



**OPTATIVA**  
**Recuperación Avanzada de**  
**Información**

---

Dr. J. Federico Medrano

*@jfedemedrano*

Unidad N° 2 – Parte 2

# Temas a desarrollar

---

- ~~Diseño de un spider/crawler. Estrategias.~~
- ~~Web scraping.~~
- ~~Recolección y representación estructural de sitios web/dominios.~~
- ~~Tratamiento de grafos.~~
- Minería de textos. Minería de datos.
- Recolección de repositorios. Protocolos de recuperación.
- Diseño de un motor de búsqueda. Diseño de un indexador de documentos.
- Optimización en motores de búsqueda web (SEO/SEM)

Minería de textos. Minería de  
datos

---

# Minería de datos

---

- El **Data Mining** es un conjunto de técnicas y tecnologías que permiten explorar grandes bases de datos, de manera *automática o semiautomática*, con el objetivo de encontrar patrones repetitivos que expliquen el comportamiento de estos datos.
- La minería de datos surgió con la **intención o el objetivo de ayudar a comprender una enorme cantidad de datos**, y que estos, pudieran ser utilizados para extraer conclusiones para contribuir en la mejora y crecimiento de las empresas, sobre todo, por lo que hace a las ventas o fidelización de clientes.
- Su intención es la de aportar información valiosa a las empresas para así, ayudarlas en la toma de decisiones futuras.

# Minería de datos

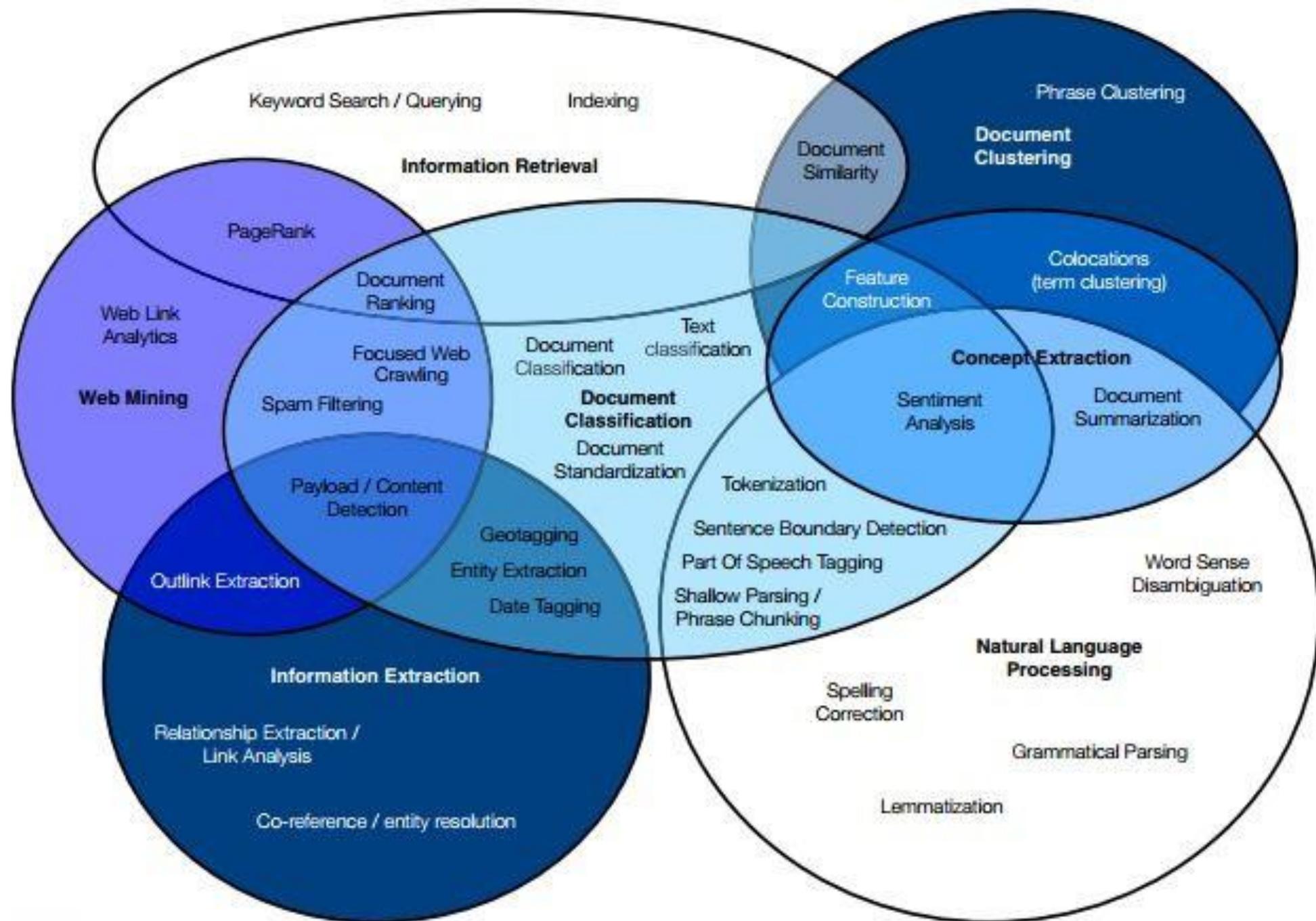
---

- Es un mecanismo de explotación, consistente en la **búsqueda de información valiosa en grandes volúmenes de datos**. Está muy relacionada con las bodegas de datos que proporcionan la información histórica con la cual los algoritmos de minería tienen la información necesaria para la toma de decisiones.
- Según *Fallad* y sus coautores (1996): "La minería de datos es un proceso no trivial de **identificación** válida, novedosa, potencialmente útil y entendible **de patrones** comprensibles que se encuentran **ocultos en los datos**".
- Según *Molina* y sus colaboradores (2001): "Es la integración de un conjunto de áreas que tienen como propósito la **identificación** de un **conocimiento** obtenido a partir de las bases de datos que aporten un sesgo hacia la toma de decisión".

