

Lezione 6: Modelli lineari - classificazione

In questa lezione consideriamo modelli di classificazione nella forma

$$p(y=c/\vec{x}; \vec{\theta}) = \frac{p(\vec{x}/y=c; \vec{\theta}) p(y=c; \vec{\theta})}{\sum_c p(\vec{x}/y=c; \vec{\theta}) p(y=c; \vec{\theta})}$$

Il termine $p(y=c; \vec{\theta})$ è il prior sulle class labels ed il termine $p(\vec{x}/y=c; \vec{\theta})$ è chiamato CLASS CONDITIONAL DENSITY per la classe c . Il modello è chiamato "GENERATIVE CLASSIFIER" perché specifica un modo per generare le features \vec{x} \forall classe c campionando $p(\vec{x}/y=c; \vec{\theta})$. Invece, un "DISCRIMINATIVE CLASSIFIER" modella direttamente la classe posteriore $p(y/\vec{x}; \vec{\theta})$. Scegliendo le probabilità condizionate in una certa maniera, il posteriore è una funzione lineare in \vec{x} ($\ln p(y=c/\vec{x}; \vec{\theta}) = \vec{w}^T \vec{x} + \text{const.}$). Il metodo è chiamato LINEAR DISCRIMINANT ANALYSIS (LDA).

5.1 GAUSSIAN DISCRIMINANT ANALYSIS (GDA)

Consideriamo che la distribuzione di probabilità condizionate sulle classe sia una gaussiana

$$p(\vec{x}/y=c; \vec{\theta}) = \mathcal{N}(\vec{x}/\vec{\mu}_c, \vec{\Sigma}_c)$$

Il corrispondente posteriore è pure gaussiano

$$p(y=c/\vec{x}; \vec{\theta}) \propto \pi_c \mathcal{N}(\vec{x}/\vec{\mu}_c, \vec{\Sigma}_c)$$

dove $\pi_c = p(y=c)$ è la prior probability per la classe c . Prendendo il \ln :

$$\ln p(y=c/\vec{x}; \vec{\theta}) = \ln \pi_c - \frac{1}{2} \ln |2\pi \vec{\Sigma}_c| - \frac{1}{2} (\vec{x} - \vec{\mu}_c)^T \vec{\Sigma}_c^{-1} (\vec{x} - \vec{\mu}_c) + \text{costante}$$

Questa funzione è chiamata FUNZIONE DISCRIMINANTE. Il DECISION BOUNDARY tra due classi (p.es. c e c') è una funzione quadratica di \vec{x} . Questo modello è noto come QUADRATIC DISCRIMINANT ANALYSIS (QDA). Consideriamo ora un caso specifico di GDA in cui le matrici di covarianza siano condivise tra diverse classi, cioè $\vec{\Sigma}_c = \vec{\Sigma}$. Se $\vec{\Sigma}$

è indipendente da c , riscriviamo:

$$\begin{aligned} \ln p(y=c/\bar{x}^2, \bar{y}^2) &= \ln \pi_c - \frac{1}{2} (\bar{x}^2 - \bar{\mu}_c^2)^T \bar{\Sigma}_c^{-1} (\bar{x}^2 - \bar{\mu}_c^2) + \text{const} \\ &= \underbrace{\ln \pi_c - \frac{1}{2} \bar{\mu}_c^2^T \bar{\Sigma}_c^{-1} \bar{\mu}_c^2}_{\gamma_c} + \underbrace{\bar{x}^2^T \bar{\Sigma}_c^{-1} \bar{\mu}_c^2}_{\beta_c} + \underbrace{\text{const} - \frac{1}{2} \bar{x}^2^T \bar{\Sigma}_c^{-1} \bar{x}^2}_{K} \\ &= \gamma_c + \bar{x}^2^T \beta_c + K \quad \text{LINEAR DISCRIMINANT ANALYSIS (LDA)} \end{aligned}$$

