

La tecnica del conteggio approssimato permette di stimare in modo efficiente la numerosità degli elementi più frequenti in una collezione. Vediamo ora come la tecnica delle proiezioni casuali ci permette di stimare in modo efficiente le distanze fra coppie di punti nello spazio Euclideo d -dimensionale quando d è grande. Ricordiamo che la distanza Euclidea fra due punti $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ è calcolata come

$$\|\mathbf{x} - \mathbf{x}'\| = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}.$$

In molte applicazioni i dati possono essere rappresentati come vettori di numeri. Due esempi importanti sono le immagini (ogni coordinata è un pixel) e i testi (ogni coordinata è una parola del dizionario e il valore della coordinata è la frequenza con la quale la parola compare nel testo). Se consideriamo l'elenco dei film disponibili su Netflix come un dizionario, allora anche un utente di Netflix può essere visto come un vettore di numeri dove ogni coordinata è un film e il valore della coordinata rappresenta una valutazione del film da parte dell'utente.

In tutti questi casi, possiamo interpretare la vicinanza di due punti in \mathbb{R}^d come una misura della similarità fra gli elementi (immagini, testi, utenti) che i punti rappresentano. Quindi, la capacità di calcolare in modo efficiente qual è il punto in un insieme più vicino ad un dato punto (*nearest neighbor*) diventa fondamentale per, ad esempio, suggerire film a nuovi utenti basandosi sui film apprezzati da utenti che hanno un profilo simile (ovvero, le loro codifiche in \mathbb{R}^d sono vicine in termini di distanza Euclidea).

Problema **Nearest neighbor**.

Istanza: Un insieme finito $S \subset \mathbb{R}^d$ e un punto $\mathbf{x} \in \mathbb{R}^d$.

Soluzione: $\operatorname{argmin}_{\mathbf{x}' \in S} \|\mathbf{x} - \mathbf{x}'\|$.

Putroppo, trovare il nearest neighbor in d dimensioni diventa computazionalmente costoso quando $d \gg 1$, come di solito succede nelle applicazioni interessanti. Per esempio, se $|S| = n$ e voglio risolvere il problema nearest neighbor calcolando le distanze fra \mathbf{x} e i punti di S impiegherò un tempo dell'ordine di nd . Se devo risolvere il problema ogni volta che aggiungo un nuovo utente ad S impiegherò quindi un tempo dell'ordine di $\sum_{t=1}^n (td) = \Theta(n^2d)$.

Per ovviare a questo problema mostriamo che per ogni $0 < \varepsilon, \delta < 1$ esiste $k = \mathcal{O}\left(\frac{1}{\varepsilon^2} \ln \frac{|S|}{\delta}\right)$ ed esiste una classe \mathcal{F} di funzioni $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$ tale che

$$(1 - \varepsilon) \|\mathbf{x} - \mathbf{x}'\|^2 \leq \|f(\mathbf{x}) - f(\mathbf{x}')\|^2 \leq (1 + \varepsilon) \|\mathbf{x} - \mathbf{x}'\|^2 \quad \mathbf{x}, \mathbf{x}' \in S$$

con probabilità almeno $1 - \delta$ rispetto all'estrazione di $f \in \mathcal{F}$.

Per dimostrare questo risultato utilizziamo una tecnica simile al conteggio approssimato. Ovvero, usiamo k funzioni casuali che sono analoghe alle funzioni hash del conteggio approssimato.

Queste funzioni sono rappresentate da k vettori casuali $\mathbf{Z}_1, \dots, \mathbf{Z}_k \in \mathbb{R}^d$ estratti in un modo che spiegheremo tra breve. La funzione casuale associata al vettore \mathbf{Z}_j è definita come

$$f_j(\mathbf{x}) = \mathbf{Z}_j^\top \mathbf{x} = \sum_{i=1}^d Z_{j,i} x_i .$$

Il prodotto scalare $\mathbf{Z}_j^\top \mathbf{x}$ calcola la lunghezza della proiezione di \mathbf{x} su \mathbf{Z}_j moltiplicata per la lunghezza di \mathbf{Z}_j . Quindi, usando la funzione f_j possiamo approssimare la distanza $\|\mathbf{x} - \mathbf{x}'\|$ fra i vettori \mathbf{x} e \mathbf{x}' con la differenza $|f_j(\mathbf{x}) - f_j(\mathbf{x}')|$ fra numeri reali $f_j(\mathbf{x})$ e $f_j(\mathbf{x}')$. Per ridurre l'errore di approssimazione utilizziamo k funzioni indipendenti invece di una sola.

