

robots.txtを効果的に設定しよう

クローリングが不要な部分は robots.txt で回避する

“robots.txt”とは、検索エンジンにアクセスしクローリングしてほしい部分と、そうでない部分を伝えるためのファイルです (①)。

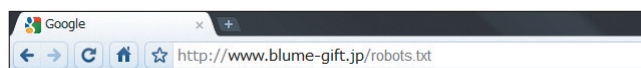
このファイルは必ず“robots.txt”というファイル名でサイトのルートディレクトリに置く必要があります (②)。

Googleウェブマスターツールをご利用いただくと、robots.txt ファイルを簡単に作成することができます。詳細は、ヘルプセンターの [robots.txt ファイルを使用してページをブロックまたは削除する](#) をご確認ください。サブドメインを持つサイトで、ある特定のサブドメイン内のページをクローリングさせないようにするには、そのサブドメイン用に別のrobots.txt ファイルを用意する必要があります。

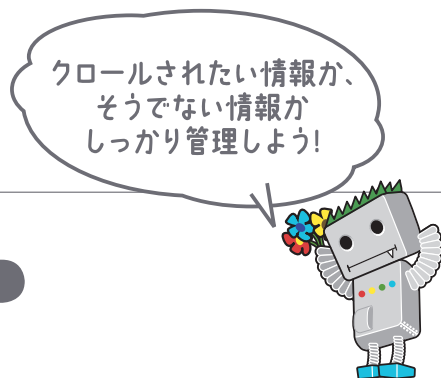
検索結果にコンテンツを表示させない方法は他にも、“NOINDEX”をrobotsメタタグに追加、[.htaccess](#) を使ってディレクトリにパスワードを設定、Googleウェブマスターツールを使ってすでにクローリングされたコンテンツを削除するなどがあります。

```
User-agent: *
Disallow: /image/
Disallow: /search
```

① Robots Exclusion Standard に準拠している検索エンジンのロボットすべて（「*」というワイルドカードのシンボルで表現される）に対し、/image/以下にあるコンテンツ、もしくは/searchで始まるURLにあるコンテンツに、アクセスもクローリングもさせない場合の例



② フラワーギフトショップのrobots.txtファイルのアドレス



ポイント

慎重に扱うべきコンテンツにはより安全な方法を使用しよう

機密事項や慎重に扱うべきコンテンツがクローリングされないようにするには、robots.txt の設置だけでは十分ではありません。その理由の1つは、クローリングできないように設定したURLであっても、そのURLへのリンクがインターネット上のどこか（例えば [リファラーログ](#) など）に存在する場合、検索エンジンはそのURLを参照できるからです。また、Robots Exclusion Standard に準拠しない検索エンジンや不正な検索エンジンなどは、robots.txt の指示に従わないかもしれません。そしてもう1つ、好奇心の強いユーザーの中には、robots.txt にあるディレクトリやサブディレクトリを見て、見られたくないコンテンツのURLを推測する人がいるかもしれません。コンテンツの暗号化や.htaccess を使ってパスワードをかけて保護する方が、より確実に安全な手段だといえます。

- 検索結果のようなページはクローリングさせない**
※検索結果のページから、さほど価値が変わらない別の検索結果のページへ飛んでも、ユーザーの利便性を損なうだけです
- 同一か、ほとんど違いがない自動生成されたページを大量にクローリングさせないようにする**
※「重複コンテンツに近いこれら100,000ものページはインデックスされるべきだろうか？」と問い直してみましょう
- プロキシサービスによって生成されたURLはクローリングさせないようにする**

.htaccess
ウェブサーバーの動作環境を制御するアクセス環境設定ファイル

リファラーログ
アクセスログに記載されているリファラー情報。これをたどっていくと閲覧者がどこのサイトから来たかなどを調べられる

プロキシサービス
内部ネットワークと外部ネットワークを接続する場合に接続を代行するコンピュータ、またはそのための機能を持ったソフトウェアのこと

参考ページ

ウェブマスター向けヘルプセンター
<http://www.google.co.jp/support/webmasters/>
↳ [robots.txtファイルを使用してページをブロックまたは削除する](#) 検索