Verifichiamo ora come fitting un modello LDA usando la MLE. La likelihood ha la forma

$$p(\mathcal{D}/\bar{y}^2) = \prod_{n=1}^{N_d} \text{Cat}(y_n/\bar{\pi}^2) \prod_{c=1}^C \mathcal{N}(\bar{x}_n^2/\bar{\mu}_c^2, \bar{\Sigma}_c^2)^{\mathbb{1}(y_n=c)}$$

$$\ln p(\mathcal{D}/\bar{y}^2) = \left[\sum_{n=1}^{N_d} \sum_{c=1}^C \mathbb{1}(y_n=c) \ln \pi_c \right] + \sum_{c=1}^C \left[\sum_{n: y_n=c} \left[\bar{\Sigma}_c^2 \ln \mathcal{N}(\bar{x}_n^2/\bar{\mu}_c^2, \bar{\Sigma}_c^2) \right] \right]$$

quindi è possibile ottimizzare separatamente π e $(\mu_c, \bar{\Sigma}_c^2)$. Sappiamo dalle lezioni precedenti che

$$\hat{\mu}_c = \frac{1}{N_{d,c}} \sum_{n: y_n=c} \bar{x}_n^2 \quad \hat{\Sigma}_c^2 = \frac{1}{N_{d,c}} \sum_{n: y_n=c} (\bar{x}_n^2 - \hat{\mu}_c^2)(\bar{x}_n^2 - \hat{\mu}_c^2)^T$$

Osservazioni:

a) se $\bar{\Sigma}_c^2 = \bar{\Sigma}^2$, segue che $\hat{\Sigma}_c^2 = \frac{1}{N_d} \sum_{c=1}^C \sum_{n: y_n=c} (\bar{x}_n^2 - \hat{\mu}_c^2)(\bar{x}_n^2 - \hat{\mu}_c^2)^T$

b) se forziamo $\bar{\Sigma}_c^2$ ad essere diagonale, riduciamo il numero di parametri da $O(CD^2)$ a $O(CD)$. Questo è chiamato naive Bayes model.

Questa assunzione è alquanto restrittiva, ma funziona bene per C, D grandi

c) se assumiamo un prior uniforme sulle classi, possiamo scegliere la class label più probabile come segue

$$\hat{y}(\bar{x}^2) = \arg\max_c \ln p(y=c/\bar{x}^2, \bar{y}^2) = \arg\min_c (\bar{x}^2 - \bar{\mu}_c^2)^T \bar{\Sigma}_c^{-1} (\bar{x}^2 - \bar{\mu}_c^2)$$

Questo metodo è chiamato NEAREST CENTROID CLASSIFIER poiché assegna \bar{x}^2 alla classe con il più vicino $\bar{\mu}_c^2$

6.2 NAIVE BAYES CLASSIFIERS (NBC)

Supponiamo che le features siano condizionalmente indipendenti date le class labels. L'assunzione corrisponde ad usare una probabilità condizionata delle forme

$$p(\vec{x}/y=c, \vec{\theta}) = \prod_{d=1}^D p(x_d/y=c, \vec{\theta}_{dc})$$

dove $\vec{\theta}_{dc}$ sono i parametri per la distribuzione di probabilità relative alla classe c ed alle feature d . Quindi:

$$p(y=c/\vec{x}, \vec{\theta}) = \frac{p(y=c/\pi) \prod_{d=1}^D p(x_d/y=c, \vec{\theta}_{dc})}{\sum_c p(y=c/\pi) \prod_{d=1}^D p(x_d/y=c, \vec{\theta}_{dc})}$$

dove π_c è il prior per la classe c e $\vec{\theta} = (\pi, \{\vec{\theta}_{dc}\})$ sono i parametri. Dobbiamo specificare la forma delle varie distribuzioni di probabilità.

- Nel caso di features binarie, $x_d \in \{0, 1\}$ possiamo usare la distribuzione di Bernoulli: $p(\vec{x}/y=c, \vec{\theta}) = \prod_{d=1}^D \text{Ber}(x_d/\theta_{dc})$ con θ_{dc} la probabilità che $x_d=1$ nella classe c .
- Nel caso di features categoriche, $x_d \in \{1, \dots, K\}$, possiamo usare la distribuzione categorica $p(\vec{x}/y=c, \vec{\theta}) = \prod_{d=1}^D \text{Cat}(x_d/\vec{\theta}_{dc})$ con $\theta_{dc,k}$ le probabilità che $x_d=k$ dato $y=c$.
- Nel caso di features a valori reali, $x_d \in \mathbb{R}$, possiamo usare la distribuzione gaussiana $p(\vec{x}/y=c, \vec{\theta}) = \prod_{d=1}^D \mathcal{N}(x_d/\mu_{dc}, \sigma_{dc}^2)$ con μ_{dc} le medie delle feature d quando la class label è c e σ_{dc}^2 è la sua varianza (questo caso è equivalente alla GDA con covarianze diagonali).

