

RETE NEURALE RICORRENTE PER LA FATTORIZZAZIONE NON NEGATIVA DI MATRICI

Renzo Perfetti¹, Giovanni Costantini², Massimiliano Todisco²

¹Dipartimento di Ingegneria, Università di Perugia

²Dipartimento di Ingegneria Elettronica, Università di Roma Tor Vergata

In diversi ambiti relativi all'analisi dei dati, sia supervisionata (es. classificazione) sia non supervisionata (es. clustering), un approccio molto interessante è rappresentato dai *subspace methods*. Tali metodi proiettano il vettore delle features in un sottospazio, rappresentato da una base opportuna, in modo da ottenere una rappresentazione più significativa o più adatta alla successiva elaborazione. In questo modo lo spazio originale delle features, con un numero elevato di dimensioni, risulta proiettato in un sottospazio vettoriale di dimensioni ridotte, in cui la classificazione o la ricerca di somiglianze tra i dati risulta facilitata.

Diversi criteri sono stati applicati per ottenere le basi degli spazi a dimensione ridotta. Il più noto consiste nel minimizzare l'errore quadratico di rappresentazione e conduce alla ben nota tecnica denominata PCA (*Principal Component Analysis*). La NMF (*Nonnegative Matrix Factorization*), come la PCA, rappresenta i dati come combinazioni lineari dei vettori di una base di dimensione ridotta tuttavia, a differenza della PCA, non permette valori negativi sia nelle componenti dei vettori di base sia nei coefficienti della combinazione lineare. Tale vincolo conduce a basi completamente differenti rispetto alla PCA; ad esempio, nel caso dei volti umani le basi della NMF rappresentano caratteristiche localizzate, ovvero parti significative come occhi, bocca, naso, ecc [1].

Il problema della NMF può essere formulato come un problema di ottimizzazione: data la matrice \mathbf{V} ($m \times n$) non negativa, dobbiamo determinare due matrici *non negative* \mathbf{W} ($m \times p$) e \mathbf{H} ($p \times n$), con $p < \min(m, n)$, tali che sia minima la funzione costo:

$$J(\mathbf{W}, \mathbf{H}) = \|\mathbf{WH} - \mathbf{V}\|^2 \quad (1)$$

In pratica si ha $p \ll \min(m, n)$. Dalla (1) è evidente che la NMF è in realtà una fattorizzazione non negativa *approssimata*. L'intero p è detto rango della fattorizzazione.

Nelle applicazioni pratiche le n colonne di \mathbf{V} sono i dati (i vettori delle features con m componenti) e le colonne di \mathbf{W} sono i vettori di base; le colonne di \mathbf{H} rappresentano i coefficienti che esprimono i dati nel nuovo spazio con p dimensioni.

Il problema di ottimizzazione della NMF non è convesso pertanto presenta il noto repertorio di inconvenienti, in particolare esistenza di soluzioni locali multiple e dipendenza della soluzione trovata dalle condizioni iniziali quando si utilizzano metodi iterativi. Inoltre è affetto da una intrinseca indeterminazione; infatti, la funzione costo non cambia considerando al posto di \mathbf{W} e \mathbf{H} le matrici \mathbf{WD} e $\mathbf{D}^{-1}\mathbf{H}$, essendo \mathbf{D} una matrice qualsiasi purché invertibile e non negativa. Infine, è noto che la NMF senza ulteriori vincoli produce generalmente rappresentazioni poco interessanti dei dati, in particolare non garantisce una scomposizione in parti significative [2-4]. Per evitare tali inconvenienti, alcuni autori hanno suggerito di normalizzare le colonne di \mathbf{W} (i vettori di base) con la norma L_1 per ottenere soluzioni 'sparse'. Poiché le componenti delle matrici sono non negative la norma L_1 è semplicemente la somma delle componenti. Nel caso del clustering spesso si ottengono risultati migliori normalizzando le righe della matrice \mathbf{H} (i coefficienti della rappresentazione).

In [5] si propone una rete neurale ricorrente che risolve il problema della NMF vincolata. Le variabili di stato della rete rappresentano le componenti di \mathbf{W} ed \mathbf{H} . La rete minimizza una funzione Lagrangiana che combina la (1) con i vincoli di sparsità sulle colonne di \mathbf{W} o sulle righe di \mathbf{H} . I moltiplicatori di Lagrange, che rappresentano i pesi dei vincoli, sono calcolati direttamente dalla rete durante l'evoluzione dinamica, evitando il ricorso alla validazione o al tuning empirico. Nell'articolo si studia la convergenza della rete, mostrando che le soluzioni locali del problema NMF corrispondono ad altrettanti stati di equilibrio asintoticamente stabili. La validità della rete proposta è illustrata da numerose simulazioni riguardanti sia l'estrazione di features sia il clustering. A titolo di esempio, in Fig. 1 sono mostrate quattro delle 256 immagini del dataset Swimmer e, in Fig. 2, i vettori di base (colonne di \mathbf{W}) calcolati dalla rete neurale, assumendo $p = 20$. La rete individua correttamente le 16 possibili posizioni delle braccia e quattro repliche del torso.



Fig. 1

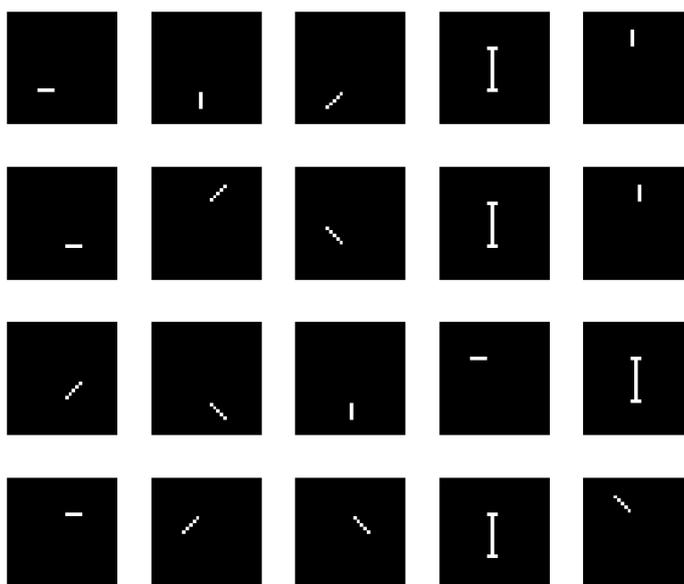


Fig. 2

BIBLIOGRAFIA

- [1] Berry M. W., Browne M., Langville A. N., Pauca V. P., and Plemmons R. J., 2006. "Algorithms and Applications for Approximate Nonnegative Matrix Factorization", Computational Statistics and Data Analysis, vol. 52, pp. 155-173.
- [2] D. Donoho, V. Stodden, "When does non-negative matrix factorization give a correct decomposition into parts?", 16 (2003) 1141–1148Adv. Neural Inf. Process. Syst. 16, 2003, 1141–1148.
- [3] Hoyer P., 2004. "Non-negative Matrix Factorization with Sparseness Constraints". J. of Machine Learning Research 5, 1457–1469.
- [4] Choi S., 2008. "Algorithms for Orthogonal Nonnegative Matrix Factorization", 2008. IEEE International Joint Conference on Neural Networks, Hong Kong, 1-8 June, pp. 1828-1832.
- [5] G. Costantini, R. Perfetti, M. Todisco, "Recurrent neural network for approximate nonnegative matrix factorization", Neurocomputing vol. 138, 2014, pp. 238-247.
<http://dx.doi.org/10.1016/j.neucom.2014.02.007>