# Minería de textos

---

- Debido a que el 80% de la información de una compañía está almacenada en forma de documentos, las técnicas como la categorización de texto, el procesamiento de lenguaje natural, la extracción y recuperación de la información o el aprendizaje automático, entre otras, apoyan al Text Mining (minería de texto).
- En ocasiones se confunde el Text Mining con la recuperación de la información (Information Retrieval o IR) (Hearst, 1999). Esta última consiste en la recuperación automática de documentos relevantes mediante indizaciones de textos, clasificación, categorización, etcétera. Generalmente, se utilizan palabras clave para encontrar una página relevante.
- En cambio, el Text Mining se refiere al examen de una colección de documentos y el **descubrimiento de información no contenida en ningún documento** individual de la colección; en otras palabras, trata de obtener información sin haber partido de algo. (Nasukawa y otros, 2001).



# Recolección de repositorios.

---

Protocolos de recuperación

# Interoperabilidad de los sistemas de información

---

- *Interoperabilidad*. Esfuerzo requerido para acoplar un sistema con otro (Pressman, 2010).
- *Interoperabilidad* es la posibilidad de que distintos tipos de ordenadores, redes, sistemas operativos, y aplicaciones trabajen juntos de forma eficaz, sin comunicación previa, de tal forma que puedan intercambiar información de manera útil y con sentido (Gomez, 2007)

# Niveles de interoperabilidad

---

- **Sintáctica (De Giusti)**
- Hace referencia a todo lo necesario para que dos sistemas sean capaces de establecer una comunicación e intercambiar información.
- Esto incluye:
  - protocolos de comunicación y transferencia
  - codificación de caracteres
  - formatos de datos

# Niveles de interoperabilidad

---

## Sintáctica

- Elementos que corresponden a la interoperabilidad sintáctica pueden ser, por ejemplo:
  - protocolo TCP/IP
  - protocolo HTTP
  - **protocolo OAI-PMH**
  - formato XML y esquemas XML (XSD)
  - Directrices de interoperabilidad

# Niveles de interoperabilidad

---

## **Semántica (De Giusti)**

- Hace referencia a todo lo necesario para que el sistema receptor haga una correcta interpretación de la información recibida, de forma automática.
- Se busca que el sistema receptor "**entienda**" los datos tal como los "**entiende**" el emisor.
- ***Para contar con interoperabilidad semántica, primero debe asegurarse la interoperabilidad sintáctica***

# Niveles de interoperabilidad

---

## **Semántica**

Entran en juego:

- Formatos de metadatos
- Vocabularios controlados:
  - Tesoros
  - Sistemas de clasificación
- Ontologías
- Directrices de interoperabilidad

¿Cómo lograr que dos sistemas que no se conocen puedan comunicarse/interactuar?

La respuesta es ***hablar el mismo idioma***, en este sentido, se debe implementar un protocolo de comunicación en común (intercambio de información).

# Protocolos de intercambio de información

---

## Open Archives Initiative Protocol for Metadata Harvesting ([OAI-PMH](#))

- Establece un conjunto de reglas a partir de las cuales puede realizarse el intercambio de recursos de forma exitosa.
- Se centra en la **transferencia** de metadatos de un extremo a otro, sin establecer restricciones en cuanto a los datos que se transfieren.

# OAI-PMH

---

Define dos perfiles de trabajo

- **Data Provider:** es aquél repositorio que ofrece sus recursos bajo el protocolo OAI-PMH, para que otros los recolecten mediante cosechas.
- **Service Provider:** es aquél que recolecta recursos desde distintos Data Providers y brinda un servicio a una comunidad de usuarios en base a los recursos recolectados y el valor agregado aportado sobre los mismos (deduplicación, normalización, ordenamiento, búsquedas, etc).

Diseño de un motor de  
búsqueda. Diseño de un  
indexador de documentos.

