

修士論文

ランキング学習における クラスタリングを用いた学習データ選定による 学習効率向上手法

エルゲン 瑛夏

2025年1月30日

岐阜大学大学院 自然科学技術研究科 知能理工学専攻 知能情報学領域
鈴木研究室

本論文は岐阜大学大学院自然科学技術研究科
修士（工学）授与の要件として提出した修士論文である。

エルゲン 瑛夏

指導教員：

鈴木 優 准教授

ランキング学習における クラスタリングを用いた学習データ選定による 学習効率向上手法*

エルゲン 瑛夏

内容梗概

順序推定に用いられるランキング学習の主な3手法の中で、ペアワイズ手法は推定精度と使用データの柔軟性が高い手法として知られている。しかし、学習候補データ中の二つのデータの全ての組み合わせを用いて損失計算を行う必要があるため、学習候補データの増加に伴い教師データが指数的に増加してしまうという問題がある。この問題を解決するため、我々はペアワイズ手法において教師データ数が同じ場合に通常の学習手法と比較して順序推定精度を高め、学習効率を向上させることを目指す。我々は、教師データの学習難易度の段階的上昇が学習効率向上に有効であることに着目し、ペアの中のデータが持つ特徴量同士の距離が近いペアを学習難易度が高いと考え、段階的に近づけていくことを行う。提案手法では、候補データを特徴量に基づいて二つのクラスタにクラスタリングし、異なるクラスタ間でペアデータを作成する。次に、既存の各クラスタをさらに二つにクラスタリングし、直前のクラスタが同じだったデータ同士かつ、異なる新たなクラスタに分類されたデータ同士でペアを作成する。これらの操作を繰り返すことで、ペアの中のデータ同士が持つ特徴量を近づけていくことができ、学習難易度の段階的上昇が可能だと考えた。提案手法の有効性を確認するために評価実験を行った。結果、提案手法を用いて学習した場合にモデルの学習効率が向上することを確認できた。

キーワード

ランキング学習, ペアワイズ手法, 能動学習, カリキュラム学習

*岐阜大学大学院 自然科学技術研究科 知能理工学専攻 知能情報学領域 修士論文, 学籍番号: 1234525019, 2025年1月30日.

目次

図目次		iv
表目次		v
第 1 章	はじめに	1
第 2 章	基本的事項	4
2.1	ランキング学習	4
2.1.1	Top-k 学習	5
2.2	能動学習	6
2.3	カリキュラム学習	7
2.4	クラスタリング	8
2.4.1	K-means 法	8
2.5	NDCG@k	10
2.6	t 検定	11
第 3 章	関連研究	14
第 4 章	提案手法	16
4.1	Top-2 学習	18
4.2	教師データ作成	18
4.2.1	クラスタリング	19
4.2.2	ペアデータ作成	20
4.2.3	反復処理	20
第 5 章	評価実験	22
5.1	実験手順	22
5.2	提案手法の有効性確認	23
5.2.1	実験内容	23
5.2.2	結果・考察	24

第6章 おわりに	27
謝辞	29
参考文献	30
発表リスト	31

目次

2.1	能動学習の流れ	7
2.2	K -means 法の様子	9
2.3	t 分布	13
4.1	候補データ数 4 における提案手法の全体フロー	17
4.2	教師データ追加の様子	21
5.1	学習率が 10^{-4} , 10^{-5} の場合における NDCG@5 の推移	24

表目次

2.1	二つのモデルによる予測結果についての MSE と NDCG@4 による評価の比較	11
5.1	Fold 毎に収集したクエリ数の内訳	23
5.2	学習率を 10^{-4} とした場合の NDCG@5 の推移と p 値	25
5.3	学習率を 10^{-5} とした場合の NDCG@5 の推移と p 値	25

第 1 章 はじめに

我々は、順序推定を行うための機械学習手法であるランキング学習 [1] を用いて、検索結果をはじめとした順位付きデータの表示順序を利用者にとって適したものにすることを目的とする。ランキング学習とは、順序推定を行うデータの特徴量を入力した際に、並び替えた後の相対的な順序関係が正しくなるようなスコアを出力する機械学習モデルを作成する学習手法である。

モデルの構築に使用することができる候補のデータ全てを候補データ、候補データの中に存在する二つのデータの組み合わせをペアデータ、候補データの中でも実際に損失計算に用いるデータを教師データと呼ぶ。ランキング学習の手法は、損失関数の定義によって主にポイントワイズ手法 [1]、リストワイズ手法 [1]、ペアワイズ手法 [1] の三つが提案されている。ポイントワイズ手法では、一つひとつの教師データに対する推定モデルによる予測結果と各データに付与されている正解ラベルの誤差を損失計算に用いて学習を行う。ポイントワイズ手法のアルゴリズムには回帰、順序回帰、分類がある。回帰モデルの場合は、一つひとつの教師データに対して推定モデルによって算出された予測値と各データに付与されている正解ラベルとの誤差を小さくすることを目的に学習を行う。この際、推定モデルが各データの正解ラベルに近い値を予測できている場合でも、予測値の降順で並べたデータ順序と正解値の降順で並べたデータ順序が近くなるとは限らない。このように、一つひとつの教師データについての予測値を最適化するような損失計算手法では複数の教師データに対する順序予測を最適化することは難しい。

リストワイズ手法では損失計算に候補データの中に存在するデータの順列を用いるため、候補データ全体の順序推定を最適化するように順序推定モデルを学習することができ、3手法の中で最も精度が高くなることが期待できる。しかし、候補データ全てが順序付けされており、その順序が各データに正解ラベルとして付与されている必要があるため、データセットを準備することが困難である。

ペアワイズ手法では、候補データの中に存在する二つのデータを組み合わせるペアデータを作成し、そのペアの中でどちらのデータが優れているか、つまりペアの中で順序を予測するように学習を行う。リストワイズ手法の中でも、候補データの中に存在するデータの各順列において上位 2 位までに注目して損失計算を行う手法

はペアワイズ手法と一致する。

本研究ではこれら3手法のうち、データセットの準備の容易さからペアワイズ手法に着目したが、学習時間が膨大となる点が問題である。この手法では、学習に用いるペアデータの数は候補データの中に存在する二つのデータの全ての組み合わせの数となる。候補データ数の増加に伴い、学習に用いる計算量が $O(m \cdot n_{max}^2)$ に従って増加してしまうという問題がある。ここで m は候補データ中のクエリ数、 n_{max} はクエリ中の文書データ数である。

そこで本研究ではカリキュラム学習を用いて、ペアワイズ手法において教師データ数が同じ場合に、通常の学習手法と比較して順序推定精度を高め、学習効率を向上させることを目指す。カリキュラム学習とは、学習難易度が徐々に高まるような順番にデータを並び替えて学習することによって、学習難易度の高い複雑な領域を習得する以前に必要な前提知識が構成されている状態にし、モデルの学習効率を高める学習手法である。先行研究 [2] では、学習難易度は複数の要因から生じるため、学習難易度を定義する際の基準が複数になり、採用する基準によって同じデータでも学習難易度が変わってしまうという問題があると述べている。

我々はペアワイズ学習が各ペアデータの中に存在する二つのデータの順序を予測するように学習を行っていることに着想を得て学習難易度の定義を試みる。ペアデータの中に存在する二つのデータが持つ特徴量同士の距離が近い場合に、それらのデータに対する順序推定モデルによる予測値が近くなり順序推定が難しくなるため、学習難易度が高くなると考えた。本稿では学習難易度を徐々に高めるため、ペアデータの中に存在する二つのデータが持つ特徴量同士の距離が徐々に近くなるように教師データを並び替える。この難易度の定義方法はペアワイズ学習における学習アルゴリズムの特性に基づいているため、扱うタスク毎に難易度の定義を行う必要性がなく、ペアワイズ手法を用いる場合に広く応用可能であると考え。また、この難易度の定義方法であれば特徴量を参照することのみから候補データを難易度順に並び替えることができるため、データ準備に必要な追加の計算時間を抑えることができると考えた。

