

Machine Learning per la fisica applicata e la fisica delle alte energie

Lezione 4: Ottimizzazione.

Abbiamo visto nelle lezioni precedenti che lo scopo del machine learning è la stima dei parametri $\bar{\theta}$ (e/o degli iperparametri Φ) di un modello, cioè la determinazione dei loro valori di best-fit e dei corrispondenti intervalli di credibilità. Ciò si realizza minimizzando una funzione scalare (la loss function) $\mathcal{L}: \Theta \rightarrow \mathbb{R}$ (con $\bar{\theta} \in \Theta$)

$$\boxed{\bar{\theta}^* = \underset{\bar{\theta}}{\operatorname{argmin}} \mathcal{L}(\bar{\theta})} \quad \text{o meglio} \quad \boxed{\hat{\bar{\theta}} = \underset{\bar{\theta}}{\operatorname{argmin}} \mathcal{L}(\bar{\theta})} \quad \text{per } \bar{\theta}^* \sim \hat{\bar{\theta}}.$$

Assumiamo per semplicità che $\Theta \subseteq \mathbb{R}^D$, con D il numero delle variabili ottimizzate. L'ottimizzazione è dunque CONTINUA.

Osservazioni

- 1) Lo spazio dei parametri è tale che può esistere un solo minimo globale o più minimi globali, oltre ad una serie di minimi locali. Il punto di minimo è, computazionalmente, sempre definito in termini locali, cioè $\exists \delta > 0, \forall \bar{\theta} \in \Theta / \|\bar{\theta} - \hat{\bar{\theta}}\| < \delta, \mathcal{L}(\hat{\bar{\theta}}) \leq \mathcal{L}(\bar{\theta})$.
 Un minimo locale può essere circondato da minimi simili (flat direction). Un algoritmo che converge ad un punto stazionario è detto GLOBALLY CONVERGENT.
- 2) Per loss function che siano due volte differenziabili, possiamo considerare $g(\bar{\theta}) = \nabla \mathcal{L}(\bar{\theta})$ il vettore gradiente e $H(\bar{\theta}) = \nabla^2 \mathcal{L}(\bar{\theta})$ la matrice Hessiana. Per un punto $\hat{\bar{\theta}} \in \mathbb{R}^D$ e $\hat{g} = g(\hat{\bar{\theta}})|_{\hat{\bar{\theta}}}$ e $\hat{H} = H(\hat{\bar{\theta}})|_{\hat{\bar{\theta}}}$
 - CONDIZIONE NECESSARIA: se $\hat{\bar{\theta}}$ è un minimo locale, allora $\hat{g} = \vec{0}$ e \hat{H} è positiva semi-definita
 - CONDIZIONE SUFFICIENTE: se $\hat{g} = \vec{0}$ e \hat{H} è positiva definita, allora $\hat{\bar{\theta}}$ è un minimo locale.

3) Talora potremmo avere dei vincoli nella parametrizzazione (CONSTRAINED PARAMETRIZATION). I vincoli sono generalmente categorizzati in equazioni e disequazioni (e.g. $h_k(\bar{\theta}) = 0$ per $k \in \mathcal{E}$; $g_j(\bar{\theta}) \leq 0$ per $j \in \mathcal{I}$). Il set di vincoli è dunque

$$C = \{ \bar{\theta} : g_j(\bar{\theta}) \leq 0 : j \in \mathcal{I}, h_k(\bar{\theta}) = 0 : k \in \mathcal{E} \} \subseteq \mathbb{R}^D$$

$$\hat{\theta} = \underset{\bar{\theta} \in C}{\operatorname{argmin}} \mathcal{L}(\bar{\theta}) \quad (\text{se } \bar{\theta} \in \mathbb{R}^D \text{ si parla di UNCONSTRAINED PARAM.})$$

4) Può essere utile determinare se una funzione di loss è convessa. Se $f: \mathbb{R}^n \rightarrow \mathbb{R}$ è doppio differenziabile, allora f è convessa se e solo se $H = \nabla^2 f(x)$ è positiva semi-definita $\forall x$; f si dice strettamente convessa se H è positiva definita.