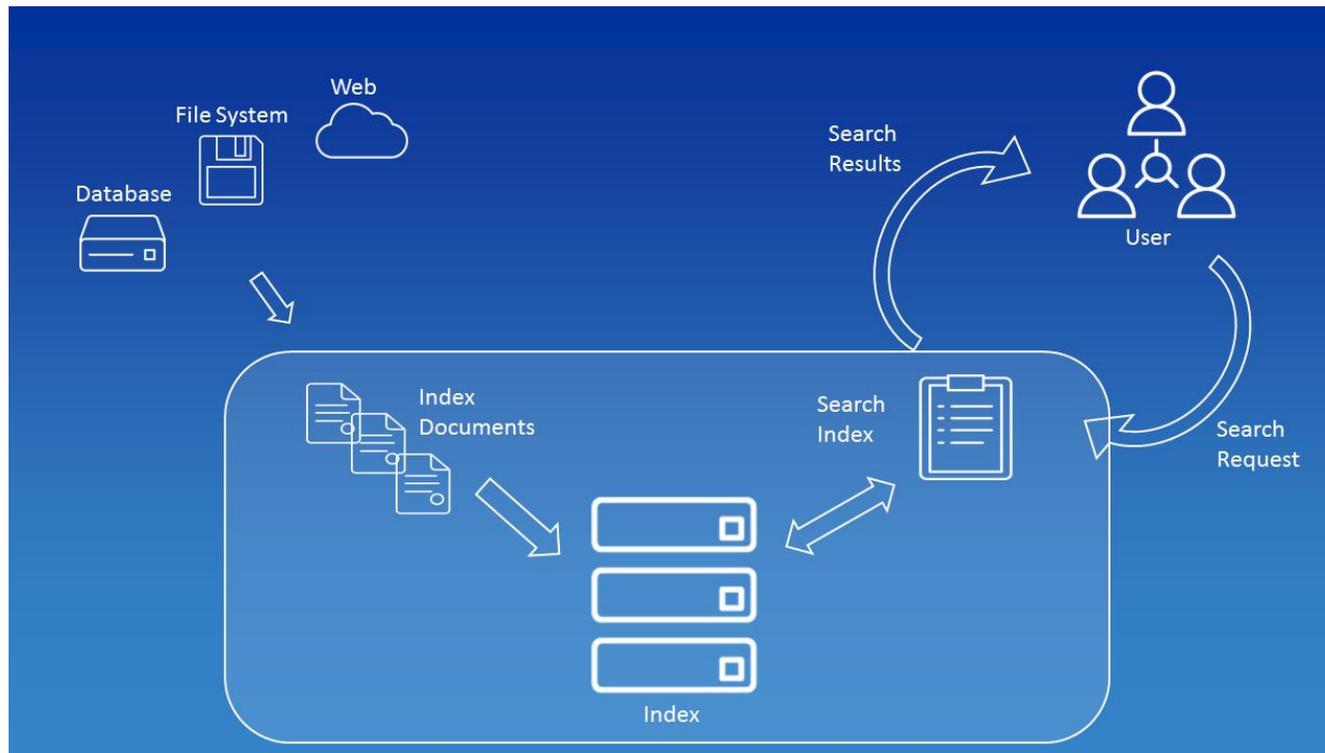
---

# Lucene

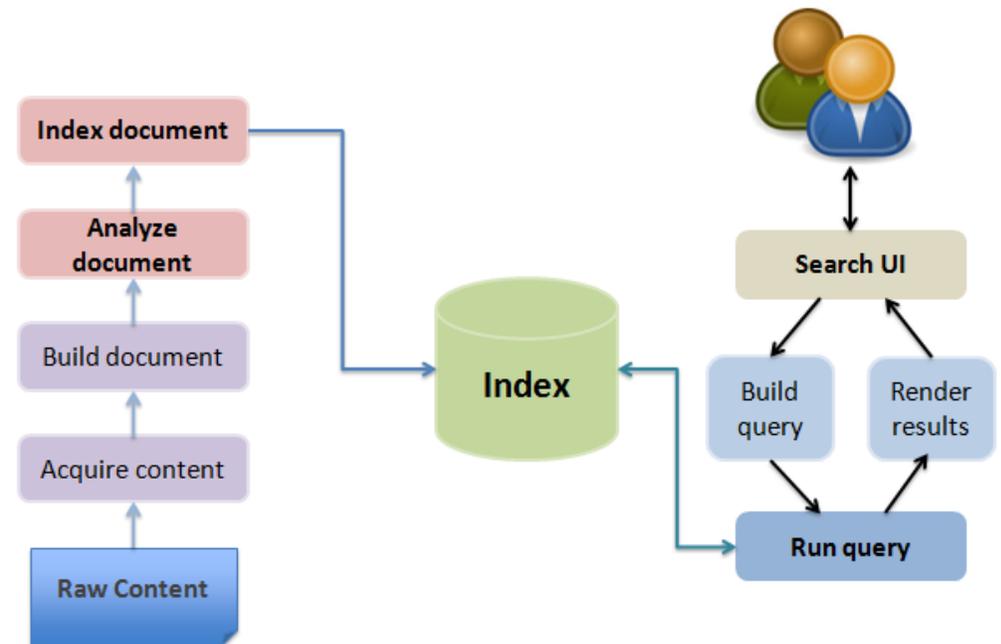


- <https://lucenenet.apache.org/> (C#)
- <https://lucene.apache.org/> (Java)
- <https://lucene.apache.org/pylucene/> (Python)

**Apache Lucene** es una API de código abierto para recuperación de información. Es útil para cualquier aplicación que requiera indexado y búsqueda a texto completo.



## Lucene Flow



# EJEMPLOS

---

- <https://www.youtube.com/watch?v=Lv2wJRvSddw>
- <https://www.youtube.com/watch?v=TnYDeJkULKc>
- <https://howtodoinjava.com/lucene/lucene-index-search-examples/>
- <https://www.baeldung.com/lucene>
- <https://pythonhosted.org/lupyne/examples.html>

# Optimización en motores de búsqueda web (SEO/SEM)

---

# SEM (*Search Engine Marketing*)

---

- El **concepto de SEM** (*Search Engine Marketing*) se refiere a la promoción de un sitio web en los buscadores mediante el uso de anuncios de pago a través de plataformas como [Google Ads](#) o [Bing Ads](#). Y en ocasiones, se amplía este concepto para referirnos también a otro tipo de publicidad mediante estas y otras plataformas de display y medios sociales, donde se suele pagar generalmente en base a los clics que nos generan los anuncios.
- Mediante esta estrategia el objetivo es dar visibilidad inmediata a nuestro sitio Web, ya que desde que configuramos las campañas y pujamos por salir, nuestros anuncios tienen la posibilidad de aparecer.