提案手法では最初に、候補データの中で特徴量が大きく離れた二つのデータからペアデータを作成するために、候補データを特徴量に基づいて二つのクラスターにクラスタリングし、異なるクラスター間でペアデータを作成する。次に、既に作成され

たペアデータと比較して、より特徴量が近いデータ同士のペアを作成する。そのために、既存の各クラスタ中のデータをさらに二つにクラスタリングし、直前のクラスタが同じだったデータ同士かつ、異なる新たなクラスタに分類されたデータ同士でペアを作成する。これらの操作を繰り返すことで、作成されるペアデータの中に存在する二つのデータが持つ特徴量同士の距離を徐々に近づけ、最終的に類似したデータ同士のペアが作成される。このデータ作成手法によって、教師データの段階的な学習難易度上昇と学習効率の向上を期待できると考えた。

提案手法の有効性を確認するために評価実験を行った。評価実験では提案手法を用いて作成した教師データと無作為な順序で作成した教師データを使用してペアワイズ学習を行い、作成された順序推定モデルの推定精度を比較した。学習効率向上の有無を確認するため、学習に用いるデータ数は教師データ全体に対して1割から10割までの範囲を1割間隔で推移させて実験を行った。

評価実験の結果、提案手法を用いてペアワイズ学習を行うことによって、学習効率が向上することを確認できた。

本稿による貢献は以下の通りである。

- ペアワイズ学習における学習難易度はペアデータの中に存在する二つのデータが持つ特徴量同士の距離の近さで適切に定義できる事を確認した。
- ペアワイズ学習において教師データの学習難易度を徐々に高めることが学習効率向上に有効であることを確認した。

第2章 基本的事項

2.1 ランキング学習

ランキング学習とは、予測スコアではなく予測順位の予測のためにモデルを学習する機械学習手法であり、検索エンジンによる検索結果の最適化を代表とした順序付きデータに用いられる。検索結果の表示順序決定には Web ページと検索ワードとの関連度をはじめ、ページの構造や内部リンクの数など多くのランキングを行うための要素を採用している。そのため、ランキングを行う際にそれぞれの要素がどの程度重要であるかを人手で判断するのが難しい。そこで機械学習を用いて各要素の重みを最適化する。

ランキング学習は損失関数の定義の仕方により大きく3種類に分類される。一つ目はポイントワイズ学習である。ポイントワイズ学習は、モデルが入力データの正解スコアを正確に予測できていれば順序も正確に予測できるという発想に基づいた手法である。損失関数は (2.1.1) 式のように表され、ある検索クエリ q 中の d 番目のデータ $x_{q,d}$ についての順序推定モデルによる予測スコア $f(x_{q,d})$ と、その正解スコア $y_{q,d}$ を用いた損失関数 $L(f(x_{q,d}), y_{q,d})$ を最小化する回帰問題を解いている。

$$L_{Pointwise} = \sum_{q,d} L(f(x_{q,d}), y_{q,d}) \quad (2.1.1)$$

二つ目はペアワイズ学習である。ペアワイズ学習とは任意の二つのデータからなるペアの優劣関係を正確に予測できれば文書群全体を正しく並び替えられるという発想に基づいて学習を行い、重みを最適化する手法である。損失関数は (2.1.2) 式のように表される。クロスエントロピー誤差などを用いて、モデルによるクエリ q 中の i 番目のデータ $x_{q,i}$ についての予測スコア $f(x_{q,i})$ と、 $x_{q,j}$ についての予測スコア $f(x_{q,j})$ を用いた損失関数 $L_{Pairwise}$ を最小化する2値分類問題を解いている。 m_q がクエリ q 中のデータ数を表し、 $y_{q,i}, y_{q,j}$ がクエリ q 中の i 番目と j 番目のデータのラベルを表す。

$$L_{Pairwise} = \sum_q \sum_{i,j: y_{q,i} > y_{q,j}}^{m_q} L(f(x_{q,i}), f(x_{q,j})) \quad (2.1.2)$$

三つ目はリストワイズ学習である。リストワイズ学習とは、モデルによる対象文書群に対する予測ランキングの正確さをそのまま損失関数として用いて学習を行う手法である。損失関数は複数クエリ中の全データから算出するため (2.1.3) のように表される。クエリ q 中のあるデータ $x_{q,1}$ についての予測スコア $f(x_{q,1})$ からクエリ q 中の全データ数 m_q のデータ x_{q,m_q} についての予測スコア $f(x_{q,m_q})$ それぞれと、それに対応した正解ラベル $y_{q,1}$ から y_{q,m_q} を用いて損失 $L_{Listwise}$ を (2.1.3) 式のように算出する。 m_q はクエリ q 中の全データ数を表す。損失関数は NDCG などのランキング問題用の評価指標を直接最小化する。

$$L_{Listwise} = \sum_q L((f(x_{q,1}), y_{q,1}), \dots, (f(x_{q,m_q}), y_{q,m_q})) \quad (2.1.3)$$

2.1.1 Top-k 学習

Top- k 学習とはリストワイズ学習における損失計算手法の一つで、候補データの中に存在するデータから作成した各順列の中に存在する上位 k 位までに注目して損失計算を行う手法である。学習に用いるデータは、検索ワードであるクエリと、クエリとの関連度に基づいた正解順位ラベルをもつデータで構成されているとする。Cao ら [3] は損失関数を (2.1.4) 式のように定義している。(2.1.4) 式では、正解順序発生確率 $p_{y^{(i)}}(x_j^{(i)})$ と、予測順序発生確率 $p_{z^{(i)}(f_\omega)}(x_j^{(i)})$ のクロスエントロピー誤差を計算している。正解順序発生確率 $p_{y^{(i)}}(x_j^{(i)})$ は i 番目のクエリ $q^{(i)}$ の正解順位リスト $y^{(i)}$ 中で j 番目のデータ $x_j^{(i)}$ が 1 位となる順序発生確率である。予測順序発生確率 $p_{z^{(i)}(f_\omega)}(x_j^{(i)})$ はニューラルネットワークモデル ω によるランキングモデル f_ω によって算出された $q^{(i)}$ 中のデータの予測スコアリスト $z^{(i)}(f_\omega)$ 中で、 $x_j^{(i)}$ が 1 位となる順序発生確率である。また、 $x_j^{(i)}$ がスコア $s(j)$ を持つ時、 $x_j^{(i)}$ が 1 位となる順序発生確率を (2.1.5) 式のように定義している。(2.1.5) 式では教師データのスコアの総和で s_j を割ることで、 $x_j^{(i)}$ が 1 位に選ばれる確率を算出している。

$$L(y^{(i)}, z^{(i)}(f_\omega)) = - \sum_{j=1}^{n^{(i)}} p_{y^{(i)}}(x_j^{(i)}) \log(p_{z^{(i)}(f_\omega)}(x_j^{(i)})) \quad (2.1.4)$$