Scriviamo ora un NBC usando la MLE

$$p(\mathcal{D}/\vec{\theta}) = \prod_{n=1}^{N_D} \text{Cat}(y_n/\pi) \prod_{d=1}^D p(x_{nd}/y_n, \vec{\theta}_d)$$

$$= \prod_{n=1}^{N_D} \text{Cat}(y_n/\pi) \prod_{d=1}^D \prod_{c=1}^C p(x_{nd}/\vec{\theta}_{dc})^{\mathbb{I}(y_n=c)}$$

$$\ln p(\mathcal{D}/\vec{\theta}) = \left[\sum_{n=1}^{N_D} \sum_{c=1}^C \mathbb{I}(y_n=c) \ln \pi_c \right] + \sum_{c=1}^C \sum_{d=1}^D \left[\sum_{n: y_n=c} \ln p(x_{nd}/\vec{\theta}_{dc}) \right]$$

Questa equazione si scompone in due contributi

$$\ln p(\mathcal{D}/\vec{\theta}) = \ln p(\mathcal{D}_y/\pi) + \sum_c \sum_d \ln p(\mathcal{D}_{cd}/\vec{\theta}_{dc})$$

$\mathcal{D}_y = \{y_n : n=1, N\}$ LABELS e $\mathcal{D}_{dc} = \{x_{nd} : y_n = c\}$ VALORI delle FEATURE
ab. per samples di classe.

Nelle lezioni precedenti, abbiamo mostrato che la MLE per π è il vettore $\hat{\pi}_c = N_c/N$. Gli altri parametri sono stimati analogamente

• features binarie: $\hat{\theta}_{dc} = N_{dc}/N_c$

• features discrete: $\hat{\theta}_{dck} = N_{dck}/N_c$

• features a valore reale:

$$\hat{\mu}_{dc} = \frac{1}{N_c} \sum_{n: y_n=c} x_{nd} \quad \hat{\sigma}_{dc}^2 = \frac{1}{N_c} \sum_{n: y_n=c} (x_{nd} - \hat{\mu}_{dc})^2$$

6.3. GENERATIVE VS DISCRIMINATIVE CLASSIFIERS

Abbiamo visto che un modello nella forma $p(\vec{x}, y) = p(y)p(\vec{x}/y)$ è detto GENERATIVE CLASSIFIER, perché genera esempi \vec{x} \forall classe y .

Un modello nella forma $p(y/\vec{x})$ è detto DISCRIMINATIVE CLASSIFIER perché può essere usato solamente per discriminare tra classi.

Ⓐ Vantaggi dei DISCRIMINATIVE CLASSIFIERS

- Migliore accuratezza predittiva. La ragione è che stimare $p(y/\vec{x})$ è più semplice che stimare la probabilità congiunta $p(y, \vec{x})$.
- Libertà nel manipolare le features. I dati di input possono essere manipolati a piacere (e.g. basis expansion).
- Probabilità ben calibrate. Alcuni generative classifiers (come NBC) fanno assunzioni particolarmente forti che possono essere sorgenti di bias.

Ⓑ Vantaggi dei GENERATIVE CLASSIFIERS

- Facili da fitare. MLE consiste nel contare e mediare.
- Possono facilmente fare a meno di features mancanti. ⓧ

- Possiamo filtrare classi diverse separatamente
- Possiamo fare a meno del fatto che tutti i dati di training abbiano una label.
- Sono robusti rispetto a features sparse.

* Questo si può fare marginalizzando sulle features non note.

Ricordiamo che la distribuzione categorica è definita come

$$\text{Cat}(y/\bar{\theta}) \triangleq \prod_{c=1}^C \theta_c^{\mathbb{1}(y=c)}$$

cioè $p(y=c/\bar{\theta}) = \theta_c$. I parametri sono tali che $0 \leq \theta_c \leq 1$ e $\sum_{c=1}^C \theta_c = 1$

Ricordiamo che la distribuzione di Bernoulli è definita come

$$\text{Ber}(y/\bar{\theta}) = \begin{cases} 1 - \theta & \text{se } y = 0 \\ \theta & \text{se } y = 1 \end{cases}$$

$$\text{Ber}(y/\bar{\theta}) \triangleq \theta^y (1-\theta)^{1-y}$$

