seja minimizante.

(iii) o algoritmo não pode ser considerado globalmente convergente, pois não existe garantia de a matriz $\nabla^2 J(\cdot)$ ser definida positiva em todo o espaço $\mathbb{R}^{N_{\text{passos}}}$.

Algumas adotadas foram propostas na tentativa de corrigir alguns destes desvantagens.

⇒ Polak-Ribière (PR):

- α_i é obtido com o auxílio de uma busca unidimensional
- β_i é calculado de forma aproximada:

$$\beta_i = \frac{g_{i+1}^T (g_{i+1} - g_i)}{g_i^T g_i} \quad \begin{matrix} \text{usa apenas informações} \\ \text{do gradiente} \\ (1^{\text{a}} \text{ ordem}) \end{matrix}$$

⇒ Fletcher-Reeves (FR):

- também usa a busca unidimensional

$$\beta_i = \frac{\|g_{i+1}\|^2}{\|g_i\|^2}$$

⇒ Gradiente Conjugado Escalonado (SCG):

- evita a busca unidimensional usando uma abordagem parecida com aquela introduzida no método Levenberg-Marquardt p/ escalarizar o tamanho do passo d_i .

O passo d_i no método do gradiente conjugado é dado por $d_i = -\frac{d_i^T \nabla^2 J(\theta_i)}{d_i^T \nabla^2 J(\theta_i) d_i}$

A ideia de Moller (1993) é estimar o termo $s_i = \nabla^2 J(\theta_i) d_i$ através de aproximação:

$$s_i = \nabla^2 J(\theta_i) d_i \approx \frac{\nabla J(\theta_i + \sigma_i d_i) - \nabla J(\theta_i)}{\sigma_i}, \quad 0 < \sigma_i \ll 1$$

Para evitar problemas c/ a positividade (implícita) da Hessiana, Moller (1993) introduziu um termo de regularização (no mesmo espírito de Levenberg-Marquardt):

$$s_i = \frac{\nabla J(\theta_i + \sigma_i d_i) - \nabla J(\theta_i)}{\sigma_i} + \lambda_i d_i$$

O denominador de d_i , $d_i^T \nabla^2 J(\theta_i) d_i$ será representado por $S_i = d_i^T s_i$. Se usarmos a matriz Hessiana, o denominador sempre será positivo. Por isso, o parâmetro λ_i será ajustado a cada iteração alterando o sinal de S_i .

No artigo, Moller (1993) mostrou como ajustar λ_i de modo a garantir $S_i > 0$.

Além disso, como d_i escolha a matriz herminia de uma maneira artifical, a aproximação quadrática da função - $J_{\text{quad}}(\theta)$ - sobre a qual o algoritmo atua pode não ser uma boa aproximação de $J(\theta)$ em alguns pontos. Por isso, mesmo quando a matriz é positiva definida, pode ser necessário ajustar d_i visando manter uma boa aproximação de $J(\theta)$.

- Moller, então, define:

$$\lambda_i = \frac{J(\theta_i) - J(\theta_i + \alpha_i d_i)}{J(\theta_i) - J_{\text{quad}}(\alpha_i d_i)} = \frac{2s_i [J(\theta_i) - J(\theta_i + \alpha_i d_i)]}{\mu_i^2},$$

onde $\mu_i = -d_i^T \nabla J(\theta_i)$.

λ_i é um medida de quão boa é a aproximação quadrática $J_{\text{quad}}(\alpha_i d_i)$ em relação a $J(\theta + \alpha_i d_i)$. Quanto mais próximo de 1, melhor a aproximação. Por isso, Moller propõe:

$$\lambda_i > 0,75, \text{ então } x_i = \frac{1}{2} x_i$$

$$\lambda_i < 0,25, \text{ então } x_i = 4x_i$$

⇒ Gradiente Conjugado Escalonado Modificado (SCGM)

- em vez de trabalhar com a menor aproximação (s_i) do produto do herminio pelo direcção d_i , utiliza-se o operador diferencial proposto por Peardmutter (1994) para obter o vetor produto de maneira exata.

Considera a expansão de Taylor do vetor gradiente $\nabla J(\cdot)$ em torno de θ :

$$\nabla J(\theta + \Delta\theta) = \nabla J(\theta) + \nabla^2 J(\theta) \Delta\theta + O(\|\Delta\theta\|^2)$$

Para $\Delta\theta = \alpha v$, v a constante e próxima de zero e $v \in \mathbb{R}^{N_{\text{pesos}}}$ um vetor arbitrário:

$$\nabla J(\theta + \alpha v) = \nabla J(\theta) + \nabla^2 J(\theta) \alpha v + O(\alpha^2)$$

$$\therefore \nabla^2 J(\theta) \cdot v = \frac{\nabla J(\theta + \alpha v) - \nabla J(\theta)}{\alpha} + O(\alpha)$$

Fazendo o limite $\alpha \rightarrow 0$:

$$\nabla^2 J(\theta) \cdot v = \lim_{\alpha \rightarrow 0} \frac{\nabla J(\theta + \alpha v) - \nabla J(\theta)}{\alpha} = \frac{\partial \nabla J(\theta + \alpha v)}{\partial \alpha} \Big|_{\alpha=0}$$

OPERADOR $R\{\cdot\}$: $R_v\{f(\theta)\} = \frac{\partial f(\theta + \alpha v)}{\partial \alpha} \Big|_{\alpha=0}$

$$\text{Logo, } R_v\{\nabla J(\theta)\} = \nabla^2 J(\theta) \cdot v$$

Como $R\{\cdot\}$ é um operador diferencial, ele obedece às regras usuais de diferenciação.