$$P_s(j) = \frac{\exp(s_j)}{\sum_{m=1}^{n^{(i)}} \exp(s_m)} \quad (2.1.5)$$

Top-1 学習の利点は教師データの数が少なく、学習時間が短く抑えられることである。 $q^{(i)}$ 中に $(x_1^{(i)}, x_2^{(i)}, x_3^{(i)})$ の 3 データがある場合を例に挙げる。このとき、 $q^{(i)}$ 中の順列は $[x_1^{(i)}, x_2^{(i)}, x_3^{(i)}]$, $[x_1^{(i)}, x_3^{(i)}, x_2^{(i)}]$, $[x_2^{(i)}, x_1^{(i)}, x_3^{(i)}]$, $[x_2^{(i)}, x_3^{(i)}, x_1^{(i)}]$, $[x_3^{(i)}, x_1^{(i)}, x_2^{(i)}]$, $[x_3^{(i)}, x_2^{(i)}, x_1^{(i)}]$ となる。 $x_1^{(i)}$ が 1 位となる順序発生確率を求めるには、順序 $[x_1^{(i)}, x_2^{(i)}, x_3^{(i)}]$ と順序 $[x_1^{(i)}, x_3^{(i)}, x_2^{(i)}]$ の発生確率の和を取る必要がある。そのため、 $x_1^{(i)}$, $x_2^{(i)}$, $x_3^{(i)}$ それぞれが 1 位になる確率を求めるためには順列全てを計算に用いる必要があり、教師データの総数は候補データの中に存在するデータ数の階乗となる。ここで、順序 $[x_1^{(i)}, x_2^{(i)}, x_3^{(i)}]$ と $[x_1^{(i)}, x_3^{(i)}, x_2^{(i)}]$ の発生確率の和は $x_1^{(i)}$ が 1 位になる順序発生確率と等しいため、候補データに存在するデータの中で $x_1^{(i)}$ が 1 位になる順序発生確率のみから算出可能である。そのため Top-1 学習では候補データの中に存在するデータの順列全てではなく、候補データ中の 1 データの順列のみを教師データとして用いる。候補データの中に存在するデータ数を D とすると、Top-1 学習の教師データ数は ${}_D P_1$ となる。例では、 $(x_1^{(i)}, x_2^{(i)}, x_3^{(i)})$ それぞれが 1 位になる順序発生確率を計算するだけで良いため、教師データの総数は 3 になる。

2.2 能動学習

能動学習 [4] とは機械学習の学習手法の一つで、モデルの学習に有効なデータから優先して学習に使用する手法である。優先して学習に使用するデータの選定は学習アルゴリズムによって行われる。この学習方法によって学習効率の向上が期待できるため、データセット作成や計算コストを削減することができる。代表的な能動学習の型の一つにプールベース能動学習がある。この方法では、ラベル無しデータは簡単に収集できるが、ラベル付けのコストが大きい場合を問題設定として能動学習を行う。図 2.1 に能動学習の流れを示す。能動学習は以下の四つの段階から構成される。手順 4 における終了条件は、モデルの精度や学習に使用したデータ数などから設定される。

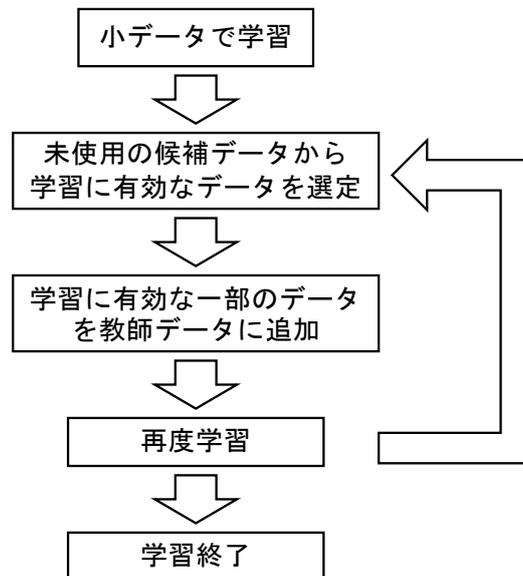


図 2.1 能動学習の流れ

1. 候補データ中に存在する一部のデータを用いてモデルを学習する.
2. 構築されたモデルや何らかのアルゴリズムを用いて, 未使用の候補データから学習に有効なデータを選定する.
3. 選定されたデータを教師データに追加して再度モデルを学習する.
4. 2. と 3. の手順を終了条件を満たすまで繰り返す.

2.3 カリキュラム学習

カリキュラム学習とは機械学習手法の一つで, モデルにとっての学習難易度が低いデータから学習を始めて, 徐々に学習させるデータの学習難易度を高くしていく学習手法である. この学習手法は人間の学習プロセスから着想を得ており, 子供が数学を学ぶ際に足し算や引き算といった学習難易度の低いものから学習を始め, 徐々に方程式や微積分といった難易度の高いものを学習するという学習手法に似て

いる。この過程を辿ることで、モデルが学習難易度の高い複雑な領域を学習する以前に必要な前提知識が構成されている状態にし、モデルの学習効率を高める効果を期待できる。この学習方法では無作為なデータの順序で学習を行う場合に比べ、小データで高い予測精度のモデルを作成することができ、学習効率を高めることができる。一方で、扱うタスクによって難易度を定義する際の基準が異なるため、難易度の定義が難しいという問題がある。

2.4 クラスタリング

クラスタリングとは教師なし学習の一つで、データをクラスタと呼ばれる似た特徴量を持つグループに分割する手法である。グループ分けに際してデータが持つ特徴量のみを必要とするため人間によるラベル付けを必要とせず、自動的にデータをグループに分割することができる。

2.4.1 K-means 法

K-means 法 [5] とはクラスタリングの代表的な手法の一つである。この手法は四つの異なる段階から構成される。K-means 法の様子を図 2.2 に示す。散布図中の白い記号はクラスタリング対象のデータを表し、黒い記号は各クラスタの重心を表す。1 段階目ではアルゴリズムがデータセットに対してランダムにクラスタの割り当てを行う。この際のクラスタ数は事前に設定することができる。図 2.2 では三つのクラスタにクラスタリングする場合を例にあげている。2 段階目にアルゴリズムが各クラスタの重心を算出する。3 段階目では各データと各クラスタの重心との距離が計算される。4 段階目では各データが最も近いクラスタに再配置される。この第 3 段階と第 4 段階の操作は、再配置後のデータのクラスタに変化が無くなるまで繰り返される。