# SEO (*Search Engine Optimization*)

---

- El concepto de **SEO** (*Search Engine Optimization*) se refiere al trabajo de optimización y de aumento de la popularidad de un sitio web, con el objetivo de que dicho sitio sea **rastreable** por los motores de búsqueda, **indexado correctamente** y suficientemente relevante para que algunas o muchas de las páginas sean mostradas en las **primeras posiciones** de los buscadores para determinadas consultas de búsqueda de los usuarios.
- Por lo tanto, se trata de conseguir aparecer en los primeros resultados (lo ideal es en la primera página, y a ser posible, en las cinco primeras posiciones) de un buscador para un conjunto de búsquedas que nos interesan, pero sin tener que pagar un coste directo publicitario por cada visita, gracias a que somos muy relevantes y/o populares.

# Web Spam

---

- Técnicas de SEO *Black Hat* o Sombrero Negro.
- Una práctica para conseguir una posición elevada, pero injustificada, en los resultados de los motores de búsqueda, utilizando técnicas para engañar a los algoritmos de clasificación.
- **Spamdexing** es una práctica poco ética de SEO que puede terminar con penalizaciones de los motores de búsqueda. Esto significa que los sitios web que practican spamdexing incluyen una gran cantidad de información en las páginas con el fin de indexarse y posicionarse más arriba en el resultado de la página del motor de búsqueda, pero no proporcionan información relevante o cualitativa para los usuarios.

# Spamdexing

---

- Agregado automático de BackLinks (Enlaces entrantes que apuntan a un determinado sitio web o página que proviene de otro sitio web.)
- Keyword stuffing
- Link farms
- Spam blog
- Cloaking
- Doorway page



*Puedes comprar **zapatos rojos** en nuestra tienda online. Si necesitas comprar **zapatos rojos** para una ocasión especial, aprovéchate ahora de nuestra oferta en **zapatos rojos**. No encontrarás **zapatos rojos** de tanta calidad a este precio.*

El Keyword Stuffing es una técnica de [Black Hat SEO](#) que consiste en el uso excesivo de palabras clave dentro de un texto con el objetivo mal enfocado de darle más relevancia a esta palabra. Google penaliza con mucha frecuencia este tipo de sobre-optimización.

## Types of Firepits

December 21st, 2012

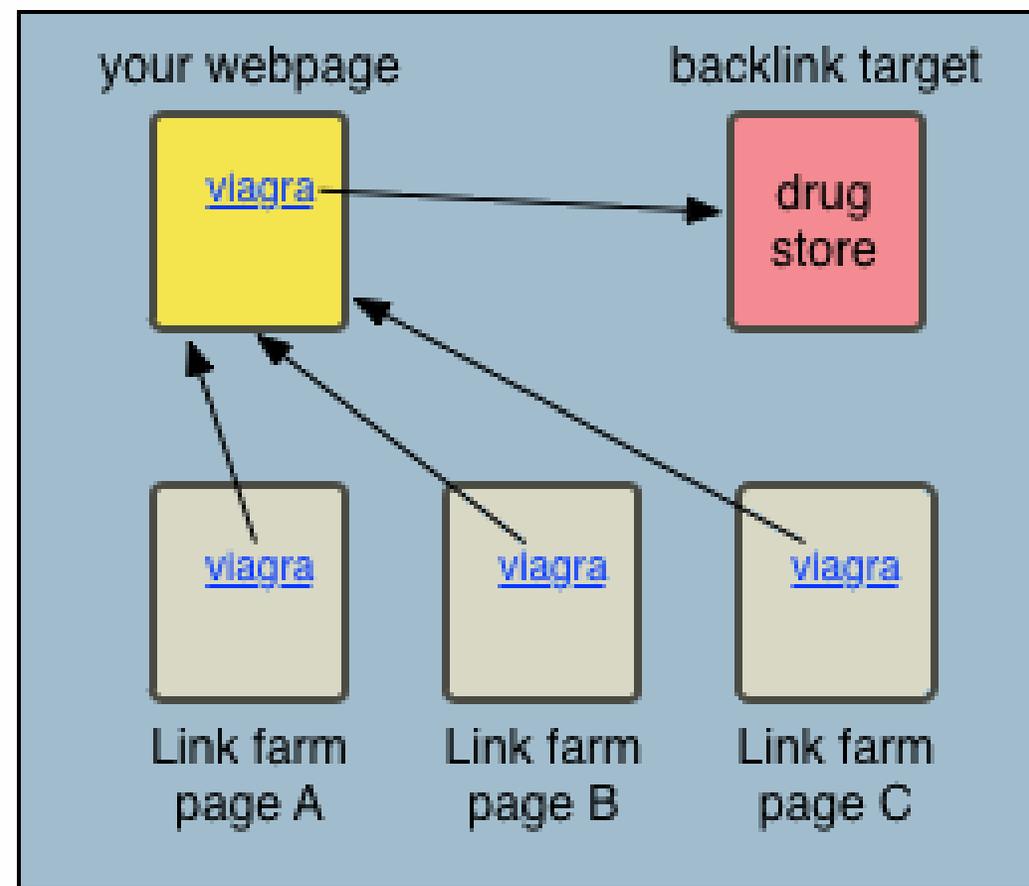
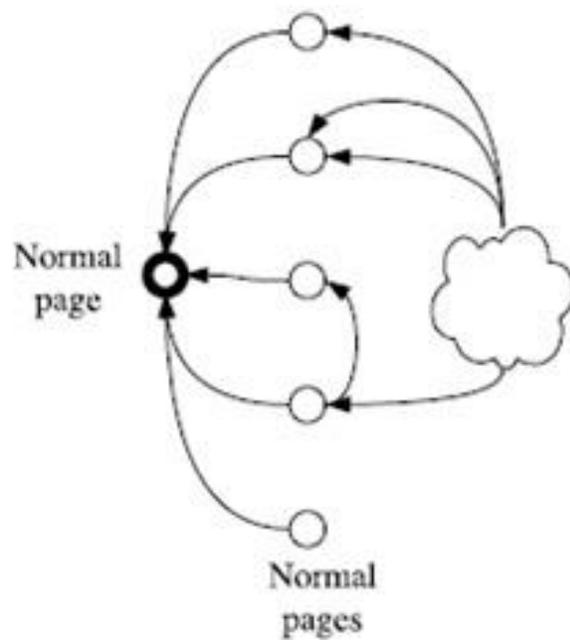
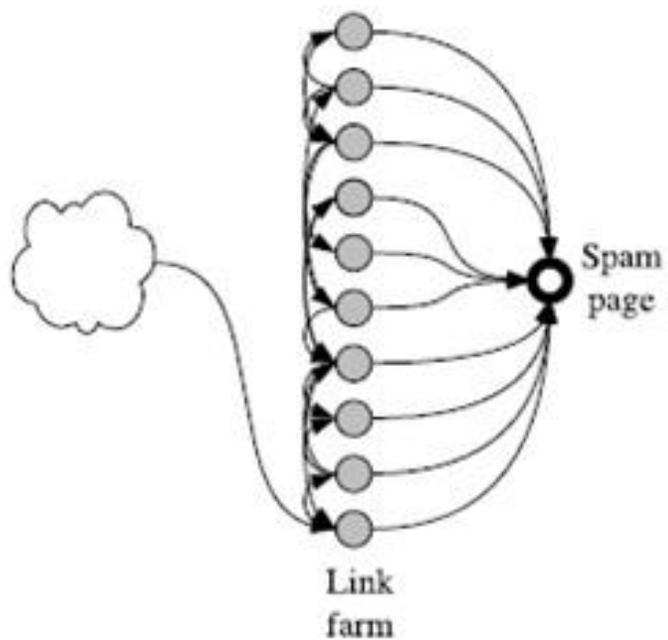
Firepits can create a beautiful ambiance for your backyard. Firepits can be created in just about every area of your home; however backyard firepits are the most common. No matter what climate you live in your home will be able to benefit from a firepit.