I vettori \mathbf{Z}_j sono ottenuti generando ciascuna componente $Z_{j,i}$ per $i = 1, \dots, d$ con estrazioni indipendenti da una distribuzione di probabilità con media zero e varianza uno, cioè

$$\mathbb{E}[Z_{j,i}] = 0 \quad \text{e} \quad \text{Var}[Z_{j,i}] = 1 \quad j = 1, \dots, k \quad i = 1, \dots, d .$$

Quindi,

$$\begin{aligned} \mathbb{E}\left[(f_j(\mathbf{x}) - f_j(\mathbf{x}'))^2\right] &= \mathbb{E}\left[\left(\sum_{i=1}^d (x_i - x'_i) Z_{j,i}\right)^2\right] \\ &= \mathbb{E}\left[\sum_{r=1}^d \sum_{s=1}^d (x_r - x'_r)(x_s - x'_s) Z_{j,r} Z_{j,s}\right] \\ &= \mathbb{E}\left[\sum_{i=1}^d (x_i - x'_i)^2 Z_{j,i}^2\right] + \mathbb{E}\left[\sum_{r,s:r \neq s} (x_r - x'_r)(x_s - x'_s) Z_{j,r} Z_{j,s}\right] \\ &= \sum_{i=1}^d (x_i - x'_i)^2 \mathbb{E}[Z_{j,i}^2] + \sum_{r,s=1}^d (x_r - x'_r)(x_s - x'_s) \mathbb{E}[Z_{j,r}] \mathbb{E}[Z_{j,s}] \\ &= \sum_{i=1}^d (x_i - x'_i)^2 \text{Var}[Z_{j,i}] \\ &= \sum_{i=1}^d (x_i - x'_i)^2 = \|\mathbf{x} - \mathbf{x}'\|^2 . \end{aligned} \tag{1}$$

dove (1) vale perché le $Z_{j,i}$ hanno media zero e quindi

$$\text{Var}[Z_{j,i}] = \mathbb{E}\left[(Z_{j,i} - \mathbb{E}[Z_{j,i}])^2\right] = \mathbb{E}[Z_{j,i}^2] .$$

Questo dimostra che posso usare $(f_j(\mathbf{x}) - f_j(\mathbf{x}'))^2$ per stimare la distanza quadrata $\|\mathbf{x} - \mathbf{x}'\|^2$.

Definiamo ora la proiezione casuale $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$,

$$f(\mathbf{x}) = \left(\frac{f_1(\mathbf{x})}{\sqrt{k}}, \dots, \frac{f_k(\mathbf{x})}{\sqrt{k}}\right) .$$

Si noti che $f(\mathbf{x}) = M\mathbf{x}$ dove M è la matrice casuale $k \times d$ avente $\mathbf{Z}_1/\sqrt{k}, \dots, \mathbf{Z}_k/\sqrt{k}$ come righe. Questo implica che f è una trasformazione lineare, ovvero $f(a\mathbf{x} + b\mathbf{x}') = af(\mathbf{x}) + bf(\mathbf{x}')$ per ogni $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ e $a, b \in \mathbb{R}$. Quindi,

$$\mathbb{E} \left[\|f(\mathbf{x}) - f(\mathbf{x}')\|^2 \right] = \mathbb{E} \left[\|f(\mathbf{x} - \mathbf{x}')\|^2 \right] = \frac{1}{k} \sum_{j=1}^k \mathbb{E} [f_j(\mathbf{x} - \mathbf{x}')^2] = \|\mathbf{x} - \mathbf{x}'\|^2$$

dato che ciascun f_j è uno stimatore della distanza quadrata.

Ora, detto $\mathbf{v} = \mathbf{x} - \mathbf{x}'$ e usando sempre il fatto che f è lineare,

$$\frac{\|f(\mathbf{x}) - f(\mathbf{x}')\|^2}{\|\mathbf{x} - \mathbf{x}'\|^2} = \left\| f \left(\frac{\mathbf{v}}{\|\mathbf{v}\|} \right) \right\|^2.$$

Quindi, se vogliamo dimostrare che

$$(1 - \varepsilon) \|\mathbf{x} - \mathbf{x}'\|^2 \leq \|f(\mathbf{x}) - f(\mathbf{x}')\|^2 \leq (1 + \varepsilon) \|\mathbf{x} - \mathbf{x}'\|^2 \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$$

possiamo equivalentemente dimostrare che $1 - \varepsilon \leq \|f(\mathbf{v})\|^2 \leq 1 + \varepsilon$ per ogni $\mathbf{v} \in \mathbb{R}^d$ tale che $\|\mathbf{v}\| = 1$.

A questo punto ci serve fare assunzioni sulla distribuzione delle variabili casuali $Z_{j,i}$. Assumiamo quindi che le $Z_{j,i}$ abbiano una distribuzione Normale (ovvero, Gaussiana con media zero e varianza uno). Per le proprietà della Normale, vale che per ogni $\varepsilon, \delta > 0$ fissati e per ogni $\mathbf{v} \in \mathbb{R}^d$ di norma unitaria,

$$\mathbb{P} \left(\left| \|f(\mathbf{v})\|^2 - 1 \right| > \varepsilon \right) \leq \delta \quad \text{per } k = \mathcal{O} \left(\frac{1}{\varepsilon^2} \ln \frac{1}{\delta} \right) \quad (2)$$

dove la probabilità è calcolata rispetto all'estrazione delle $\{Z_{j,i} : j = 1, \dots, k, i = 1, \dots, d\}$.

