

MAESTRÍA EN CIENCIA DE DATOS

TIPOLOGÍA Y CICLO DE VIDA DE LOS DATOS

Práctica 1: Web scraping

Cristhyan Leonardo Naranjo Puertas

Profesor
Xavier Vivancos Garcia

9 de abril de 2020

ÍNDICE

1. Marco de trabajo	2
2. OBSERVACIONES	7

1. MARCO DE TRABAJO

La siguiente es una descripción del trabajo realizado y su justificación

1. **Contexto.** Explicar en qué contexto se ha recolectado la información. Explique por qué el sitio web elegido proporciona dicha información.

Solución. El mercado laboral es un indicador muy importante en el estado económico de un país, por lo que con este proyecto se pretende obtener un registro administrativo a partir de las ofertas de empleo ofrecidas en la página “empleo.com” en donde al realizar una toma continua de muestras se puede generar una base longitudinal que permita ver el ciclo de vida de una oferta laboral, los intervalos de salario con mayor oferta, los puestos laborales con mayor y menor oferta entre otros datos de interés para investigación. Además estamos frente a un momento histórico con las cuarentenas causadas por la pandemia del coronavirus que seguramente afectará la situación laboral de muchas personas, por lo que esta base de datos permitirá hacer un seguimiento a la evolución del mercado laboral durante y después de las cuarentenas, lo que permitirá realizar políticas laborales que conduzcan a un mejor manejo en situaciones futuras.

2. **Definir un título para el dataset.** Elegir un título que sea descriptivo.

Recolección de información actualizada del mercado laboral usando web scripting sobre la página “el empleo.com”

3. Descripción del dataset. Desarrollar una descripción breve del conjunto de datos que se ha extraído (es necesario que esta descripción tenga sentido con el título elegido).

Nombre	Tipo	Descripción	Ejemplo
Empleo	Character	Nombre del empleo ofrecido	Gerente general
Sueldo	Character	Intervalo de sueldo ofrecido	18 a 21 millones
Empresa	Character	Organización que ofrece el empleo	Empresa confidencial
Ciudad	Character	Ciudad en donde se ofrece el empleo	Cali
Vacante	Numeric	Cantidad de puestos ofrecidos para ese empleo	1
Fecha	Character	Fecha de publicación de la oferta	Publicado 16 Mar 2020
Experiencia	Character	Cantidad de años mínimos para poder aplicar a la oferta	10 años de experiencia
Contrato	Character	Tipo de contrato ofrecido	Contrato indefinido
Área	Character	segmento económico donde se ubica la empresa ofertante	Salud
Sub área	Character	sub area del segmento económico donde se ubica la empresa	Medicina
Nivel	Character	Nivel del cargo	Gerente
Sector	Character	Sección o departamento dentro de la empresa	Químicos
Nivel educativo	Character	Nivel educativo mínimo para poder aplicar	Especialización
Descripción	Character	requisitos necesarios para poder aplicar al empleo	Importante multinacional alemana requiere...
Serial	Numeric	Serial de la oferta dentro de la pagina	1884488870
Link	Character	Dirección web de la oferta	https : //www.elemp...

4. **Representación gráfica. Presentar una imagen o esquema que identifique el dataset visualmente**

La siguiente imagen muestra una la estructura de una publicación de empleo estándar en donde se pueden ver las diferentes variables que son capturadas

Informador horario adicional termino fijo - davivienda - bucaramanga

\$ Salario a convenir

📍 Bucaramanga

📅 Publicado 4 Abr 2020

🔧 Administrativa y Financiera

👥 1 Vacante

🎓 Bachillerato Académico

✓ Aplicar a oferta



Davivienda
Sector de la vacante:
Financiero

Cuenta creada unicamente para la presencia de marca dentro del portal

Descripción general

Somos un Banco que se interesa en hacer sentir a nuestros clientes satisfacción por el servicio que les brindamos, por eso si tienes interés genuino por la gente y te gusta hacer las cosas sencillas para los demás, este trabajo es para ti.
Davivienda está en busca de un Informador Horario Adicional y que cumpla con las siguientes responsabilidades:

Figura 1.1: Ejemplo de una publicación

El siguiente diagrama ilustra el proceso que se debe realizar para la captura de la información, dentro del programa se deberá seleccionar si se obtendrá la información de todo el país o de una ciudad específica, luego se debe seleccionar si se van a capturar todos los o solo algún intervalo de sueldos específico, después se inicia el ciclo según la cantidad de ofertas disponibles almacenando las 16 variables que se muestran.

