

# Reducción de costos de infraestructura en la nube

Alexander Baquix

Universidad Galileo de Guatemala  
Departamento de Investigación de Operaciones  
{alexqbq}@galileo.edu

## Abstract

Hoy en día, es un hecho que un gran porcentaje de sistemas informáticos han migrado de despliegues *on prem* a despliegues en la nube. El impacto de este cambio es evidente en la rapidez de construcción de nuevos sistemas, pero, también trae una gran responsabilidad: el control de costos.

A continuación encontraremos un caso real, de cómo usando métodos de programación lineal, podríamos optimizar recursos computacionales de una forma determinista que pueda servir de ayuda para optimizar el uso de recursos y por consiguiente reducir el costo de mantener estos sistemas corriendo.

## 1 Introducción

Olay es una compañía enfocada en comercio conversacional, que diariamente procesa más de 500M de mensajes. Para ser capaces de manejar todo este tráfico, la compañía creó varias docenas de servicios que funcionan en conjunto para proveer los servicios que Olay ofrece al mundo.

Cada servicio debe asegurar alta disponibilidad por lo cuál cada uno posee varias replicas. Esto da como resultado que Olay tenga que manejar cientos de aplicaciones y alojarlos sabiamente en la infraestructura disponible.

Actualmente, esta *elección sabia*, es manual. Un operador decide qué infraestructura es suficiente para alojar todas las aplicaciones, en un inicio era factible pero debido al crecimiento de la compañía esto se ha vuelto más complejo de tratar.

Es por ello que decidí proponer una posible forma de calcular esto de forma determinista, buscando suplir la infraestructura necesaria para alojar todas las aplicaciones y al mismo tiempo reducir el costo que esto conlleva.

## 2 Contexto

El equipo de ingeniería de Olay sigue el patrón bien conocido de *contenedores* y usa *Kubernetes*(K8[1]) como un orquestador de todos estos componentes. K8 es un sistema construido por Google que se encarga de alojar correctamente estos contenedores en máquinas virtuales. Lo que K8 no puede hacer con totalidad es decidir cuántas máquinas virtuales

están disponibles. Y es justo esta tarea la que un operador realiza, indicándole al orquestador las máquinas disponibles y los recursos que esta puede usar.

### 2.1 Contenedores

Para darle más contexto al lector, describiremos superficialmente que es un contenedor en el contexto de software.

Un contenedor es análogamente como un paquete aislado y autosuficiente. Este paquete contiene dependencias y configuraciones aisladas que hacen que una aplicación específica pueda funcionar dentro de ella.

La tecnología detrás de los contenedores, hace que sea posible simular completo aislamiento de las aplicaciones haciéndolas parecer que están corriendo en máquinas completamente diferentes, a pesar de que en realidad estén corriendo en una misma máquina.

La belleza de esto radica en que, cada aplicación puede estar corriendo en sistemas operativos diferentes, con dependencias diferentes, sin afectar al resto de contenedores corriendo en la misma máquina.

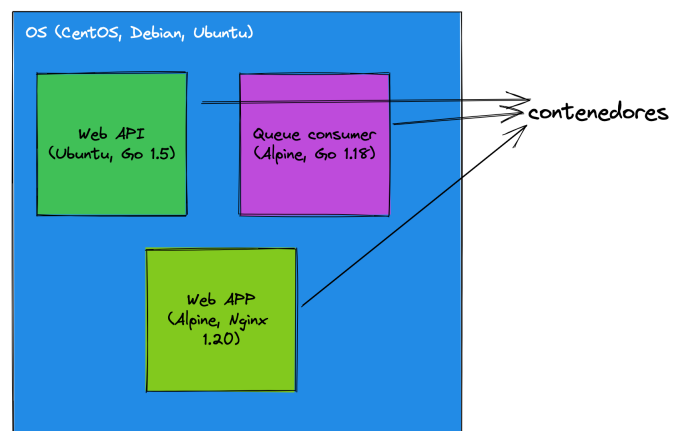


Figure 1: Containers

### 2.2 K8

Kubernetes nació luego de la aparición de los contenedores, con el fin de proveer a los operadores una herramienta declarativa para automatizar la organización de los contenedores



que no existen ese tipo de VMs en el mercado, uno puede elegir entre tipos definidos de VMs, como aquellas que tiene 10 VMs pero solo 2 CPUs.

### 3 Planteamiento del problema

Ya con el contexto en lugar, podrá inferir que lo que tratamos de optimizar es el costo de alquilar VMs pero al mismo tiempo suplir los requisitos de cada aplicación en todo el ecosistema de Olay.

#### 3.1 Recopilación de datos

Olay posee varios cientos de contenedores corriendo en la nube, estos son cerca de 1000, en conjunto se estima que la cantidad total de memoria requerida es de 2122 GB y la cantidad de CPUs necesarios es 520.

#### Costos

Las máquinas disponibles a considerar son las siguientes:

Nombre	Memory (GB)	CPU	Costo mes (\$)
e2-highmem-2	16	2	65.99
e2-standard-4	8	4	97.84
e2-highcpu-8	8	8	144.45

#### Índices

Tipos de máquinas disponibles

$$j = 1, 2, 3 \quad (1)$$

Tipos de recursos disponibles, en este caso sólo consideramos memoria y CPU

$$i = 1, 2 \quad (2)$$

#### Parámetros

Parámetro	Representación
$b$	Cantidad de recurso necesario
$c_j$	costo por máquina j
$a_{ij}$	cantidad de recurso i en maquina j

#### Variables

$x_j$  cantidad de maquinas de tipo j a ser utilizadas

#### Función objetivo

$$\min c_1 * x_1 + c_2 * x_2 + c_3 * x_3 \quad (3)$$

#### Restricciones

$$\begin{aligned} a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in} * x_n &\geq b \\ x_j &\geq 0 \end{aligned} \quad (4)$$

### 3.2 Resultados

Dado los datos descritos anteriormente, podemos proceder a definir el siguiente problema de minimización

$$\min 65.99 * x_1 + 97.84 * x_2 + 144.45 * x_3$$

Sujeto a:

$$16 * x_1 + 8 * x_2 + 8 * x_3 \geq 2122$$

$$2 * x_1 + 4 * x_2 + 5 * x_3 \geq 520$$

$$x_1, x_2, x_3 \geq 0$$

Lo cuál nos da un resultado de:  $z^* = 14258.345$ , con un valor de  $x_1 = 90.17$ ,  $x_2 = 84.92$  y  $x_3 = 0$

## 4 Conclusiones

Con este simple proceso de recolección y análisis, para los operadores de Olay será ahora más determinista elegir qué cantidad de qué tipo de máquinas habilitar para que Kubernetes pueda alojar contenedores de una forma eficiente.

Este es un ejemplo del beneficio de la programación lineal en ámbitos que tal vez parezca inusual usarlo.

## References

- [1] Google. *Kubernetes*. URL: <https://kubernetes.io/>.