

Regressione lineare come strumento per il Machine Learning (Apprendimento Automatico)

AI artificiale / Intelligenza

Disciplina che si occupa della progettazione di software / hardware che possono prendere decisioni in maniera autonoma.

Machine Learning (ML) è una branca dell'AI.
(TOM MITCHELL)

ML si preoccupa di progettare algoritmi che si perfezionano in maniera automatica con l'esperienza.

L'idea è quella di fare acquisire alla macchina / al software la capacità di eseguire un

Compito prefissato (Decidere se un'email è spam o meno; predire il costo di un appartamento...)

Un primo modo per raggiungere il nostro obiettivo è quello di

1. Scrivere delle regole ben precise da seguire.
2. Fornire dei dati al software e lasciare che il software impari in maniera autonoma dai dati a disposizione.

affron^{ta}
ML ha due problemi principali

a) CLASSIFICAZIONE

Es. Spam / Non Spam

Melanoma / Non melanoma

Frode / Non Frode

b) REGRESSIONE

(Prendere il valore di un dato)

Es Costo di un appartamento

REGRESSIONE LINEARE

(Altezza padre)

Dati: $x_1, \dots, x_{20}, \dots, x_N$

In generale N è
molto grande

(Altezza studente)

$y_1, \dots, y_{20}, \dots, y_N$

Cerchiamo l'equazione di
una retta

$$y = mx + q$$

che approssimi l'andamento
dei dati.

Dalla Google Sheet infatti:
ci è sembrato di riconoscere

una dipendenza di tipo
lineare delle y dalle x
(Scatter plot)

In realtà $\forall i=1, \dots, 20$

$$(*) \quad y_i = mx_i + q + \underline{e_i}$$

e_i : errore per l' i -esimo dato

obiettivo : cercare m e q

IDEA : li cerco in modo tale
da minimizzare gli errori
commessi.

Def Sia

$$E = \frac{1}{20} \sum_{i=1}^{20} (y_i - (mx_i + q))^2 =$$

$$\frac{1}{20} \sum_{i=1}^{20} e_i^2 \quad \begin{matrix} \text{(Uso l'eq. ne)} \\ (*) \end{matrix}$$

È si dice errore quadratico
medio.

OSSERVAZIONE

Non consideriamo gli ϵ_i con i loro segni per evitare che errori con segno opposto si annullino. Prendiamo il quadrato (invece che ad esempio $|\epsilon_i|$) per pesare di più gli errori più grandi.

REGRESSIONE LINEARE AI MINIMI QUADRATI : Scegliiamo m e q in modo tale da minimizzare \bar{E} (errore quadratico medio)

Procedura: si considera \bar{E} come funzione di m e q
 $\bar{E} = \bar{E}(m, q)$ e se ne calcolano le derivate parziali rispetto a m e q .

Invece
Noi utilizziamo delle definizioni di statistica.

Def $x_1 - x_{20}$

si dice media di $x_1 - x_{20}$

$$\bar{x} = \frac{1}{20} \sum_{i=1}^{20} x_i$$

(La media di $y_1 - y_{20}$: $\bar{y} = \frac{1}{20} \sum_{i=1}^{20} y_i$)

Def Si dice varianza di $x_1 - x_{20}$

$$s_x^2 = \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})(x_i - \bar{x})$$

s_x^2 è una misura di quanto gli x_i si discostano dalla media \bar{x} .

Def Si dice covarianza di

$x_1 - x_{20}, y_1 - y_{20}$

$$s_{xy} = \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})(y_i - \bar{y})$$

$$\underline{\text{Se}} \quad Y_i = mx_i + q \quad i=1, \dots, 20$$

(se gli x_i fossero tutti nulli,
se (x_i, Y_i) sono tutti sulla retta
di equazione $Y = mx + q$)

$$\begin{aligned} 1. \bar{Y} &= \frac{1}{20} \sum_{i=1}^{20} Y_i = \frac{1}{20} \sum_{i=1}^{20} (mx_i + q) = \\ &= \frac{1}{20} \sum_{i=1}^{20} mx_i + \frac{1}{20} \sum_{i=1}^{20} q = \\ &= m \cdot \frac{1}{20} \sum_{i=1}^{20} x_i + q = m \bar{x} + q \end{aligned}$$

$$\begin{aligned} 2. S_{XY} &= \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})(Y_i - \bar{Y}) = \\ &\quad Y_i = mx_i + q \\ &\quad \bar{Y} = m \bar{x} + q \\ &= \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})(mx_i + q - m \bar{x} - q) = \\ &= \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})m(x_i - \bar{x}) = \\ &= m \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})^2 = \\ &= m S_x^2 \end{aligned}$$

Dalla 2. $S_{XY} = m S_x^2$

$$\Rightarrow m = \frac{S_{XY}}{S_x^2}$$

Dalla 1. $\bar{Y} = m \bar{x} + q$

\Rightarrow Posso ricavare q da \bar{Y}, m, \bar{x}

$$q = \bar{Y} - m \bar{x}$$

Se $y_i = mx_i + q$ possiamo ricavare
 m e q da $\bar{x}, \bar{y}, S_x^2, S_{XY}$

TEOREMA $(x_i, y_i), i=1, \dots, 20$

La retta di regressione lineare
(quella che minimizza l'errore
medio quadratico) è

$$Y = mx + q$$

con
$$m = \frac{s_{xy}}{s_x^2}$$
 |,
$$q = \bar{y} - m\bar{x}$$

dove $\bar{x}, \bar{y}, s_{xy}, s_x^2$ sono state definite precedentemente.

Con i nostri dati $\bar{x} \approx 179.85$

$$s_x^2 \approx 42.63$$

$$s_{xy} \approx 37.63$$

$$\bar{y} \approx 145.2$$

$$m = \frac{s_{xy}}{s_x^2} \approx 0.8828$$

$$q = \bar{y} - m\bar{x} \approx 145.2 - 0.8828 (179.85)$$

$$\approx -13.56$$

$$Y = m_1 x_1 + m_2 x_2 + q$$

Si risolve un opportuno sistema

lineare per trovare m_1, m_2 e q

(MINIMI QUADRATI
REGRESSIONE MULTILINEARE)

cinzia.elia@uniba.it