

AJAX / Escalabilidad de aplicaciones web

Jesús Arias Fisteus

Computación Web (2024/25)

uc3m | Universidad **Carlos III** de Madrid
Departamento de Ingeniería Telemática

Parte I

Peticiones HTTP asíncronas en JavaScript

El término **AJAX**¹ se refiere a la capacidad de un *script* que se ejecuta en el lado del cliente en el contexto de una página Web para **cargar datos** desde el servidor o **enviar datos** al servidor **sin necesidad de recargar la página completa**.

¹Originalmente, acrónimo de *Asynchronous JavaScript and XML*.

El intervalo de tiempo desde que el *script* envía una petición HTTP hasta que recibe su respuesta es impredecible (p.e., podría haber congestión en la red o el servidor, etc.).

Una **petición síncrona** bloquearía el hilo de ejecución hasta la llegada de la respuesta, impidiendo mientras tanto la actualización de la página o la ejecución de sus *scripts*.

Una **petición asíncrona** no bloquea el hilo de ejecución, y la página continúa reaccionando normalmente mientras no llega la respuesta.

Es **necesario**, por tanto, que los *scripts* **envíen sus peticiones asíncronamente**.

Con peticiones asíncronas, se procesará la respuesta cuando esta llegue, mediante el modelo de **programación asíncrona** de JavaScript (basado en el objeto **Promise**) o mediante **funciones de callback**.

Las APIs `fetch` y `XmlHttpRequest`

Los programas JavaScript pueden enviar peticiones asíncronas mediante dos APIs alternativas: `fetch` y `XmlHttpRequest`.

El API `fetch` proporciona una solución más moderna y sencilla, y se basa en el objeto `Promise`.

El API `XmlHttpRequest` es más antigua, compleja y se basa en funciones de *callback*.

JSON (JavaScript Object Notation) es actualmente el formato más habitual para representar los datos intercambiados entre cliente y servidor.

```
1 {  
2   "name": "mary",  
3   "email": "mary@example.com",  
4   "roles": ["player", "moderator"],  
5   "level": 12,  
6   "location": {  
7     "city": "Leganés",  
8     "region": "Comunidad de Madrid",  
9     "country": "Spain"  
10  }  
11 }
```

Ejemplo con fetch (GET)

```
1 fetch("/last_msgs")
2   .then((response) => response.json())
3   .then((json) => console.log(json))
4   .catch((err) => console.error("error:", err));
```

Ejemplo con fetch (POST)

```
1 fetch("/post_msg", {
2   method: "POST",
3   headers: {
4     "Content-Type": "application/json",
5   },
6   body: JSON.stringify(message)
7 }).then((response) => response.json())
8   .then((json) => console.log(json))
9   .catch((err) => console.error("error:", err));
```

Ejemplo con JQuery (GET)

```
1 $.getJSON("/last_msgs").done(function(messages) {
2     console.log(messages);
3 }).fail(function(jqxhr, textStatus, error) {
4     var err = textStatus + ",␣" + error;
5     console.log("Request␣Failed:␣" + err);
6 });
```

Ejemplo con JQuery (POST)

```
1 const message = new Message(...);
2 $.ajax({
3   type: "post",
4   url: "/post_msg",
5   data: JSON.stringify(message),
6   contentType: "application/json",
7   success: function (confirmation) {
8     console.log(confirmation);
9   },
10  error(function(jqxhr, textStatus, error) {
11    var err = textStatus + ",␣" + error;
12    console.log("Request␣Failed:␣" + err);
13  })
14 });
```

Por motivos de seguridad, **los navegadores aplican restricciones a las peticiones** enviadas desde el origen de la página a recursos en otros orígenes (otros dominios).

En general, un *script* en el contexto de una página de un origen dado puede enviar peticiones asíncronas, con algunas restricciones, a otros orígenes, pero **no puede leer sus respuestas**.

Sin esta restricción, sería fácil que un script ejecutándose en una pestaña robase datos privados del usuario de aplicaciones abiertas en otras pestañas.

Una aplicación podría necesitar, para funcionar correctamente, deshabilitar esta protección, si se considera que esto es seguro, mediante **CORS (Cross-Origin Resource Sharing)**.

Añadir la cabecera siguiente en la respuesta HTTP a una petición desde otro origen, sea cual sea, permite al script leer esta respuesta:

```
Access-Control-Allow-Origin: *
```

Añadir la cabecera siguiente en la respuesta HTTP permite al script leer esta respuesta siempre y cuando este se esté ejecutando en el contexto de una página con origen en `https://www.example.com`:

```
Access-Control-Allow-Origin: https://www.example.com
```

- ▶ Fetch API en MDN.
- ▶ Cross-Origin Resource Sharing (CORS) en MDN.

Parte II

Escalabilidad en aplicaciones web

El término **escalabilidad** se refiere a la capacidad de la infraestructura de servidor de una aplicación web de dar servicio a una carga creciente (tasa de peticiones, número de usuarios, etc.), o capacidad de ser ampliable para ello.

El uso de **memoria caché** reduce la carga del servidor Web y del gestor de bases de datos.

No es necesario volver a construir páginas HTML que han sido construidas recientemente sobre datos que no han cambiado.

No es necesario volver a solicitar datos que se han consultado recientemente a la base de datos si no han cambiado.