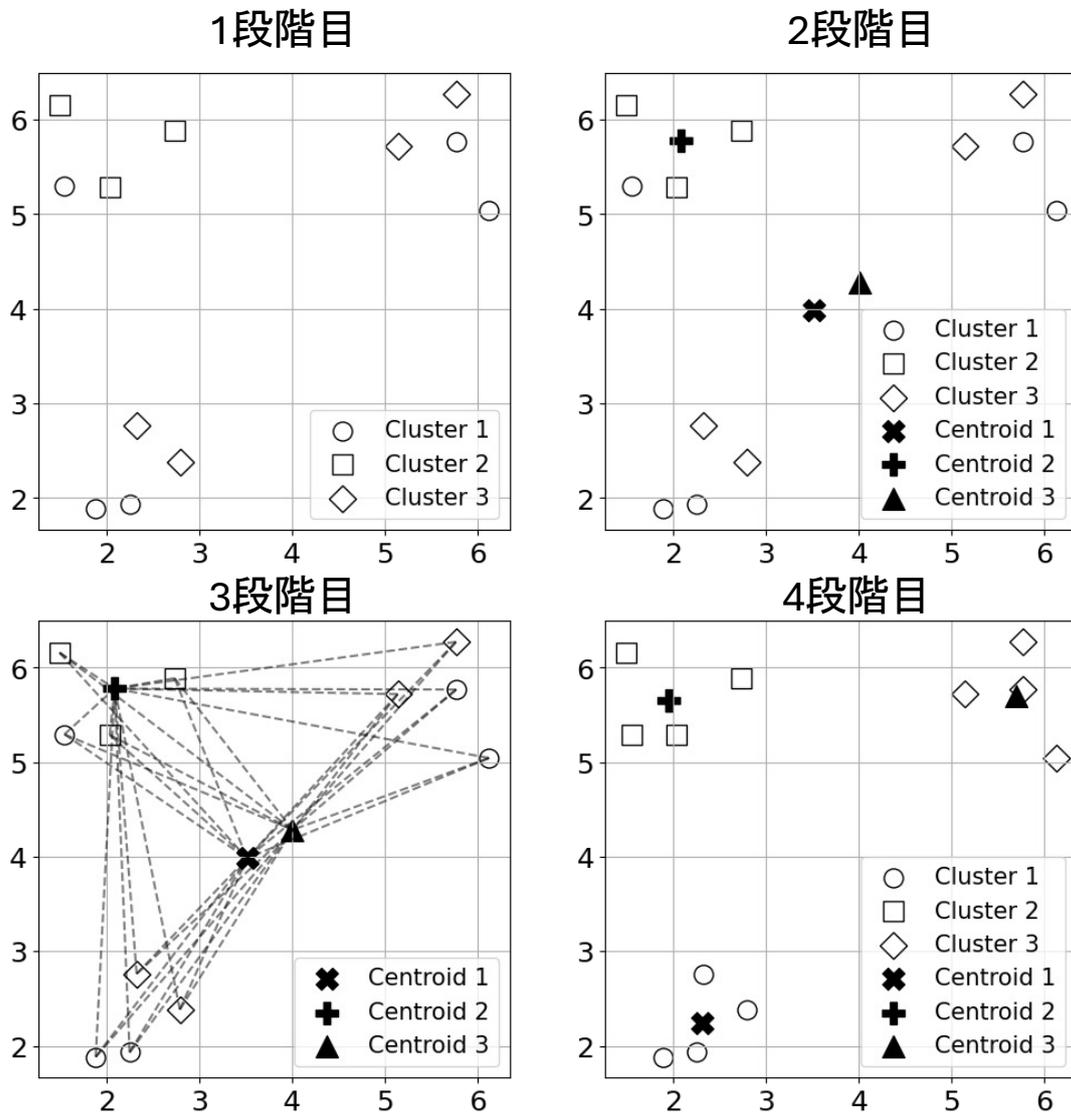


図 2.2 K -means 法の様子

2.5 NDCG@k

NDCG(Normalized Discounted Cumulative Gain)とは、順位に応じたスコア割引を行う評価指標である DCG(Discounted Cumulative Gain) を予測スコアについて算出した (DCG_{pred}) を、理想の DCG(DCG_{ideal}) で割ることで正規化したものである。理想の DCG とは、あるデータ群を正解順序スコアの降順に基づいて並び替えた際の順序と、予測順序スコアの降順に基づいて並び替えた際の順序が一致している場合の DCG である。

NDCG@k を (2.5.1) 式に示す。NDCG@k はスコアの降順について上位 k 位までを計算する際に使用する NDCG のことである。0 から 1 の値を取り、1 に近い方が精度が高い。

DCG@k を (2.5.2) 式に示す。DCG@k はスコアの降順について上位 k 位までを計算する際に使用する DCG のことである。予測順位が r 番であるデータの予測スコア $g(r)$ を順位に応じたスコア割引である $\log_2(r + 1)$ で割ることで算出した値を、 r が 1 から k 番目について和をとることで算出する。ランキング学習の目的である上位データの推定精度を評価するために、下位の順位にスコア割引を行うこの評価指標を採用した。

$$NDCG@k = \frac{DCG_{pred}@k}{DCG_{ideal}@k} \quad (2.5.1)$$

$$DCG@k = \sum_{r=1}^k \frac{g(r)}{\log_2(r + 1)} \quad (2.5.2)$$

表 2.1 に NDCG 計算の具体例を示す。表 2.1 では文書 1 から 4 までの四つのデータについて、モデル A とモデル B の 2 モデルを用いてスコア予測を行った結果と、スコアを降順で並べた順位を () 内に示している。また、それぞれのモデルの予測結果を MSE(Mean Squared Error) と NDCG@4 を用いて評価した場合の値を示し、優れているモデルの値を太字で示している。モデル A の予測結果を回帰タスクで広く用いられている評価指標である MSE を用いて評価すると 5 であり、モデル B による予測結果を評価すると 10.935 をとる。データのスコアを正しく予測

表 2.1 二つのモデルによる予測結果についての MSE と NDCG@4 による評価の比較

	正解スコア (順位)	モデル A による 予測スコア (順位)	モデル B による 予測スコア (順位)
文書 1	5(1)	2(4)	0.5(1)
文書 2	4(2)	3(3)	0.4(2)
文書 3	3(3)	4(2)	0.3(3)
文書 4	2(4)	5(1)	0.2(4)
MSE	-	5	10.935
NDCG@4	-	0.625	1

することを目的とする回帰タスクであれば、モデル A が優れたモデルであると評価できる。一方で NDCG を用いて評価すると、モデル A の予測結果が 0.625、モデル B による予測結果は 1 をとる。データの順序を正しく予測することを目的とするランキング学習であれば、モデル B が優れたモデルであると評価できる。

2.6 t 検定

t 検定とは、二つのグループの平均値の違いが偶然なのか、統計的に有意な差があるのかを判定するための統計的検定手法である。 t 検定には大きく分けて 2 種類の検定手法が存在する。一つ目は対応のある t 検定で、同じ検定対象についてある操作を加える前後での比較を行うものである。二つ目は対応のない t 検定で、独立した二つのグループの比較を行う。

t 検定では最初に帰無仮説として「検定で比較している二つのグループに統計的な差がない。」と仮説を設定する。また、帰無仮説が棄却された際に成立する対立仮説を「検定で比較している二つのグループに有意な差がある。」と設定する。次に有意水準という、帰無仮説を棄却する基準となる確率を設定する。有意水準には一般的に 5% が用いられる。ここで設定した有意水準を p 値が下回った場合に帰無仮説は棄却される。 p 値とは、帰無仮説が正しいと仮定した場合に、観測されたデータ以上の極端なデータが得られる確率のことで、(2.6.1) 式を用いて算出する t 値を用いて導出される。(2.6.1) 式では \bar{X}_1 がグループ 1 の平均値、 \bar{X}_2 がグループ

2 の平均値を表し、 s_1^2 がグループ 1 の不偏分散、 s_2^2 がグループ 2 の不偏分散、 n_1 がグループ 1 のサンプル数、 n_2 がグループ 2 のサンプル数を表す。不偏分散 s^2 は (2.6.2) 式で導出される。 n がサンプルサイズ、 x_i が各データ、 \bar{x} が平均値を表す。 p 値は t 分布を元に導出される。 t 分布を図 2.3 に示す。横軸に t 値、縦軸に確立密度を示す。また、黒塗りの範囲が帰無仮説の棄却域を示す。図 2.3 では両側検定を行っている。両側検定とは、 t 値が極端に大きい場合と極端に小さい場合どちらでも帰無仮説を棄却する検定手法である。両側検定では設定した有意水準を 2 等分し、 t 分布の左右に棄却域として設定する。ここで求められた p 値と有意水準を比較し、帰無仮説を棄却するかを決定する。

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (2.6.1)$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2.6.2)$$