Si noti che

$$\|f(\mathbf{v})\|^2 = \frac{1}{k} \sum_{j=1}^k (\mathbf{Z}_j^\top \mathbf{v})^2$$

e inoltre

$$\mathbb{E} \left[(\mathbf{Z}_j^\top \mathbf{v})^2 \right] = \mathbf{v}^\top \mathbb{E} [\mathbf{Z}_j \mathbf{Z}_j^\top] \mathbf{v} = \mathbf{v}^\top I \mathbf{v} = \|\mathbf{v}\|^2 = 1$$

dove abbiamo usato il fatto che la matrice $M = \mathbb{E} [\mathbf{Z}_j \mathbf{Z}_j^\top]$ ha componenti $M_{r,s} = \mathbb{E} [Z_{j,r} Z_{j,s}]$ tali che

$$M_{r,s} = \begin{cases} 0 & \text{se } r \neq s \\ 1 & \text{altrimenti} \end{cases}$$

Quindi le variabili casuali $V_j = (\mathbf{Z}_j^\top \mathbf{v})^2$ per $j = 1, \dots, k$ sono i.i.d. con media $\mu = 1$ e (2) può essere riscritta come

$$\mathbb{P} \left(\left| \frac{1}{k} \sum_{j=1}^k V_j - \mu \right| > \varepsilon \right) \leq e^{-\mathcal{O}(k\varepsilon^2)}$$

Questa diseguaglianza è analoga al Lemma di Chernoff-Hoeffding, con l'unica differenza che qui le V_j non hanno valori limitati. La formula (2) ci dice quindi che un risultato analogo al Lemma di Chernoff-Hoeffding vale anche per variabili casuali del tipo $(\mathbf{Z}_j^\top \mathbf{v})^2$ dove \mathbf{Z}_j sono Normali multivariate e $\|\mathbf{v}\| = 1$.

Per capire i prossimi passaggi ricordiamo che, per qualsiasi insieme di eventi A_1, \dots, A_N vale che

$$\mathbb{P}(\exists i : A_i) = \mathbb{P}(A_1 \cup \dots \cup A_N) \leq \sum_{i=1}^N \mathbb{P}(A_i).$$

Nel nostro caso, ci interessano gli eventi

$$A_{\mathbf{x}, \mathbf{x}'} = \left| \frac{\|f(\mathbf{x}) - f(\mathbf{x}')\|^2}{\|\mathbf{x} - \mathbf{x}'\|^2} - 1 \right| > \varepsilon$$

per ognuna delle $N = \binom{n}{2} \leq n^2$ coppie di punti distinti $\mathbf{x}, \mathbf{x}' \in S$. Allora, dato un qualunque insieme $S \subset \mathbb{R}^d$ di n punti,

$$\mathbb{P}(\exists \mathbf{x}, \mathbf{x}' \in S : A_{\mathbf{x}, \mathbf{x}'}) = \mathbb{P}\left(\bigcup_{\mathbf{x}, \mathbf{x}' \in S} A_{\mathbf{x}, \mathbf{x}'}\right) \leq \sum_{\mathbf{x}, \mathbf{x}' \in S} \mathbb{P}(A_{\mathbf{x}, \mathbf{x}'}) \leq \sum_{\mathbf{x}, \mathbf{x}' \in S} \delta \leq n^2 \delta$$

per $k = \mathcal{O}\left(\frac{1}{\varepsilon^2} \ln \frac{1}{\delta}\right)$.

Da questo ne deduciamo che, per $k = \mathcal{O}\left(\frac{1}{\varepsilon^2} \ln \frac{n}{\delta}\right)$ vale

$$(1 - \varepsilon) \|\mathbf{x} - \mathbf{x}'\|^2 \leq \|f(\mathbf{x}) - f(\mathbf{x}')\|^2 \leq (1 + \varepsilon) \|\mathbf{x} - \mathbf{x}'\|^2 \quad \text{per ogni } \mathbf{x}, \mathbf{x}' \in S \quad (3)$$

con probabilità almeno $1 - \delta$ rispetto all'estrazione delle $\{Z_{j,i} : j = 1, \dots, k, i = 1, \dots, d\}$.

Se ci accontentiamo di un errore nella stima delle distanze del 10% con probabilità del 99% rispetto all'estrazione di tutte le $Z_{j,i}$, allora ε e δ sono costanti e quindi $k = \mathcal{O}(\log n)$. Il costo per mappare i punti di S in \mathbb{R}^k è $ndk = nd \ln n$ e il costo per calcolare le coppie di distanze fra un \mathbf{x} e i punti in S è $n \ln n$. Se devo risolvere il problema nearest neighbor approssimato n volte impiegherò quindi un tempo dell'ordine di $nd \ln n + n^2 \ln n \leq n^2 d$ quando $n = \mathcal{O}(2^d)$.

Se avessimo al più $s < k$ valori non nulli in ciascuna colonna di M , allora il costo per mappare un punto di S in \mathbb{R}^k sarebbe ds . È possibile dimostrare che (3) vale per $k = \mathcal{O}\left(\frac{1}{\varepsilon^2} \ln \frac{n}{\delta}\right)$ e $s = \mathcal{O}\left(\frac{1}{\varepsilon} \ln \frac{n}{\delta}\right)$.