

Cloud capacity planning

Per offrire servizi cloud ci sono due tipi di server. Il primo è a domanda: si possono accendere e spegnere questi server quando si vuole, a seconda delle necessità. Ogni volta che se ne usa uno, esso ha un costo orario noto. Il secondo tipo di server è invece su prenotazione. All'atto della prenotazione si paga un costo fisso per un anno mentre per l'uso del server prenotato si paga per tutto l'anno un costo orario molto più basso. I costi di uso e di prenotazione sono noti e ci sono diversi contratti possibili, più o meno convenienti a seconda dell'uso che si prevede di fare. Ricavare la curva dei costi minimi in funzione della capacità complessiva richiesta (ore di calcolo richieste in un anno).

Formulare il problema e classificarlo.

Risolvere l'esempio descritto nel seguito e discutere l'ottimalità e l'unicità della soluzione per ogni dato valore della capacità richiesta.

Esempio

Contratto	Costo annuo	Costo orario
Heavy utilization:	\$ 1560	\$ 0.128
Medium utilization:	\$ 1280	\$ 0.192
Light utilization:	\$ 552	\$ 0.312

Tabella 1: Costi dei server su prenotazione.

Costi orario dei server a domanda: \$ 0.640.

Soluzione

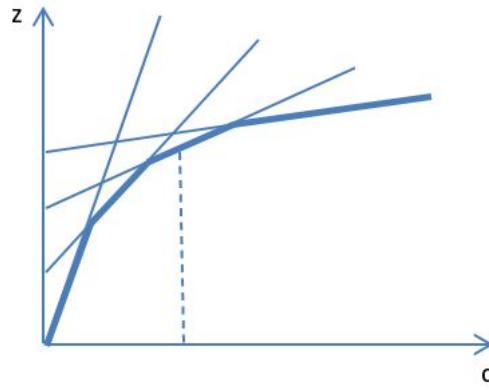
Il problema si può formulare con variabili binarie y che rappresentano la necessità di pagare i costi di prenotazione e variabili continue x che rappresentano la capacità effettivamente usata per ogni tipo di servizio. Il servizio a domanda può essere trattato semplicemente come un caso particolare di servizio a prenotazione in cui il costo di prenotazione è nullo ed il costo variabile è molto alto. In questo modo la funzione obiettivo da minimizzare è data dalla somma su tutti i possibili tipi di servizio dei costi fissi moltiplicati per la corrispondente variabile y e dei costi variabili moltiplicati per la corrispondente variabile x . La relazione tra le due variabili per ogni servizio j è data dal vincolo

$$x_j \leq qy_j$$

dove q rappresenta la capacità complessiva richiesta. La somma delle variabili x deve eguagliare la domanda complessiva:

$$\sum_j x_j = q.$$

Il modello che si ottiene in questo modo è di PLI. Esso quindi consente di calcolare una soluzione con garanzia di ottimalità (non di unicità) per ogni dato valore della capacità richiesta q . Tuttavia non si presta ad eseguire l'analisi parametrica sul parametro q , come è richiesto, per ricavare i costi in funzione della domanda complessiva. Per poter eseguire l'analisi parametrica bisogna riformulare il problema come modello di programmazione lineare con variabili continue. A questo scopo basta osservare che il problema richiede di determinare una funzione lineare a tratti, come in figura.



Ciò equivale ad osservare che è sempre solo uno il servizio da scegliere (una sola variabile x è diversa da zero all'ottimo). Per valori bassi di q risultano convenienti i servizi con minor costo fisso e maggior costo variabile, mentre al crescere di q diventano più convenienti i servizi che hanno alti costi fissi e bassi costi variabili. Pertanto, per ogni dato valore di q il corrispondente valore dei costi, $z(q)$, è dato dal minimo tra i valori di costo corrispondenti ai diversi servizi per quel valore di q . Per trovarlo si può quindi risolvere un problema di massimizzazione lineare nel continuo:

$$\text{maximize } z$$

$$\text{s.t. } z \leq f_j + c_j q \quad \forall j \in S$$

Facendo l'analisi parametrica su q si ottiene facilmente quanto richiesto.

Dall'analisi parametrica si può osservare che il contratto 2 non è mai conveniente.