Now you might be asking yourself how is a firepit built and do I need to set in on anything? The perfect addition to any firepit is a fire pit table. Fire pit tables can be made in different materials and patterns which will enable you to effortlessly blend your firepit and fire pit table in with your backyard landscape.

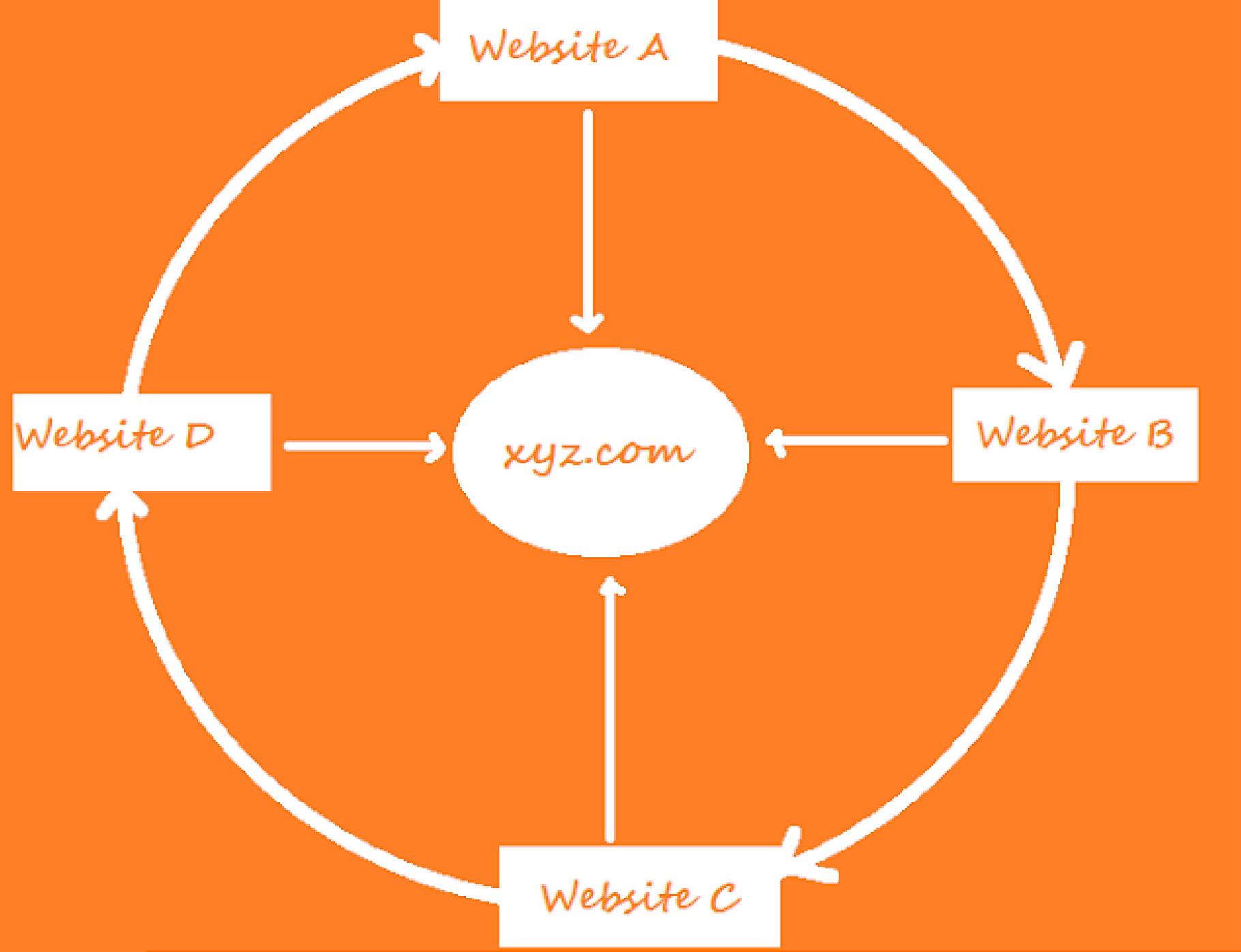
A very common type of fire pit table is the oriflammé fire pit table. The major benefits of these fire pit tables and the reason that they have become so popular is because they are lightweight, portable, and completely original. These tables are made in the USA and you can get your very own, one of a kind fire pit table! A major benefit of these tables is that they are easy to assemble, most will take approximately ten minutes to put together and the assembly process does not require any tools.