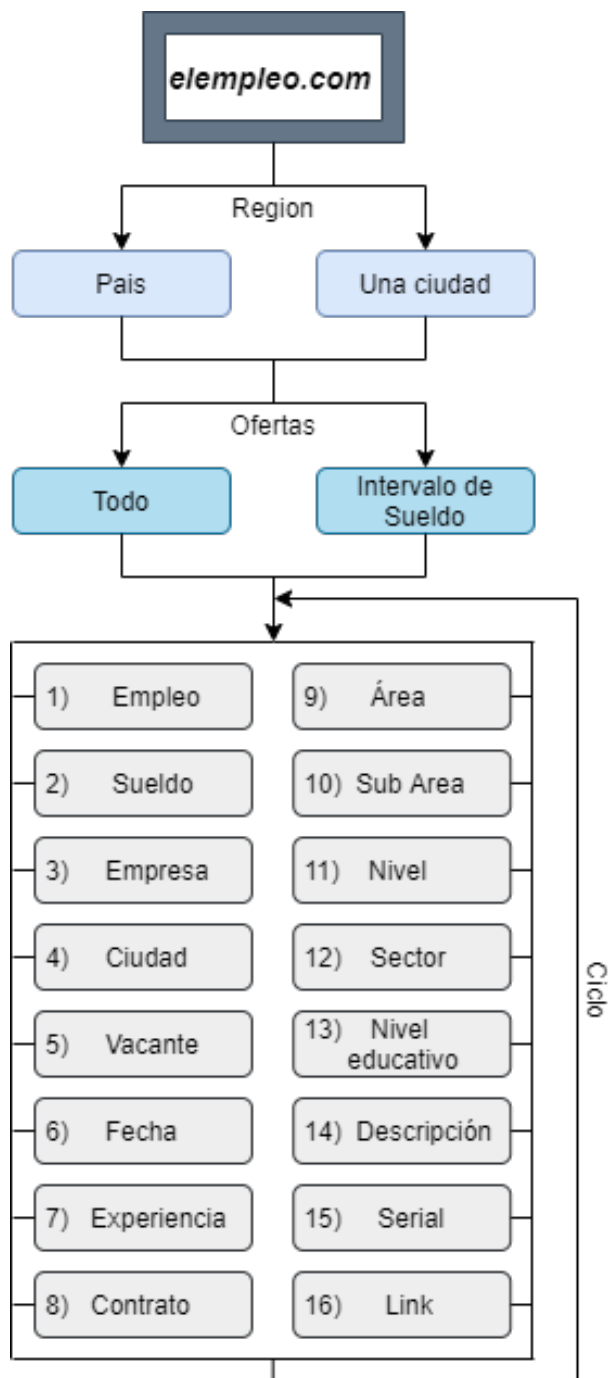


Figura 1.2: Diagrama

5. Contenido. Explicar los campos que incluye el dataset, el periodo de tiempo de los datos y cómo se ha recogido.

El data set llamado *Empleo_2020-04-05_.csv* fue realizado en una única toma el 5 de abril del 2020, seleccionando todas las ofertas de empleo ofrecidas para la ciudad de Cartagena Colombia. El data set contiene los siguientes datos: *Empleo, Sueldo, Empresa, Ciudad, Vacante, Fecha, Experiencia, Contrato, Área, Sub área, Nivel, Sector, Nivel, educativo, Descripción, Serial y Link*, con las características que se describieron en la tabla del tercer punto. En general los datos contienen la información detallada de la oferta de empleo ofrecida, destacando el serial que permite un seguimiento de la oferta desde su creación hasta su desaparición.

6. Agradecimientos. Presentar al propietario del conjunto de datos. Es necesario incluir citas de investigación o análisis anteriores (si los hay).

Agradezco a la empresa Leadearsearch S.A.S propietaria del producto *empleo.com* (COPY-RIGHT © 2019)

Esta investigación no tiene fines comerciales y se realiza acorde a los términos y condiciones de la pagina, cito (la siguiente cita no busca un respaldo por el creador del contenido, solo se usa para ilustrar el punto)

USO DEL MATERIAL¹

La Compañía lo autoriza a usted a consultar, revisar y usar el material que se encuentra en el Web Site, únicamente para su uso personal y no comercial y según lo establecido en estos Términos y Condiciones. El contenido de este Web Site, incluyendo pero sin limitarse a los textos, gráficas, imágenes, logotipos, iconos, software y cualquier otro material (el 'Material') están protegidos bajo las leyes colombianas de derechos de autor, leyes de propiedad industrial y otras leyes aplicables. Todo el Material es de propiedad de la Compañía o de sus proveedores o clientes. El uso no autorizado del material puede constituir una violación de las leyes colombianas o extranjeras sobre derechos de autor, leyes de propiedad industrial u otras leyes. Usted no podrá vender o modificar el Material en manera alguna ni ejecutar o anunciar públicamente el Material ni distribuirlo para propósitos comerciales. Usted no podrá copiar o adaptar el código HTML que la Compañía crea para generar sus páginas, ya que el mismo está protegido por los derechos de autor de la Compañía.

¹<https://www.empleo.com/co/terminos-condiciones>

7. **Inspiración.** Explique por qué es interesante este conjunto de datos y qué preguntas se pretenden responder.