5) Quando la loss function non sia continua (in un punto) allora si può definire i sottogradienti: $\bar{g} \in \mathbb{R}^n$ è un sottogradiente di f in \bar{x} , con $\bar{x} \in \operatorname{Dom}(f)$, se $\forall \bar{z} \in \operatorname{Dom}(f)$, $f(\bar{z}) \geq f(\bar{x}) + \bar{g}^T(\bar{z} - \bar{x})$.

4.1 METODI DEL PRIMO ORDINE

Questi metodi sono basati sulle derivate prime delle loss function, cioè guardano quali sono le direzioni negative nello spazio dei parametri. Un punto di partenza DEVE essere specificato, ed è detto $\bar{\theta}_0$.

Per ogni iterazione t , i parametri vengono aggiornati come

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \eta_t \bar{d}_t \quad \begin{array}{ll} \eta_t : \text{STEP SIZE} & \bar{d}_t : \text{DESCENT} \\ & (\text{o LEARNING RATE}) \end{array} \quad \text{DIRECTION}$$

$$\bar{g}_t = \nabla_{\bar{\theta}} \mathcal{L}(\bar{\theta})|_{\bar{\theta}_t} \sim \bar{d}_t \quad \text{STEEPEST DESCENT}$$

L'algoritmo viene iterato finché un punto stazionario viene raggiunto. Formalmente, richiederemo che esista $\eta_{\max} > 0$ / $\mathcal{L}(\bar{\theta} + \eta \bar{d}) \leq \mathcal{L}(\bar{\theta})$ $\forall 0 < \eta \leq \eta_{\max}$. Il gradiente \bar{g}_t punta nella direzione di massima variazione di $\mathcal{L}(\bar{\theta}_t)$. La scelta ovvia sembra quindi essere $\bar{g}_t = -\bar{d}_t$. La step size identifica, in sequenza $\{\eta_t\}$ il LEARNING RATE SCHEDULE.

La maniera più semplice è di scegliere η_t costante. Tuttavia, se η_t è troppo grande, il metodo potrebbe non convergere, se è troppo piccolo la convergenza potrebbe essere troppo lenta.

In linea generale è meglio scegliere η_t in modo adattivo, in modo che comporti la massima riduzione della loss function lungo la direzione scelta.

$$\eta_t = \underset{\eta > 0}{\operatorname{argmin}} \Phi_t(\eta) = \underset{\eta > 0}{\operatorname{argmin}} \mathcal{L}(\bar{\theta}_t + \eta \bar{d}_t)$$

Questo approccio è definito LINE SEARCH perché cerchiamo lungo la direzione identificata da \bar{d}_t .

Inoltre, vogliamo definire algoritmi che convergano rapidamente. Nel caso in cui la loss function sia convessa, si può mostrare che la gradient descent converge con tasso lineare.

$$|\mathcal{L}(\bar{\theta}_{t+1}) - \mathcal{L}(\hat{\theta}^*)| \leq \mu |\mathcal{L}(\bar{\theta}_t) - \mathcal{L}(\hat{\theta}^*)| \quad \text{con } \mu \text{ il tasso di convergenza}$$

$0 < \mu < 1$

Il metodo di gradient descent può essere molto inefficiente nel caso di direzioni piatte. Una tecnica, note come metodo dei momenti o delle palle pesanti, consiste nello spostarsi velocemente lungo direzioni "buone" (in base alle iterazioni precedenti) e lentamente in direzioni ove il gradiente cambia rapidamente.

$$\bar{m}_t = \beta \bar{m}_{t-1} + \bar{g}_{t-1} \quad \text{e} \quad \bar{v}_t = \bar{v}_{t-1} - \eta_t \bar{m}_t$$

dove \bar{m}_t è la "quantità di moto" e $0 < \beta < 1$ (tipicamente $\beta \sim 0,9$)

$$\bar{m}_t = \beta \bar{m}_{t-1} + \bar{g}_{t-1} = \beta^2 \bar{m}_{t-2} + \beta \bar{g}_{t-2} + \bar{g}_{t-1} = \dots = \sum_{\tau=0}^{t-1} \beta^\tau \bar{g}_{t-\tau-1}$$

Se tutti i gradienti nel passato sono costanti, allora

$$\bar{m}_t = \bar{g} \sum_{\tau=0}^{t-1} \beta^\tau \quad \text{e lo scaling factor è una serie geometrica}$$

$$\sum_{\tau=0}^{t-1} \beta^\tau \xrightarrow{t \rightarrow \infty} \frac{1}{1-\beta} \quad \text{se } \beta \sim 0,9, \text{ moltiplichiamo il gradiente per } 10$$

4.2 METODI DEL SECOND'ORDINE

Il classico metodo alle derivate seconde è il metodo di Newton.

