

ANALISI DEI DATI PER LE IMPRESE INDUSTRIALI (5)

Apprendimento supervisionato



Nell'ultimo articolo abbiamo affrontato l'apprendimento **non supervisionato**, ora iniziamo quello **supervisionato**, che per la vastità del

tema, non si esaurirà in una sola trattazione.

Gli algoritmi di apprendimento supervisionato utilizzano **dati che hanno già il risultato per istruire il sistema**. Il modello creato è poi usato per predire il valore delle nuove osservazioni. Esiste un'ampia varietà di algoritmi di apprendimento supervisionato che possono essere raggruppati in due categorie principali, **la regressione e la classificazione**.

A seconda delle caratteristiche dei dati, ci può essere un'attività di regressione, per cui si ricerca un valore numerico (pubblicità e relativa vendita...), o qualitativo per la classificazione (Maschio-Femmina, Sì-No, ...).

REGRESSIONE

Se vogliamo ipotizzare il prezzo di una casa, guardiamo tutte le offerte simili a quella di nostro interesse e raccogliamo le informazioni più importanti come l'anno di costruzione, il numero degli ambienti e dei bagni, ..., e i prezzi. In questo modo, facendo lavorare l'algoritmo e inserendo i dati in questione dell'appartamento che ci interessa, capiremo se il prezzo ipotizzato sarà in linea con il mercato o meno. La regressione lineare **cerca la relazione tra una variabile target continua (il prezzo) con una o più variabili indipendenti** (anno, ambienti, bagni, ...), adattando un'equazione lineare ai dati stessi.

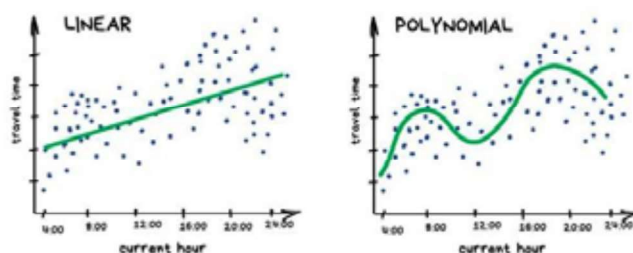
Quando la linea che rappresenta la re-

lazione è retta, è una **regressione lineare**. In casi più complessi, quando la linea è una curva si chiama **polinomiale**. Questi sono i due tipi principali di regressione.

Il grafico a dispersione di fianco (lineare) mostra una correlazione positiva tra una variabile indipendente (asse x, ad esempio spesa per pubblicità) e una variabile dipendente (asse y, volume vendite). All'aumentare di uno, aumenta anche l'altro.

Come si vede dalla figura 2 (sotto), un **modello di regressione lineare cerca di adattare al meglio la retta di regressione ai punti rappresentativi dei dati**. La tecnica più comune che si utilizza è quella dei minimi quadrati. Con questo metodo la migliore retta di regressione si trova minimizzando la

La classificazione oggi è molto utilizzata per filtri antispam, sentiment analysis, per scoprire le frodi, ad esempio con le carte di credito, e per molti altri scopi, anche in ambito scientifico



e di ricerca, prevedendo **le categorie di appartenenza**, due o più.

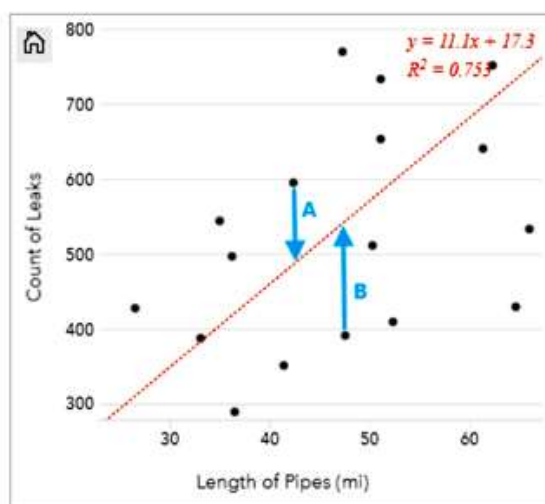
Gli algoritmi più popolari sono: **Support Vector Machine, Naive Bayes, Decision Tree & Random Forest, Logistic Regression, K-NN, ...** Oltre a quelli citati, che sono la base storica, vi sono tante evoluzioni, con specialità mirate e con prestazioni vicine a quelle delle reti neurali per dati di non

enormi dimensioni. Tutti questi algoritmi erano noti da molti anni, ma solo i recenti sviluppi nei processori e nelle memorie hanno permesso il loro uso in maniera estesa.

SUPPORT VECTOR MACHINE

Le Support Vector Machine o SVM (in italiano Macchine a Vettori di Supporto) sono modelli di classificazione il cui obiettivo è quello di trovare **la**

linea di separazione delle classi che massimizza il margine tra le classi stesse, dove con margine si intende la distanza minima dalla linea ai punti delle due classi. Con questo sistema si usa un numero di punti ridotto rispetto a tutti i dati perché lavoriamo solo con quelli di frontiera, che sono quelli più



somma dei quadrati della distanza tra i punti dati e la retta di regressione stessa. Queste distanze sono gli "errori" e si deve usare il quadrato per avere solo valori positivi.

CLASSIFICAZIONE

