



**UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO**

MOE: Mixture Of Experts

Analisi della letteratura e delle hierarchical MOE

Modellazione Statistica - Data Science

Ivan Diliso - 761053

Indice

1	Ricerca iniziale delle informazioni	3
2	Dettagli tecnici	3
2.1	Apprendimento	4
2.2	Hierarchical Mixture of Experts	4

1 Ricerca iniziale delle informazioni

MOE studiati nel campo delle reti neurali rispetto ad altri metodi ensemble come combining classifiers e ensemble of weak learners. Si utilizza un metodo dividi et impera per addestrare una serie di modelli parametrici e unirli per avere una soluzione. Se nei classici ensemble ogni learner è addestrato sullo stesso task in ME utilizza il divide et impera per dividere un task complesso in sottotask e ogni esperto è addestrato su task differenti, il modello di gating serve ad unire le soluzioni. A differenza dei modelli ensemble non c'è la necessità di rendere i learner individuali diversi in quanto ogni learner è addestrato per un task diverso. Il problema da risolvere è infatti trovare una divisione naturale dei dati. Una metodologia base è di targettare ogni esperto ad una diversa distribuzione specificata dalla funzione di gating, rispetto che apprendere la distribuzione originale dei dati.

Descrizione Tecnica di ensemble learning

Funzionamento di base Decompongo il task del modello predittivo in più sotto task, addestrare un esperto di quello specifico task su ogni task per poi sviluppare un gating model in grado di apprendere che esperto richiamare in base all'input e come combinare le predizioni. Posso suddividere il feature space di input in più feature space e addestrare un modello su ognuno di essi. Approccio divide et impera. I problemi possono essere sovrapponibili, non sovrapponibili e esperti su problemi simili collegati tra loro possono contribuire agli esempi che sono fuori dalla loro area di esperienza.

Questo approccio associa quindi un diverso peso ad ogni esperto, questa tecnica può essere vista come una forma di voting dei modelli ensemble, dove però la capacità di voto può cambiare al variare dell'input.

I pesi determinati dal gating network sono assegnati dinamicamente al variare dell'input, MOE quindi apprende che porzione del feature space è learned da ogni esperto dell'ensemble. I classificatori individuali sono addestrati per diventare esperti in una porzione del feature space. La funzione di gating quindi seleziona il classificatore, pesato con la sua expertise utilizzare per ogni istanza.

- POOLING: Utilizzo solo il classificatore con il peso più alto
- COMBINING: Utilizzo una somma pesata degli output di tutti i classificatori

2 Dettagli tecnici

Combinazione dei modelli, tipologia multi expert (diversi learner che lavorano in parallelo) con approccio locale (learner selection) si utilizza un modello di gating che guarda l'input e sceglie che modello è responsabile per generare l'output.

Sia $x \in \mathbb{R}^n$ vettore di input e si T il numero di esperti modello e h_1, \dots, h_T gli esperti del modello e y variabile target. Dati W_i parametri dell' i -esimo esperto, questo prova ad approssimare la distribuzione di y

$$h_i(y|x; W_i)$$

La funzione di gating produce un set di coefficienti che pesano il contributo degli esperti, sia v_i vettore dei pesi della funzione di gating relativa all' i -esimo esperto e α parametro del modello di gating, insieme dei pesi relativi ad ogni esperto, il set di coefficienti prodotti dal gating:

$$\pi_i(x; \alpha) : \sum_{i=1}^T \pi_i(x; \alpha) = 1$$

Sulla base di queste probabilità partizioniamo lo spazio di input, diverse partizioni appartengono a diversi esperti. L'output del modello sarà quindi:

$$H(y|x; \psi) = \sum_{i=1}^T \pi_i(x; \alpha) \cdot h_i(y|x; W_i)$$

Nella fase di training il valore $\pi_i(x; \alpha)$ indica la probabilità che l'istanza x appaia nel training set dell' i -esimo esperto. Mentre nella fase di testing definisce il contributo che h_i dà alla predizione finale. L'output della funzione di gating può essere espresso tramite una softmax

$$\pi_i(x; \alpha) = \frac{e^{v_i x}}{\sum_{l=1}^k e^{v_l x}}$$

2.1 Apprendimento

Può avvenire tramite:

- GRADIENT DESCENT
- EXPECTATION MAXIMIZATION

2.2 Hierarchical Mixture of Experts

Rimpiazzo ogni esperto con un sistema completo MOE in modo ricorsivo. Si decide la profondità della ricorsione, il tipo di esperto e il tipo di modello di gating. Questo sviluppo ricorsivo crea una struttura ad albero. Può essere interpretato come un albero di decisione con i gating model che definiscono i nodi di decisione. Questa tipologia di albero viene definita “soft decision tree” in quanto i gating model ritornano una distribuzione di probabilità sugli esperti vengono quindi esplorate tutte le path dell'albero con differenti probabilità prendendo poi una somma pesata a livello di foglie dove il prodotto è uguale al prodotto dei valori di gating di ogni path per arrivare alla foglia. In questa tipologia di apprendimento ogni nodo implementa un modello lineare (o regressione logistica) invece del valore costante di un albero CART. Nodi terminali chiamati esperti e nodi non terminali sono i nodi di gating. L'idea è che ogni esperto dà una opinione sulla predizione e queste sono combinate dal modello di gating.

Vantaggi Le boundaries tra regioni di foglie non sono più “hard” ma c'è una transizione graduale da una all'altra, portando ad uno smoothing della risposta.

L'uso di soft split permette di catturare situazioni in cui la transizione da una risposta alta a bassa è graduale.

Riferimenti bibliografici

- [1] Hastie, T., Tibshirani, R., & Friedman, J. (2009). The elements of statistical learning: data mining, inference, and prediction. Springer Science & Business Media.
- [2] Alpaydin, E. (2020). Introduction to machine learning. MIT press.
- [3] Zhou, Z. H. (2012). Ensemble methods: foundations and algorithms. CRC press.
- [4] Zhang, C., & Ma, Y. (Eds.). (2012). Ensemble machine learning: methods and applications. Springer Science & Business Media.
- [5] Bishop, C. M., & Nasrabadi, N. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: springer.