Esso consiste nell'aggiornare i parametri come:

$$\vec{\vartheta}_{t+1} = \vec{\vartheta}_t - \eta_t H_t^{-1} \vec{g}_t \quad \text{con} \quad H_t \triangleq \nabla^2 \mathcal{L}(\vec{\vartheta})|_{\vec{\vartheta}_t} = H(\vec{\vartheta}_t) \quad \begin{matrix} \text{MATRICE} \\ \text{HESSIANA} \end{matrix}$$

Poiché l'inversa H_t^{-1} regolarizza la "skewness" locale delle curvature, il metodo di Newton è vantaggioso.

- 1 Inizializzare $\vec{\vartheta}_0$
- 2 finché $t=1, 2, \dots$, si arriva a convergenza
 - 3 valutare $\vec{g}_t = \nabla \mathcal{L}(\vec{\vartheta}_t)$
 - 4 valutare $H_t = \nabla^2 \mathcal{L}(\vec{\vartheta}_t)$
 - 5 risolvere $H_t \vec{d}_t = -\vec{g}_t$ per \vec{d}_t
 - 6 usare il line search per determinare η_t lungo \vec{d}_t
 - 7 $\vec{\vartheta}_{t+1} = \vec{\vartheta}_t + \eta_t \vec{d}_t$

Il metodo segue dalla Taylor expansion di $\mathcal{L}(\vec{\vartheta})$ attorno a $\vec{\vartheta}_t$

$$\mathcal{L}_{\text{quad}}(\vec{\vartheta}') = \mathcal{L}(\vec{\vartheta}_t) + \vec{g}_t^T (\vec{\vartheta}' - \vec{\vartheta}_t) + \frac{1}{2} (\vec{\vartheta}' - \vec{\vartheta}_t)^T H_t (\vec{\vartheta}' - \vec{\vartheta}_t)$$

Il minimo di $\mathcal{L}_{\text{quad}}$ è in $\vec{\vartheta}' = \vec{\vartheta}_t - H_t^{-1} \vec{g}_t$.

Esistono poi metodi "quasi-Newton" che generano un'approssimazione iterativa delle matrice Hessiana. Uno di questi è il metodo BFGS (Broyden, Fletcher, Goldfarb, Shanno) che approssima la matrice Hessiana con $B_t \approx H_t$

$$B_{t+1} = B_t + \frac{\vec{y}_t \vec{y}_t^T}{\vec{y}_t^T \vec{s}_t} - \frac{(B_t \vec{s}_t)(B_t \vec{s}_t)^T}{\vec{s}_t^T B_t \vec{s}_t} \quad \left\{ \begin{array}{l} \vec{s}_t = \vec{\vartheta}_t - \vec{\vartheta}_{t-1} \\ \vec{y}_t = \vec{g}_t - \vec{g}_{t-1} \end{array} \right.$$

e l'inversa

$$H_t^{-1} \approx C_t \quad \text{con} \quad C_{t+1} = \left(\mathbb{I} - \frac{\vec{s}_t \vec{y}_t^T}{\vec{y}_t^T \vec{s}_t} \right) C_t \left(\mathbb{I} - \frac{\vec{y}_t \vec{s}_t^T}{\vec{y}_t^T \vec{s}_t} \right) + \frac{\vec{s}_t \vec{s}_t^T}{\vec{y}_t^T \vec{s}_t}$$