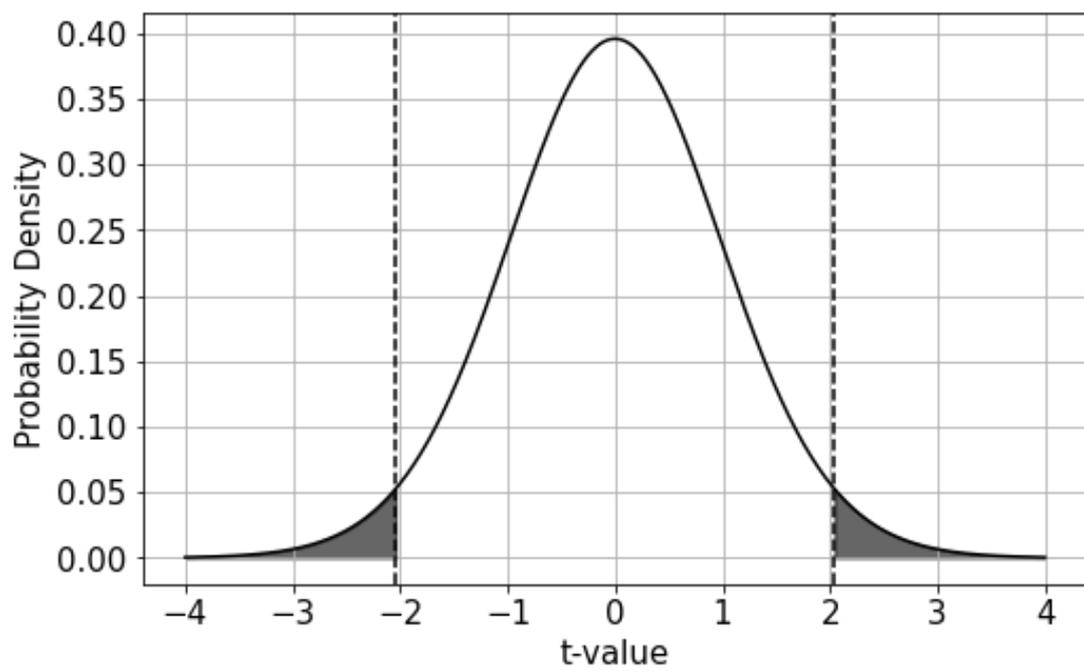


图 2.3 t 分布

第 3 章 関連研究

Cao ら [3] はリストワイズ学習のための Top- k 学習を提案した. Top- k 学習は候補データの中に存在するデータの各順列のうち上位 k 位までに着目して損失計算を行う手法である. k の増加に伴って学習済みモデルによる順序推定精度が高まることを示唆していたが, 学習時の計算量が膨大になることを避けるため Top-1 学習を行っている. 著者らは $k = 2$ の時に Top- k 学習はペアワイズ学習と一致すると述べており, 我々はリストワイズ学習への応用を考えて $k = 2$ とした Top- k 学習を用いてモデルの学習を行う.

上位の予測精度を重視したランキング学習についての研究がある. Liang ら [6] は, ランキング学習において下位データを学習に加えることは上位データの推定精度低下に繋がると考え, 実用的でないとした. そこで Liang らはランク付けされたリストの上位 N 個のみを考慮して損失計算を行い, 下位データの影響を排除した. これによって精度の向上が確認された. 本稿では, 候補データ中に存在する各クエリから上位 16 データずつを収集し学習に用いる. この操作によって, ランキング学習の目的である上位データの順序推定精度の向上を実現できると考えた.

教師データの学習難易度を徐々に高めることによる学習効率向上手法についての研究もある. Bengio ら [7] は, 学習時のデータがランダムに提示されるのではなく, 徐々に多くの概念や複雑な概念を示すような順序で提示することが学習効率を高めることを示した.

Elman ら [8] は, 言語処理タスクの学習難易度は複数の要因から生じると述べている. 著者らは言語処理タスクの解決のためには, 節や句といった文の階層構造や文の中で離れた位置にある単語同士の文法的な関係性を理解する必要性, 人称や時制といった文法要素間の一致の必要性を挙げている. 言語処理タスクにおける学習難易度は特にデータによる記憶負荷に関係しており, モデルが長い文や複雑な構文について全体の内容を保持しながら正しく処理する必要がある際に学習難易度が高まると述べている. また, データの多様性によっても学習難易度は高まると述べている. Zaremba ら [9] は入力したプログラムを読み取り, 出力結果予測を行うモデルを作成する際にカリキュラム学習を用いた. この際に Zaremba らは数値や文字の桁数, 入れ子構造の深さ, 非線形計算の有無などを学習難易度の構成要素として

挙げている。Graves ら [2] はタスクの学習難易度は複数の要因から生じるため、学習難易度を定義する際の基準が複数になり、採用する基準によって同じデータでも学習難易度が変わってしまうという問題があると述べている。そこで、あるサンプルデータを学習させる前後での損失関数の減少量や、モデルの複雑性の増加量から学習難易度を定義している。

我々は、ペアワイズ学習に用いる各ペアデータの中に存在する二つのデータが持つ特徴量同士の距離の近さを学習難易度として定義し、カリキュラム学習を行った。ペアワイズ学習では教師データ中の各ペアデータについて、ペアの中に存在する二つのデータの予測値と正解ラベルの差を損失計算に用いる。よって、ペアの中に存在する二つのデータが持つ特徴量同士の距離が近い場合に予測値が近くなり順序推定が難しくなるため、学習難易度が高くなると考えた。カリキュラム学習についての先行研究ではタスク毎に難易度の定義を試みているものが多いが、本稿では機械学習に使用するアルゴリズムに基づいた難易度の定義をしている。この難易度の定義方法であれば扱うタスク毎に難易度を定義する必要がなく、ペアワイズ学習を学習アルゴリズムに採用する場合に広く応用できると考えた。また、候補データの中に存在するデータのクラスタリング結果からデータの学習順序を決めることができ、カリキュラム設計に必要な計算量を小さく抑えることができると考えた。

第4章 提案手法

我々はランキング学習の1手法であるペアワイズ学習を用いて順序推定モデルを作成することを考えた。この際に膨大になる教師データを、カリキュラム学習を用いることで、教師データ数が同じ場合に通常の学習手法と比較して順序推定精度を高め、学習効率を向上させることを目指す。カリキュラム学習とは教師データの学習難易度が徐々に高まるような順序で学習を行う手法である。

我々は、ペアワイズ学習が各ペアデータの中に存在する二つのデータの順序を予測するように学習を行っていることに着想を得て、学習難易度の定義を試みる。ペアデータの中に存在する二つのデータが持つ特徴量同士の距離が近い場合に、それらのデータに対する順序推定モデルによる予測値が近くなり順序推定が難しくなるため、学習難易度が高くなると考えた。本稿では学習難易度を徐々に高めるため、ペアデータの中に存在する二つのデータが持つ特徴量同士の距離が徐々に近くなるように教師データを並び替える。

提案手法では最初に、候補データの中で特徴量が大きく離れた二つのデータからペアデータを作成するために、候補データを特徴量に基づいて二つのクラスターにクラスターリングし、異なるクラスター間でペアデータを作成する。次に、既に作成されたペアデータと比較して、より特徴量が近いデータ同士のペアを作成する。そのために、既存の各クラスター中のデータをさらに二つにクラスターリングし、直前のクラスターが同じだったデータ同士かつ、異なる新たなクラスターに分類されたデータ同士でペアを作成する。これらの操作を繰り返すことで、作成されるペアデータの中に存在する二つのデータが持つ特徴量同士の距離を徐々に近づけ、最終的に類似したデータ同士のペアが作成される。このデータ作成手法によって、教師データの段階的な学習難易度上昇と学習効率の向上を期待できると考えた。

