

AND 検索が行われなかった時に使用時間間隔から求められた関連度よりも、AND 検索が行われた時に求められた関連度の方が、検索語間の関連度に及ぼす影響が大きいと考えられる。単に使用時間間隔のみから関連度を求めるのでは、本来関連のある検索語群（使用時間間隔が 300 秒以上であったとしても明らかに関連があると思われる検索語群^[1]）に対して高い関連度を付与できない可能性がある。逆に、関連の低い検索語群に対して誤って高い関連度を付与してしまう可能性もある。

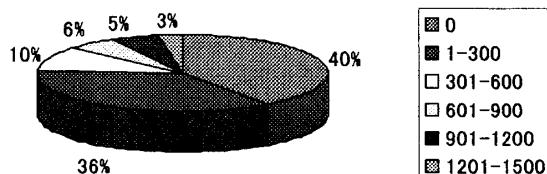


図3 2006/09/28の使用時間間隔ごとの検索回数の割合

4. 検索パターンの分類

3で述べた問題点を解決するために、検索語間の使用時間間隔のみから関連度を定義するのではなく、検索パターンも踏まえて関連度を定義することを考える。検索ログデータ内の各検索において、前後の検索語の変化を分析し、以下の5種類（グループ）に分類した。検索ログデータ内（利用者約40万人、全検索回数約120万回）において、5種類の検索パターンの存在する回数を求めたグラフが図4である。

グループ1：A-A（AND検索なし、前後で検索語が同一のパターン）

短い時間間隔の場合には検索語を入力後、結果として表示されるWebページのURLの一覧を早く得たいと思うことから、ダブルクリックを行ってしまっていると考えられる。長い時間間隔の場合には検索結果やWebページを閲覧した後、“次へ”ボタンにより次のページへ切り替えていると考えられる。

グループ2：A-B（AND検索なし、前後で検索語が異なるパターン）

短い時間間隔の場合には検索語の入力ミスに気づき、再入力している場合が考えられる。長い時間間隔の場合には検索結果やWebページを閲覧した後、検索語を変更した方が良いと判断し再入力を行っている、もしくは以前とは異なる情報を得るための新たな検索を行っていると考えられる。

グループ3：AB-AB（AND検索あり、前後で構成されている検索語が全て同一のパターン）

グループ1と同様な検索の特徴が表れる。

グループ4：AB-AC（AND検索あり、前後で検索語の一部分が同一のパターン）

検索語の一部にグループ2と同様な検索の特徴が表れる。

グループ5：AB-CD（AND検索あり、前後で構成されている検索語が全て異なるパターン）

短い時間間隔で新たな検索語を2語再入力することは難しいと考えられる。長い時間間隔の場合にはグループ2の長い時間間隔の場合と同様な検索の特徴が表れる。

1のパターンに関しては、同一の検索語となるため、関連度は定義しない。3のパターンに関しては、AとBの関

連度はAとBをAND検索した際の関連度として定義される。4のパターンに関しては、BとCは共にAとAND検索されている。従って、AB-AC間の使用時間間隔が大きくてもBとCの関連度は高くなると考えられる。このパターンの場合の関連度は、使用時間間隔に依存しない、一定の値として定義すべきであると考えられる。2、5のパターンに関しては、関連がある場合とない場合が考えられる。そこで、確実に関連があると考えられる3、4から利用者が同じ概念を表す情報を得るための検索にかかる平均検索時間を求める。この値を閾値とし、使用時間間隔が閾値以下であれば、関連度を使用時間間隔の単調減少関数として定義し、使用時間間隔が閾値以上であれば関連度を0と定義すべきであると考えられる。

例えば、A(t1)-A(t2)-AB(t3)-BC(t4)のように、順番に検索語が入力されたとする。ここで、A、B、Cは検索語、t1、t2、t3、t4はそれぞれの検索語が入力された時間を示している。同一の検索語Aが含まれるA(t1)-A(t2)-AB(t3)までを同じ概念を表す情報を得るための検索とし、この検索にかかった検索時間(t3-t1)を求める。検索ログデータからこの検索を抽出し、利用者1人当たりの平均検索時間を求め、全ての利用者の平均検索時間を求めた。平均検索時間は約1571秒であった。10年前には平均検索時間を約300秒としており、使用時間間隔が300秒以降の検索語群の関連度を0と定義していた。^[1]現在の平均検索時間は300秒より短くなるだろうと予想していたが、結果より非常に長いことが分かった。これは、明確な目的を定めずに検索を行っている利用者が増えているためであると考えられる。

検索回数(回)

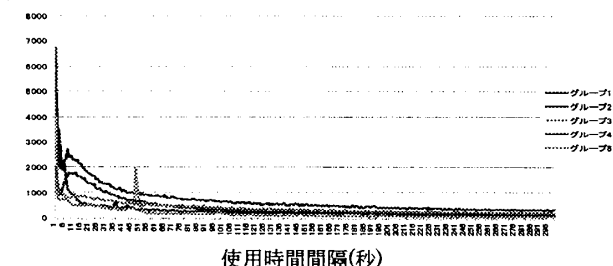


図4 2006/09/28の各検索パターンの存在する回数

5. 今後の課題

本稿で述べた関連度の定義を基にモデル化を行う。モデル化を行う上では、AND検索時における関連度の値、5種類の検索パターンの存在する割合からそれぞれのパターンの関連度の定義を全体のモデルにどう反映させるかを考慮する。

今後の展望として、新たなモデルを作成し、関連度の精度を向上させる。さらに、関連度を可視化することにより、最初に入力した検索語が最終的にどのような検索語となって利用者が目的の情報に辿り着いたのかを調査する。そして、これらを詳細な情報ニーズの抽出に結び付けたい。

参考文献

- [1] 大久保雅且, 井上孝史, 杉崎正之, 田中一男: www 検索ログに基づく情報ニーズの抽出, 情報処理学会論文誌, Vol. 39, No. 7, 1997.
- [2] 柳阿礼, 河村春雄, 徳永幸生, 杉崎正之, 池田成広: Web 検索ログの検索時間間隔を用いた利用者の行動パターンの分析, 第69回情報処理学会全国大会, 2007.