4.3. STOCHASTIC GRADIENT DESCENT (SGD)

Lo scopo del metodo è di minimizzare la media

$$\mathcal{L}(\vec{\theta}) = \mathbb{E}_{q(\vec{z})} [\mathcal{L}(\vec{\theta}, \vec{z})]$$

dove \vec{z} è una variabile casuale usata in input. Ad ogni iterazione cerchiamo di osservare $\mathcal{L}_t(\vec{\theta}) = \mathcal{L}(\vec{\theta}, \vec{z}_t)$ con $\vec{z}_t \sim q$ (una funzione) se la distribuzione $q(\vec{z})$ è indipendente dai parametri che stiamo ottimizzando, possiamo usare $\vec{g}_t = \nabla_{\vec{\theta}} \mathcal{L}_t(\vec{\theta}_t)$. L'algoritmo è;

$$\vec{\theta}_{t+1} = \vec{\theta}_t - \eta_t \nabla \mathcal{L}(\vec{\theta}_t, \vec{z}_t) = \vec{\theta}_t - \eta_t \vec{g}_t$$

Ricordiamo che molte procedure di fit sono basate sull'empirical risk minimisation che implicano la minimizzazione di

$$\mathcal{L}(\vec{\theta}_t) = \frac{1}{N} \sum_{n=1}^N \ell(\vec{y}_n, f(\vec{x}_n; \vec{\theta}_t)) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(\vec{\theta}_t)$$

Segue che

$$\vec{g}_t = \frac{1}{N} \sum_{n=1}^N \nabla_{\vec{\theta}} \mathcal{L}_n(\vec{\theta}_t) = \frac{1}{N} \sum_{n=1}^N \nabla_{\vec{\theta}} \ell(\vec{y}_n, f(\vec{x}_n; \vec{\theta}_t)) \stackrel{\text{minibatch approximation}}{\approx} \frac{1}{|\mathcal{B}_t|} \sum_{n \in \mathcal{B}_t} \nabla_{\vec{\theta}} \mathcal{L}_n(\vec{\theta}_t)$$

dove $|\mathcal{B}_t|$ sono un set di esempi casuali utilizzati all'iterazione t . Una maniera per scegliere η_t è di iniziare con un valore piccolo e di incrementarlo con t , valutando la performance su minibatches.

Esempi: $\eta_t = \eta_i$ se $t_i \leq t \leq t_{i+1}$ piece-wise constant
 $\eta_t = \eta_0 e^{-\lambda t}$ exponential decay
 $\eta_t = \eta_0 (\beta t + 1)^{-\alpha}$ polynomial decay

Un'errata utile utilizzazione usa matrici di pre-condizionamento

$$\vec{\theta}_{t+1} = \vec{\theta}_t - \eta_t H_t^{-1} \vec{g}_t \quad \text{dove } H_t \text{ è positiva definita.}$$

Esempi di pre-conditioning

• ADAGRAD (adaptive gradient): utile se molti elementi del gradiente sono nulli (p. esempio se ci sono features nell'input che sono rare)

$$\vartheta_{t+1,d} = \vartheta_{t,d} - \eta_t \frac{1}{\sqrt{s_{t,d} + \epsilon}} g_{t,d}; \quad d=1, \dots, D \quad s_{t,d} = \sum_{i=1}^t g_{i,d}^2$$

ϵ piccolo; evita la divisione per zero

$$\Delta \bar{\vartheta}_t^2 = -\eta_t \frac{1}{\sqrt{s_t + \epsilon}} \bar{g}_t^2 \quad \leadsto \quad H_t = \text{diag}(\bar{s}_t^2 + \bar{\epsilon}^2)^{1/2}$$

• ADADELTA: come ADAGRAD, ma evita che il denominatore diventi troppo grande con il tempo, cosa che può compromettere l'efficienza.

$$\Delta \bar{\vartheta}_t^2 = -\eta_t \frac{\sqrt{s_{t-1} + \epsilon}}{\sqrt{s_t + \epsilon}} \bar{g}_t^2 \quad \text{con} \quad \bar{s}_t^2 = \beta \bar{s}_{t-1}^2 + (1-\beta)(\Delta \bar{\vartheta}_t^2)^2$$

• ADAM: combina ADADELTA con i momenti.

$$\begin{cases} \bar{m}_t^2 = \beta_1 \bar{m}_{t-1}^2 + (1-\beta_1) \bar{g}_t^2 \\ \bar{s}_t^2 = \beta_2 \bar{s}_{t-1}^2 + (1-\beta_2) \bar{g}_t^2 \end{cases} \quad \Delta \bar{\vartheta}_t^2 = -\eta_t \frac{1}{\sqrt{s_t + \epsilon}} \bar{m}_t^2$$

ADAM = adaptive moment estimation. $\beta_1 = 0,9$; $\beta_2 = 0,999$; $\epsilon = 10^{-6}$.