Another popular option for fire pit tables are wood burning fire tables. These tend to be smaller tables, much as you would have end tables next to your couch in the family room. These are small fire pit tables that are generally made out wrought iron. Wood burning fire pit tables are a great way to keep your back patio warm in those cooler months.

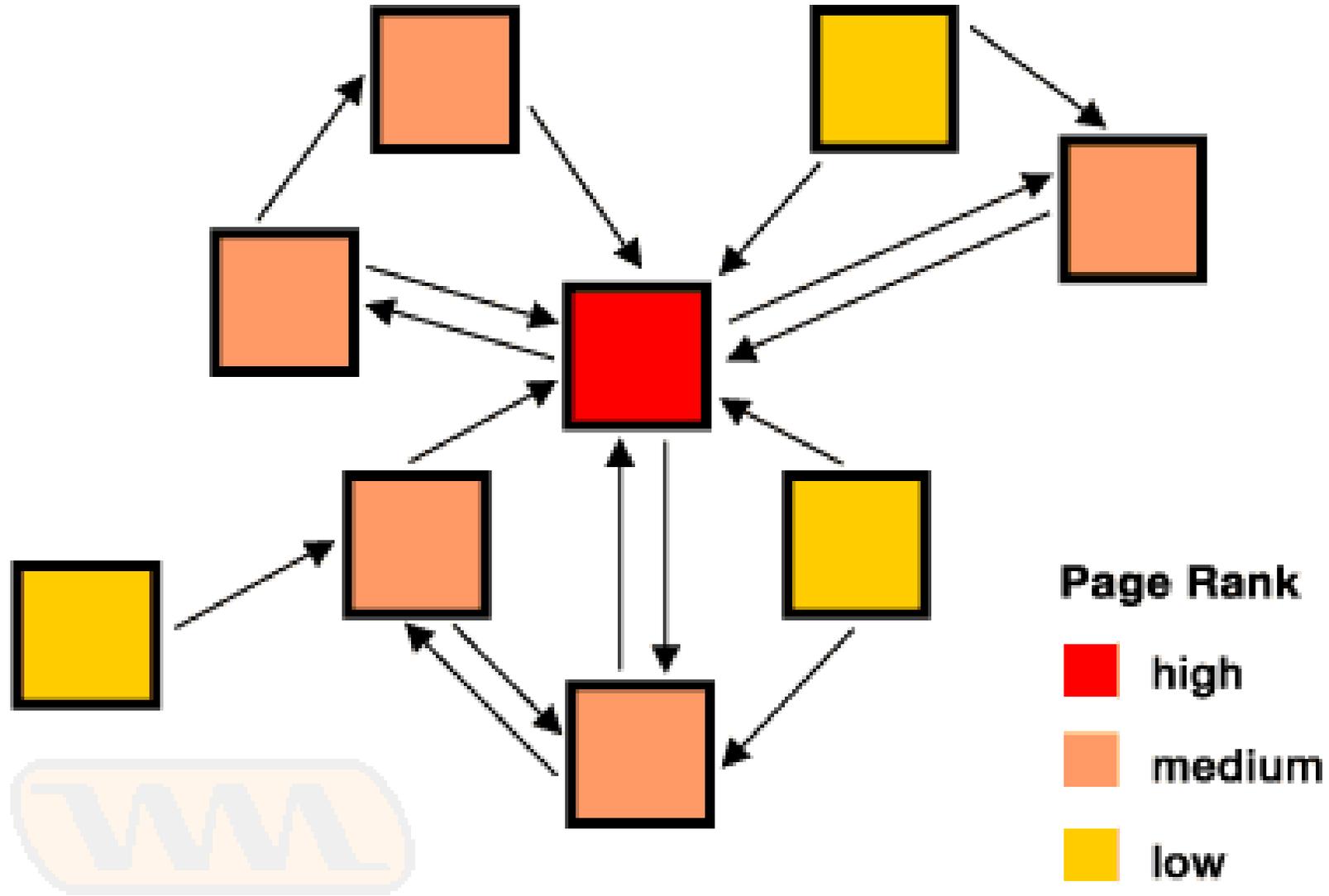
Fire pit tables also have another great function that can add some interesting décor to your backyard or front yard. The smaller fire pit tables are commonly used to enhance the beauty of their gardens and landscaping. A firepit and a fire pit table can make a grand impression on guests and can really liven up your homes garden.



**examples of Link farm backlinks**



# Links are an important factor in Page Rank



# Web blog spam

The screenshot shows a web browser window with the address bar displaying "http://www.english-jurist.edu/blog/index/1/". The page title is "Dr. B.'s Blog" and the subtitle is "A blog of classroom activities and discussions. A place where history rocks!".