提案手法の概要を図4.1に示す。以下に提案手法を用いた学習の手順を示す。

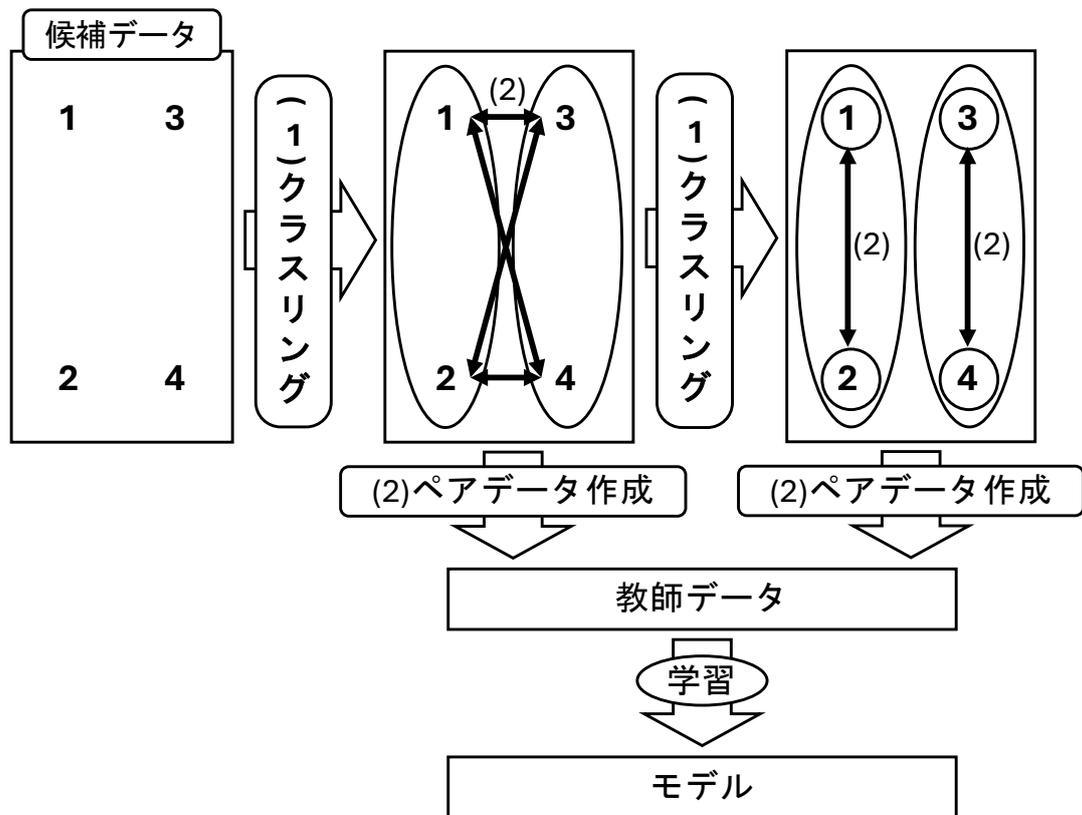


図 4.1 候補データ数 4 における提案手法の全体フロー

- (1) 候補データの中に存在するデータを二つのクラスタにクラスタリングする.
- (2) クラスタリング後の異なるクラスタ間に存在する二つのデータの順列を教師データに追加する.
- (3) (1),(2) の操作を (1) の操作前に全てのクラスタの中に存在するデータ数が一
つになるまで繰り返す.
- (4) 作成された教師データを用いて順序推定モデルを学習する.

4.1 Top-2 学習

本稿では Top- k 学習の k が 2 の場合で学習を行うため、損失関数は (4.1.1) 式を用いる. (2.1.4) 式からの差異として, j の範囲を ${}^{(i)}P_k$ に変更し, $x_j^{(i)}$ を $\pi_j^{(i)}$ に変更した. ここで ${}^{(i)}P_k$ は Top- k 学習における $q^{(i)}$ 中の教師データの総数であり, $\pi_j^{(i)}$ は $q^{(i)}$ 中に含まれるデータの順列に含まれる j 番目の順序データである. また, 順序発生確率は (4.1.2) 式を用いて計算した. $q^{(i)}$ 中に $(x_1^{(i)}, x_2^{(i)}, x_3^{(i)})$ の 3 データがある場合を例に挙げる. このとき, Top-2 学習に用いる $q^{(i)}$ 中の全教師データは $[x_1^{(i)}, x_2^{(i)}], [x_2^{(i)}, x_1^{(i)}], [x_1^{(i)}, x_3^{(i)}], [x_3^{(i)}, x_1^{(i)}], [x_2^{(i)}, x_3^{(i)}], [x_3^{(i)}, x_2^{(i)}]$ となる. $(x_1^{(i)}, x_2^{(i)}, x_3^{(i)})$ がスコア (s_1, s_2, s_3) を持つ場合を例に挙げる. このとき, 順序 $\pi = [x_1^{(i)}, x_2^{(i)}, x_3^{(i)}]$ の順序発生確率 $P_s(\pi)$ は (4.1.3) 式のように計算する.

$$L(y^{(i)}, z^{(i)}(f_\omega)) = - \sum_{j=1}^{{}^{(i)}P_k} p_{y^{(i)}}(\pi_j^{(i)}) \log(p_{z^{(i)}}(f_\omega)(\pi_j^{(i)})) \quad (4.1.1)$$

$$P_s(\pi) = \prod_{l=1}^n \frac{\exp(s_{\pi(l)})}{\sum_{m=l}^n \exp(s_{\pi(m)})} \quad (4.1.2)$$

$$P_s(\pi) = \frac{\exp(s_1)}{\exp(s_1) + \exp(s_2)} \frac{\exp(s_2)}{\exp(s_2)} \quad (4.1.3)$$

ペアワイズ学習は候補データの中に存在する二つのデータの順列を教師データとして用いるため, 計算量は訓練データのクエリ数 m とクエリ中の文書データ数 n_{max} を用いて $O(m \cdot n_{max}^2)$ と表すことができる. よって, 候補データのクエリ数やクエリ中データ数の増加に伴って計算量は指数的に増加してしまう. そこで我々はペアワイズ学習における学習効率の向上を実現し, 少ない教師データから高い予測精度をもつ順序推定モデルを作成することを考えた.

4.2 教師データ作成

提案手法では最初に, 候補データの中で特徴量が大きく離れた二つのデータからペアデータを作成するために, 候補データを特徴量に基づいて二つのクラスにク

ラスタリングし，異なるクラスタ間でペアデータを作成する．次に，既に作成されたペアデータと比較して，より特徴量が近いデータ同士のペアを作成する．そのために，既存の各クラスタ中のデータをさらに二つにクラスタリングし，直前のクラスタが同じだったデータ同士かつ，異なる新たなクラスタに分類されたデータ同士でペアを作成する．これらの操作を繰り返すことで，作成されるペアデータの中に存在する二つのデータが持つ特徴量同士の距離を徐々に近づけ，最終的に類似したデータ同士のペアが作成される．このデータ作成手法によって，教師データの段階的な学習難易度上昇と学習効率の向上を期待できると考えた．

4.2.1 クラスタリング

本節では候補データを二つのクラスタにクラスタリングする操作について説明を行う．この操作は図 4.1 の (1) で行う．

クラスタリングはクラスタ数を 2 として K -means 手法 [10] を用いて行った．クラスタリング結果について，クラスタ数が見つからなかった場合はランダムにクラスタ分割を行った．また，再現性の確保のため初期クラスタの中心を固定してクラスタリングを行った．

図 4.1 にクラスタリングの様子を示す．ある検索クエリに対する検索結果である 4 種類のウェブサイトを対象としてクラスタリングを行った場合を例に挙げる．図 4.1 では 4 種類のサイトをそれぞれ数字で表し，円に囲まれた領域がひとつのクラスタを表している．最初に，特徴量の大きく離れた二つデータ群を作成するためにデータ全体をサイト 1,2 からなるクラスタとサイト 3,4 からなる二つのクラスタにクラスタリングする．2 回目のクラスタリングでは，1 回目で作成された二つのクラスタ間よりも特徴量距離が近づいた二つのデータ群を作成するために，既存の各クラスタを対象に再度二つにクラスタリングを行う．図 4.1 ではサイト 1 からなるクラスタ，サイト 2 からなるクラスタ，サイト 3 からなるクラスタ，サイト 4 からなるクラスタの四つのクラスタにクラスタリングされる．実際の実験において全ての場合でクラスタの中に存在するデータ数が均等になるとは限らない．

4.2.2 ペアデータ作成

本節では 4.2.1 節のクラスタリング結果に基づいてペアワイズ学習に用いるペアデータの作成方法を説明する。図 4.1 の (2) の操作でクラスタリング結果からペアデータを作成する操作を行う。