4.4. CONSTRAINED MINIMISATION

Ricordiamo il caso in cui

$$\hat{\vartheta}^* = \underset{\vartheta^* \in C}{\operatorname{argmin}} \mathcal{L}(\bar{\vartheta}^2) \quad \text{con} \quad C = \{ \bar{\vartheta}^2 \in \mathbb{R}^D : h_i(\bar{\vartheta}^2) = 0, i \in \mathcal{E}, g_j(\bar{\vartheta}^2) \leq 0, j \in \mathcal{I} \}$$

\mathcal{E} set di uguaglianze; \mathcal{I} set di disuguaglianze

Assumiamo di avere solamente $h(\bar{\vartheta}^2) = 0$. Notiamo che, per ogni punto sulla superficie $h(\bar{\vartheta}^2) = 0$, $\nabla h(\bar{\vartheta}^2)$ sarà ortogonale a $h(\bar{\vartheta}^2) = 0$. Questo è ovvio se si considera un punto $\bar{\vartheta}^2 + \bar{\epsilon}^2$ su $h(\bar{\vartheta}^2)$. Se facciamo un'espansione di Taylor al primo ordine attorno a $\bar{\vartheta}^2$, otteniamo

$$h(\bar{\vartheta}^2 + \bar{\epsilon}^2) \approx h(\bar{\vartheta}^2) + \bar{\epsilon}^T \nabla h(\bar{\vartheta}^2)$$

Poiché sia $\bar{\vartheta}^2$ sia $\bar{\vartheta}^2 + \bar{\epsilon}^2$ sono su $h(\bar{\vartheta}^2) = 0$, segue che $h(\bar{\vartheta}^2) = h(\bar{\vartheta}^2 + \bar{\epsilon}^2)$, e quindi $\bar{\epsilon}^T \nabla h(\bar{\vartheta}^2) \approx 0$. Poiché $\bar{\epsilon}^2 \parallel h(\bar{\vartheta}^2)$, $\Rightarrow \nabla h(\bar{\vartheta}^2) \perp h(\bar{\vartheta}^2) = 0$.

Cerchiamo un punto $\bar{\theta}^*$ in $h(\bar{\theta}^*) = 0$ tale che $\mathcal{L}(\bar{\theta}^*)$ è minimizzato.
 Questo punto dovrebbe essere tale che $\nabla h(\bar{\theta}^*) \perp h(\bar{\theta}^*) = 0$ e deve avere
 la proprietà per cui $\nabla \mathcal{L}(\bar{\theta}^*)$ è perpendicolare a $h(\bar{\theta}^*) = 0$. Poiché
 sia $\nabla h(\bar{\theta}^*)$ e $\nabla \mathcal{L}(\bar{\theta}^*)$ sono \perp a $h(\bar{\theta}^*) = 0$ in $\bar{\theta}^*$, devono essere
 paralleli. Esiste un λ^* costante tale che ($\lambda^* \in \mathbb{R}$):

$$\nabla \mathcal{L}(\bar{\theta}^*) = \lambda^* \nabla h(\bar{\theta}^*)$$

λ^* è chiamato moltiplicatore di Lagrange. Possiamo definire

$$L(\bar{\theta}, \lambda) \triangleq \mathcal{L}(\bar{\theta}) + \lambda h(\bar{\theta})$$

Un punto stazionario è tale che

$$\nabla_{\bar{\theta}, \lambda} L(\bar{\theta}, \lambda) = 0 \iff \lambda \nabla_{\bar{\theta}} h(\bar{\theta}) = \nabla \mathcal{L}(\bar{\theta}), \quad h(\bar{\theta}) = 0.$$