**Left sidebar:**

- User login:** Username and Password fields with a "Log in" button. Links for "Create new account" and "Forgot your password".
- Links:** A list of links including "Archived blog posts", "Archived page", "Home Page", "Reading Materials Pictures", "ESL 2000 Pictures", and "ESL 2000 Books Pictures".
- Who's online:** "There are currently 1 user and 20 guests online."
- Online users:** A list containing "0/0".
- Navigation:** A list of navigation links: "Home", "Books", "Contact", "Search", "Categories", and "Application".
- Buttons and Tools:** "Feedback Profile".

**Main content area:**

**Home » Blog » Dr. B.'s Blog**

### Gaming Presentations

Submitted by [Dr. B.](#) on [Sat, 10/02/2009 - 9:11pm](#) **Conference: Game Theory - Not just another strategy topic: Research and Writing. This is what a National Game Day.**

Found out today that my papers were accepted for both the [USG Conference](#) and the [Transactions and Abstracts conference](#). The USG paper is almost done and is looking at historical representations of race and the T&A paper is on historical representations of gender.

It's exciting! USG gives me the chance to present my work to other folks who are working in the same area. Two days of papers on video games. I might just fund myself.

Feedback URL for this post:  
<http://www.english-jurist.edu/blog/feedback/2/>

**roulette**  
From roulette on Sat, 10/02/2009 - 9:11pm  
You are invited to check the sites in the field of [blackjack internet casino](#).

**riverbelle online casino**  
From riverbelle online casino on Sat, 10/02/2009 - 7:26am  
You are invited to check out some relevant information in the field of [video poker games on internet casino](#).

**premio online**  
From premio online on Fri, 10/02/2009 - 5:10pm  
You can also check the pages on [poker download video poker profiles](#).

**buy cheap online xanax**  
From buy cheap online xanax on Tue, 10/06/2009 - 9:03am  
You can also take a look at the pages on [doctor online prescription](#).

**ganar dinero internet**  
From ganar dinero internet on Mon, 10/05/2009 - 8:10am  
In your free time, check out the sites in the field of [online video poker game](#).

**discover card**  
From discover card on Sat, 10/03/2009 - 8:25am  
In your free time, check some information on [home loans mortgage rates bad credit loans](#).

**rebel strip poker downloa**  
From rebel strip poker downloa on Fri, 09/30/2009 - 3:10pm  
Please check out some relevant information dedicated to [Free Poker Tournaments](#).

**Right sidebar:**

**Current Class Blogs**

- [ENGL 304C Fall 09](#)
- [ENGL 305 Fall 09](#)

**Recent blog posts**

- [Archives of Oldergames](#)
- [Did I ever see the blood test before games?](#)
- [I Was Tired in the Court of Merit...](#)
- [I Used to Sleep in W and Wake Up in the Handmaid's Tale](#)
- [It's All About the Ethics](#)
- [Support William Shaw at the Age of 60](#)
- [Exercise Makes Me Sick!](#)
- [How to Email a Professor](#)
- [Copyright Book Released Under CC](#)
- [What's Going On](#)

[more](#)

**Recent comments**

- [My review?](#) 2 days 18 hours ago
- [Why don't you tell them what?](#) 2 days 18 hours ago
- [interesting ...](#) 2 days 18 hours ago
- [Congrat honors B and makes](#) 2 days 5 hours ago
- [Sublime!](#) 1 week 3 days ago
- [Spill please? food patterns?](#) 1 week 3 days ago
- [not practicing makes me crazy!](#)