図 4.1 にペアデータ作成の様子を示す。図中の両矢印によって結ばれた二つのデータからペアデータが作成される。ペアデータは異なるクラスタにクラスタリングされた二つのデータから作成する。この際二つのデータは総当たりによって選ばれ、教師データには順序を考慮した順列データが追加される。2 回目以降のクラスタリング結果からペアデータを作る際は、直前のクラスタが同じだったデータ同士かつ、異なる新たなクラスタに分類されたデータ同士でペアを作成する。図 4.1 のクラスタリング結果から作成される教師データをサイト番号で表すと、1 回目のクラスタリング結果から $[1,3],[3,1],[1,4],[4,1],[2,3],[3,2],[2,4],[4,2]$ の 8 データとなる。2 回目のクラスタリング結果からは $[1,2],[2,1],[3,4],[4,3]$ の 4 データとなる。

4.2.3 反復処理

本節では 4.2.1 節と 4.2.2 節で説明を行った一連の操作を繰り返し行う操作について説明を行う。提案手法では最初に候補データに対して 4.2.1 節で説明した方法でクラスタリングを行う。次にクラスタリング結果に対して 4.2.2 節で行った方法でペアデータを作成し、教師データに追加する。ペアデータの作成と教師データへの追加の完了後は 4.2.1 節で説明した 2 回目以降の方法でクラスタリングを行い、ペアデータの作成と教師データへの追加を行う。この操作をクラスタリング前に全てのクラスタの中に存在するデータ数が一つになるまで反復して行う。図 4.2 にクラスタリングを n 回繰り返した際の教師データの様子を示す。図のような順序で反復処理の各回によって作成されたペアデータが教師データに追加される。これによって徐々にペアデータの中に存在する二つのデータの特徴量同士の距離が縮まり、段階的な学習難易度上昇を実現できると考えた。

教師データ
1回目のクラスタリング から作成した教師データ
2回目のクラスタリング から作成した教師データ
⋮
n回目のクラスタリング から作成した教師データ

図 4.2 教師データ追加の様子

第 5 章 評価実験

本稿ではランキング学習の 1 手法であるペアワイズ手法において，学習難易度が徐々に高まるようにペアデータを学習していくことが学習効率の向上に有効であることを示すために評価実験を行った．評価実験では提案手法によって教師データを作成した場合と無作為な順序で教師データを作成した場合でそれぞれ順序推定モデルの学習を行い，評価指標の値を比較した．学習に使用したデータ数と精度の変化を確認するため，使用するデータ数を 1 割から 10 割まで 1 割ずつ推移させて評価実験を行った．精度に有意な差があるかを統計的検定によって示した．

5.1 実験手順

本稿ではランキング学習に関するベンチマークデータセットである LETOR4.0[11] を用いて評価実験を行った．LETOR4.0 はマイクロソフトから提供されているデータセットである．中でもリストワイズ学習用のデータセットである MQ2008-list を用いて実験を行った．MQ2008-list 中にはクエリと文書のペアがある．各行はクエリと文書の関連度，クエリ番号，46 次元の特徴量を持つ．また，MQ2008-list では同一文書でもクエリ毎に関連度が異なる．MQ2008-list には 5 交差検証用に Fold1 から Fold5 までのデータセットがあり，それぞれ訓練用・検証用・テスト用データが用意されている．我々はそれら 5 パターンのデータセットを実験に使用し，5 回検証を行った．使用するクエリ数は表 5.1 のような内訳で収集した．実験には各クエリからクエリとの関連度が上位 16 位までのデータを抽出して使用した．上位データに対する順序予測精度を高めるというランキング学習の目的達成のため，上位データを中心に学習をしたいこととバッチ学習時に設定したバッチサイズから上位 16 データの抽出を決定した．

実装はニューラルネットワークを用いて行った．1 層目に入力 46 次元，出力 10 次元の全結合層を設定し，活性化関数にシグモイド関数を設定した．2 層目に入力 10 次元，出力 1 次元の全結合層を設定した．学習速度の違いによる 2 手法の精度変化を見るため，学習率を 10^{-4} と 10^{-5} の二つの設定で実験を行った．損失関数には (4.1.1) 式を用いて学習を行った．学習は上限を 200 エポックとし，Early

表 5.1 Fold 毎に収集したクエリ数の内訳

	Fold1	Fold2	Fold3	Fold4	Fold5
訓練用	471	471	470	470	470
検証用	157	156	157	157	157
テスト用	156	157	157	157	157

Stopping は 20 エポックのうち検証用データに対する損失関数の値が減少しなかった場合とした。学習に提案手法を用いた場合と比較手法を用いた場合でそれぞれテストデータについて順序推定を行い、推定精度を比較した。

実験は 5 回検証によって行った。統計的検定により、比較する 2 手法の間で順序推定精度の平均値に有意な差があるかを確認する。帰無仮説は「比較する 2 手法の間で順序推定精度の平均値の変化に有意な差はない。」とし、対立仮説は「比較する 2 手法の間で順序推定精度の平均値の変化に有意な差がある。」とする。検定手法は対応のない 2 標本 t 検定を採用し、有意水準は 5% で両側検定を行った。 p 値が 0.05 を下回った場合に帰無仮説を棄却することができ、対立仮説が成立すると言えることができる。

5.2 提案手法の有効性確認

5.2.1 実験内容

本実験では提案手法の有効性を確認することを目的とする。提案手法を用いて作成した教師データと無作為な順序で作成した教師データを用いたペアワイズ学習を行い、順序推定精度を比較する。学習効率向上の有無を確認するため、学習に使用したデータ数が少ない区間についても精度を確認する必要があると考えた。本実験では全教師データ数に対して 1 割から 10 割の範囲で 1 割ずつ学習に使用するデータ数を推移させ、精度の変化を確認した。