Actualmente trabajo en el instituto de estadísticas de mi país (DANE) en un proyecto para crear una base de datos longitudinal a partir de la información recopilada por el sistema de seguridad social, para esto se han buscado mas registros administrativos que enriquezcan la base producidos por otras entidades publicas. Con esta actividad pude reconocer el potencial que tiene el web scraping de la pagina «El empleo» para recopilar información directamente del mercado laboral, con una vigencia única y un coste muy bajo, aun no se si se puede incluir esta información directamente en la base, pero si estoy seguro que va a ser útil para realizar estadística que permitan comparar con los datos que produzca la base.

8. **Licencia.** Seleccione una de estas licencias para su dataset y explique el motivo de su selección:

Al elegir la licencia adecuada estoy teniendo en cuenta que no soy el productor ni el propietario de la información que esta en el dataset, y que no realice ninguna alteración a los datos, solo me limite a recolectarlos y organizarlos. Por esto no me es posible autorizar un uso que valla en contra de los limites que impuso el productor y dueño de los datos. Por estos elijo la licencia **Released Under CC BY-NC-SA 4.0 License** ya que permite copiar,adaptar,transformar y redistribuir en cualquier medio o formato la información, pero se debe realizar un reconocimiento adecuado a los autores y prohíbe el uso comercial, por lo que considero que va acorde a los términos y condiciones del creador del contenido.

Contribuciones	Firma
Investigación previa	Cristhyan Naranjo
Redacción de las respuestas	Cristhyan Naranjo
Desarrollo código	Cristhyan Naranjo

2. OBSERVACIONES

Realice el programa en R porque ya había tenido experiencia manejando el paquete RSelenium , por lo que me pareció mas cómodo usar este lenguaje. Se realizaron varios ensayos para lograr un programa eficiente que capturara la mayor cantidad de información útil, pero hubo un problema que no se ha podido solucionar, resulta que de forma al parecer de manera aleatoria (al inicio pensé que se causaba por la inclusión de nuevas ofertas, luego vi que no) después de entrar a una oferta, al volver a la lista de empleos , estos aparecen ordenados de manera diferente,

luego como el programa recorre de forma ordenada las ofertas, terminan habiendo ofertas que no se consultan y otras que se consultan varias veces, sospecho que esto lo hace la página precisamente para evitar que robots consigan la información. Por lo anterior se obtuvo un dataset con registros repetidos, y algunos registros nulos, así que por el momento se está generando una base interesante para practicar técnicas de limpieza y obtener una muestra significativa de los empleos ofrecidos, Además queda abierto a que junto con la participación de la comunidad se puedan arreglar los problemas detectados. En las siguientes imágenes se puede ver como después de realizar la primer consulta, al volver a la lista de ofertas, algunos empleos cambiaron de posición.

The screenshot shows a job listing interface. On the left, there are filters for salary and location. The salary filter ranges from 'Menos de \$1' to 'A convenir', with counts for each range. The location filter is set to 'Todas las ciudades'. Below these, there is a 'Fecha de publicación' (Publication Date) filter with options: 'Hoy y ayer', 'Hace 1 semana', 'Hace 2 semanas', and 'Hace 1 mes'.

The main area displays a list of job offers. Each offer includes a company logo, the job title, the company name, the salary, the location, and the publication date. A 'Vista rápida' (Quick View) button is present for each offer.

Salario	Cantidad
Menos de \$1	255
\$1 a \$1,5	3954
\$1,5 a \$2	2195
\$2 a \$2,5	1579
\$2,5 a \$3	1219
\$3 a \$3,5	707
\$3,5 a \$4	508
\$4 a \$4,5	516
\$4,5 a \$5,5	435
\$5,5 a \$6	331
\$6 a \$8	302
\$8 a \$10	96
\$10 a \$12,5	93
\$12,5 a \$15	28
\$15 a \$18	20
\$18 a \$21	4
Más de \$21	7
A convenir	4828

Logo	Título	Empresa	Salario	Ubicación	Fecha	Acción
	Auxiliar administrativoa	TRANSPORTADORA DE VALORES DEL SUR LTDA	Salario a convenir	Quibdó	Publicado 9 Abr 2020	Vista rápida
	Guarda de seguridad	TRANSPORTADORA DE VALORES DEL SUR LTDA	Salario a convenir	Pereira	Publicado 9 Abr 2020	Vista rápida
	Cajero tipo ii	TRANSPORTADORA DE VALORES DEL SUR LTDA	Salario a convenir	Quibdó	Publicado 9 Abr 2020	Vista rápida
	Ingeniero de telecomunicaciones	GRUPO ACCION PLUS	\$4,5 a \$5,5 millones	Bogotá	Publicado 9 Abr 2020	Vista rápida
	Analista de comercio exterior inhouse- bogotá	M MUNAR CONSULTORES S.A.S	\$1,5 a \$2 millones	Bogotá	Publicado 9 Abr 2020	Vista rápida
	Jefe de proyectos	Empresa confidencial	\$4,5 a \$5,5 millones	Bogotá	Publicado 9 Abr 2020	Vista rápida
	Aprendiz universitario - davivienda - cali					

Figura 2.1: Lista de ofertas antes de la consulta