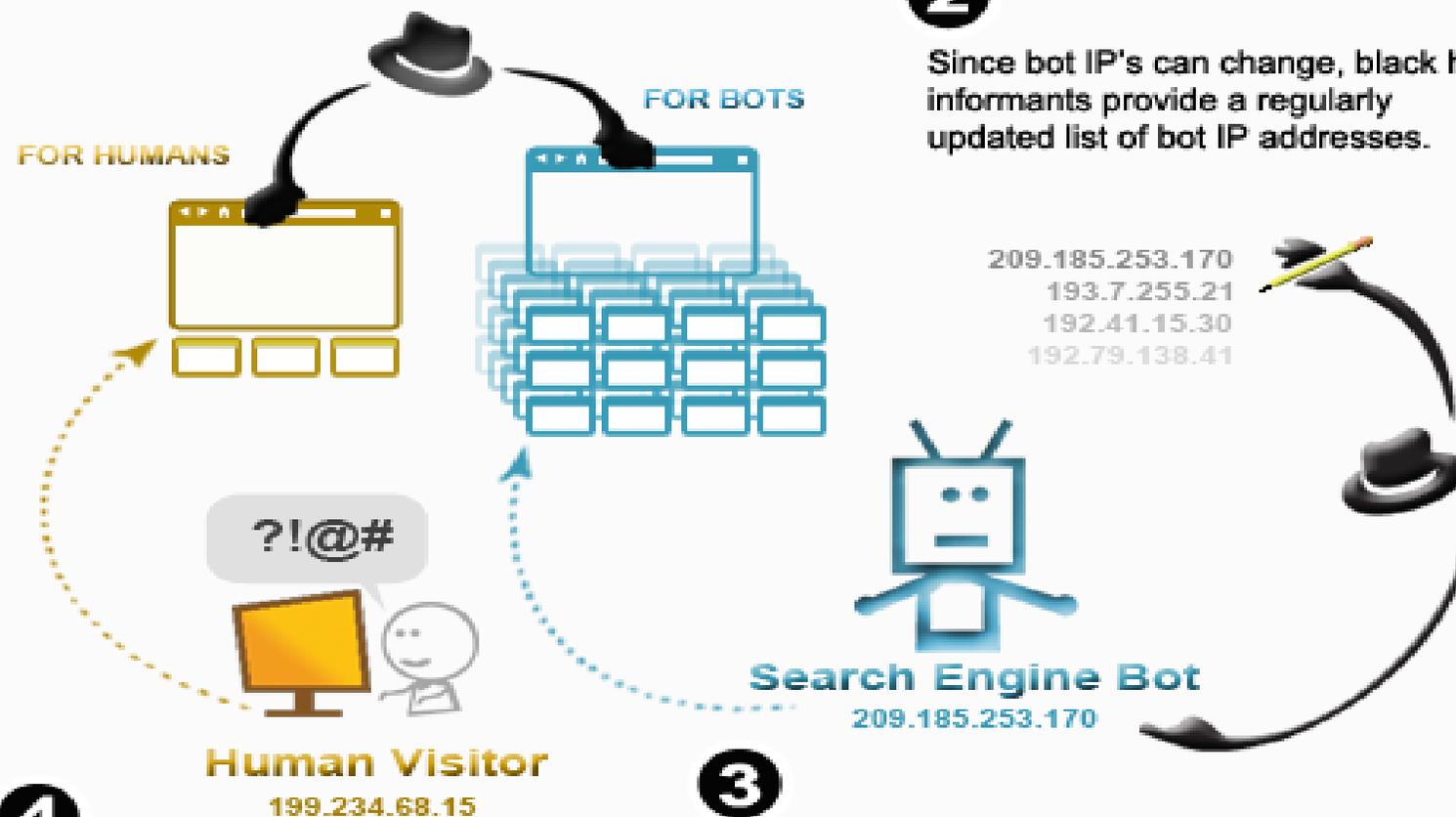
# Black Hat Cloaking Explained

1

Sites engaged in black hat SEO prepare two sets of content, one targeted for bots and the other targeted for human visitors. Bots are identified by their IP address.

2

Since bot IP's can change, black hat informants provide a regularly updated list of bot IP addresses.

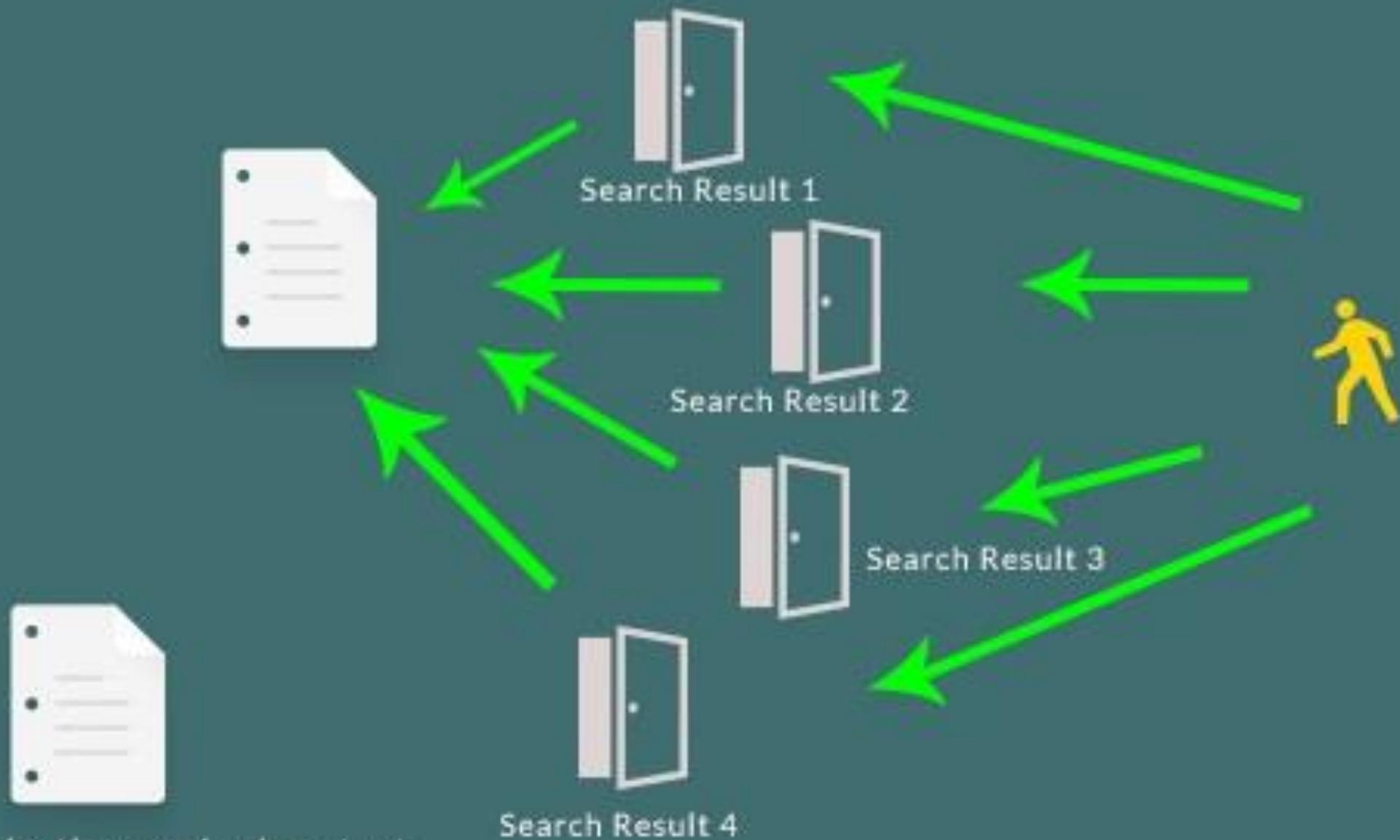


4

Human visitors often won't find the best information despite the site's high rankings.

3

Bots are served abundant fabricated content packed with targeted keywords. This false information boosts rankings.



This may be the required content