- ▶ Uso de los mecanismos de control de caché de HTTP 1.1 y versiones posteriores.
- ▶ Uso de caché de páginas delante del servidor Web (*proxy inverso*), como NGINX².
- ▶ Uso de sistemas distribuidos de almacenamiento clave-valor en RAM, como Memcached o Redis:
 - ▶ Almacenamiento de fragmentos de HTML ya construidos.
 - ▶ Almacenamiento de objetos obtenidos de la base de datos.

²<https://www.nginx.com/resources/glossary/caching/>

Cuando el uso de memoria caché no es suficiente, se pueden **incrementar los recursos de cómputo**.

La **escalabilidad vertical** consiste en mejorar los equipos que ejercen de servidor Web o gestor de base de datos (más y/o mejores CPUs, más RAM, más disco, etc.).

La **escalabilidad horizontal** consiste en replicar equipos (añadir más servidores Web o más servidores de bases de datos), sin necesariamente mejorar la capacidad de cada servidor.

Para escalar horizontalmente el servidor Web, **se despliegan servidores Web en varias máquinas** y se reparten las peticiones entre ellos:

- ▶ Balanceo de carga por DNS: un mismo nombre de dominio con varias IPs.
- ▶ Balanceo de carga con servidores de *front-end*.

Es necesario **gestionar sesiones de forma distribuida**.

En general, resulta **más complejo replicar el gestor de bases de datos** que los servidores Web.

Cuando el número de lecturas en la base de datos es considerablemente superior al de escrituras es factible replicar el gestor de bases de datos en **un servidor maestro y varios servidores esclavos**.

Los esclavos se utilizarían principalmente para lecturas, y el maestro realizaría las escrituras y propagaría los cambios a los esclavos.

En ocasiones es posible **dividir tareas y/o datos en particiones** para escalar horizontalmente.

Particionar tareas consiste en que cada servidor Web se encargue de un grupo de tareas concreto.

Por ejemplo, una tienda Web podría ejecutar las tareas de gestión de ventas en un servidor (carro de la compra, tramitación de pedidos, etc.) y las tareas de gestión de clientes en otro (creación de cuentas, configuración de cuentas, autenticación, etc.).

Particionar datos consiste en que cada gestor de bases de datos guarde un subconjunto de los datos, procurando que cada tarea se pueda llevar a cabo, en la medida de lo posible, sobre un único gestor de bases de datos.

Una posible estrategia consiste en **particionar datos por temática**.

Por ejemplo: clientes en un gestor de bases de datos, pedidos en otro, registros contables en otro, etc.

También se pueden particionar datos atendiendo a **otros criterios**.

Por ejemplo, se podrían dividir los datos por países o regiones, de forma que cada gestor de bases de datos guardase todos los datos de clientes, pedidos, registros contables, etc. de un país o región concretos.

Es posible también utilizar **bases de datos no relacionales**³ para escalar el gestor de bases de datos más eficientemente.

Las bases de datos no relacionales relajan las restricciones ACID para distribuir la base de datos de forma más eficiente.

³Conocidas habitualmente como bases de datos NoSQL.

Las bases de datos **clave-valor** almacenan cada valor asociado a una clave.

Son más eficientes en el acceso a valores por clave.

Ejemplos: Memcached⁴, Redis⁵, Amazon DynamoDB⁶, Oracle Berkeley DB⁷.

⁴<https://memcached.org/>

⁵<https://redis.io/>

⁶<https://aws.amazon.com/es/dynamodb/>

⁷[https:](https://www.oracle.com/database/technologies/related/berkeleydb.html)

[//www.oracle.com/database/technologies/related/berkeleydb.html](https://www.oracle.com/database/technologies/related/berkeleydb.html)

Las bases de datos **orientadas a columnas** almacenan de forma contigua las columnas de datos en vez de las filas.

Son más eficientes en cálculos sobre una o varias columnas que las basadas en registros (filas).

Ejemplos: Apache Cassandra⁸, Apache HBase⁹, Google Bigtable¹⁰.

⁸<https://cassandra.apache.org/>

⁹<https://hbase.apache.org/>

¹⁰<https://cloud.google.com/bigtable>

Las bases de datos **basadas en documentos** almacenan los datos de forma semiestructurada y jerárquica, sin un esquema prefijado.

Son más eficientes para almacenar datos que no siguen un esquema (tablas/columnas) prefijado.

Ejemplos: MongoDB¹¹, CouchDB¹², Google Firebase Realtime Database¹³.

¹¹<https://www.mongodb.org/>

¹²<https://couchdb.apache.org/>

¹³<https://firebase.google.com/docs/database>

Las bases de datos **basadas en grafos** almacenan los datos con estructura de grafo (vértices y aristas).

Son más eficientes para almacenar y procesar datos altamente conectados en forma de grafos como, por ejemplo, datos de la Web semántica (RDF).

Ejemplos: Neo4J ¹⁴, Openlink Virtuoso¹⁵.

¹⁴<https://www.neo4j.org/>

¹⁵<https://virtuoso.openlinksw.com/>

- ▶ Martin L. Abbott; Michael T. Fisher. “Scalability Rules: Principles for Scaling Web Sites, 2nd Edition” . Addison-Wesley Professional (2016).
 - ▶ Acceso en línea en O'Reilly
- ▶ Chander Dhall. “Scalability Patterns: Best Practices for Designing High Volume Websites” . Apress (2018).
 - ▶ Capítulos 1, 2 y 3.
 - ▶ Acceso en línea en O'Reilly