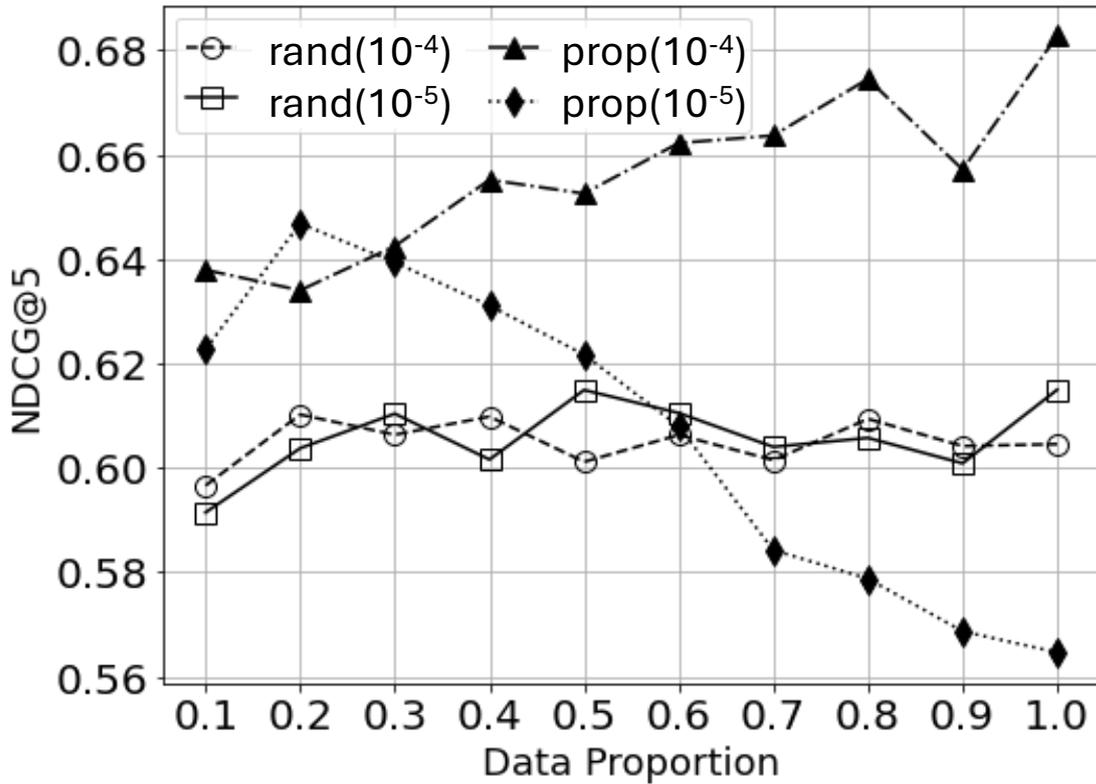


図 5.1 学習率が 10^{-4} , 10^{-5} の場合における NDCG@5 の推移

5.2.2 結果・考察

評価実験の結果を図 5.1, 表 5.2, 表 5.3 に示す. 図 5.1 の折れ線グラフでは, 学習に用いるデータの割合を全体に対して 1 割から 10 割まで推移させた場合の, 二つの手法によって作成されたモデルによるテストデータに対する順序推定結果についての NDCG@5 の変化を示している. 黒い三角が学習率を 10^{-4} とした場合の提案手法, 黒い菱形が学習率を 10^{-5} とした場合の提案手法, 白い丸が学習率を 10^{-4} とした場合のランダム手法, 白い四角が学習率を 10^{-5} とした場合のランダム手法を表している. 表 5.2, 表 5.3 では, 学習に用いるデータ割合を全体に対して 1 割から 10 割まで推移させた場合の, 二つの手法によって作成されたモデルによるテストデータに対する順序推定結果についての NDCG@5 の値と t 検定によって算出された p 値を示している. t 検定の結果から精度に有意差が確認できた場合を太字

表 5.2 学習率を 10^{-4} とした場合の NDCG@5 の推移と p 値

	データ割合				
	1	2	3	4	5
提案手法	0.63787	0.63397	0.64243	0.65513	0.65251
ランダム	0.59657	0.61017	0.60631	0.60979	0.60112
p 値	0.05427	0.01265	0.00259	0.00998	0.00004
	データ割合				
	6	7	8	9	10
提案手法	0.66218	0.66362	0.67446	0.65711	0.68272
ランダム	0.60627	0.60144	0.60935	0.60414	0.60450
p 値	0.00025	0.00049	2×10^{-6}	0.02577	3×10^{-7}

表 5.3 学習率を 10^{-5} とした場合の NDCG@5 の推移と p 値

	データ割合				
	1	2	3	4	5
提案手法	0.62272	0.64677	0.63933	0.63105	0.62165
ランダム	0.59143	0.60373	0.61034	0.60156	0.61488
p 値	0.01282	0.00614	0.03384	0.03168	0.74116
	データ割合				
	6	7	8	9	10
提案手法	0.60791	0.58424	0.57873	0.56866	0.56462
ランダム	0.61047	0.60398	0.60570	0.60084	0.61493
p 値	0.87903	0.31825	0.13221	0.13922	0.00642

で示した.

学習率を 10^{-4} とした場合は t 検定の結果, 教師データが 1 割以外の場合で p 値が 0.05 を下回り, 提案手法を用いて学習した順序推定モデルが, ランダム手法によって学習した順序推定モデルの精度を有意に上回った. この結果から, 教師データを推移させる過程で, ランダム手法では上位データの順序推定精度をあまり高め

られなかったことに対して、提案手法による学習が上位データの順序推定精度を高めることに有効であったといえる。提案手法では各クラスタリング回数ごとに、クエリの中に存在するデータから学習難易度の低いペアデータを作成した後に、他のクエリを対象にして同様にペアデータを作成することを繰り返す。このデータ作成手法により、複数クエリから少しずつ順序情報を学習することができるため、複数クエリに共通した順序基準を学習することができ、モデルの汎化性能向上に繋がったと考える。

学習率を 10^{-5} とした場合は t 検定の結果、教師データが 1 割から 4 割の場合で p 値が 0.05 を下回り、提案手法を用いて学習した順序推定モデルが、ランダム手法によって学習した順序推定モデルの精度を有意に上回った。一方で、教師データが 5 割から 9 割の場合では提案手法を用いて学習した順序推定モデルが、ランダム手法によって学習した順序推定モデルの精度を有意な差が無いことが分かった。また、教師データが 10 割の場合には提案手法を用いて学習した順序推定モデルが、ランダム手法によって学習した順序推定モデルの精度を有意に下回った。教師データが 5 割以上の場合にランダム手法の精度を有意に上回ることができなかった要因として、学習率が低いために学習難易度が低いデータについての順序情報を十分に学習することができない状態で学習難易度が高いデータを学習してしまい、学習効率の向上を実現できなかったと考える。一方で、ランダム手法を用いた学習では異なる難易度を均一に学習することができるため、安定した順序推定精度を出せていると考える。また、提案手法では学習に用いる全てのクエリの中から各クラスタリング回数ごとに一部のデータを取り出してペアデータを作成して学習に用いる。この学習過程において、学習率が低いことで全体的な学習が浅くなり、複数クエリに共通した順序情報も、特定の順序傾向を持つクエリについても十分に学習を行うことができなかったと考える。

これらの考察から、難易度の低い順序でデータを学習させることのみではなく、各難易度の段階でモデルが十分に学習をすることで、提案手法による学習効率向上の効果をより期待することができると思う。モデルの学習状況の確認のため、学習の進行度合いを確認する方法や、進行度合いに合わせた学習率の動的な調整を行うことで、より学習効率を高められると考える。

第6章 おわりに

我々は検索エンジン利用者の満足度を向上させるため、検索結果の表示順序を利用者にとって適したものになりたいと考えた。そこで我々は機械学習を用いて順序推定を行うことを考えた。本研究では、順序推定を行うための機械学習手法であるランキング学習 [1] を用いた。本研究ではランキング学習の主な3手法のうち、ペアワイズ手法に着目した。なぜなら、リストワイズ手法に必要な全データが順序づけられたデータセットを準備することが困難である一方で、ペアワイズ手法ではいくつかの順序付きのペアデータから学習することが可能であり、データセットの準備の容易さから実用性が高いためである。

ペアワイズ手法で問題となる点は、学習時間が膨大な点である。そこで本研究ではペアワイズ手法において教師データ数が同じ場合に、通常の学習手法と比較して順序推定精度を高め、学習効率を向上させることを目指した。我々はこの目標を実現するためにカリキュラム学習を用いて学習効率を高めることを考えた。

我々はペアワイズ学習が各ペアデータの中に存在する二つのデータの順序を予測するように学習を行っていることに着想を得て学習難易度の定義を試みた。ペアワイズ学習では教師データの中に存在するペアデータについて、ペアの中に存在する二つのデータの予測値と正解ラベルの差を損失計算に用いる。よって、ペアの中に存在する二つのデータの特徴量同士の距離が近い場合に予測値が近くなり、順序推定が難しく学習難易度が高くなると考えた。本稿では学習難易度を徐々に高めるため、ペアの中に存在する二つのデータの特徴量同士の距離が徐々に近くなるように教師データを並び替えることを考えた。

提案手法では最初に、候補データの中で特徴量が大きく離れた二つのデータからペアデータを作成するために、候補データの特徴量に基づいて二つのクラスターにクラスターリングし、異なるクラスター間でペアデータを作成する。次に、既に作成されたペアデータと比較して、より特徴量が近いデータ同士のペアを作成する。そのために、既存の各クラスターの中に存在するデータをさらに二つにクラスターリングし、直前のクラスターが同じだったデータ同士かつ、異なる新たなクラスターに分類されたデータ同士でペアを作成する。これらの操作を繰り返すことで、作成されるペアデータの中に存在する二つのデータが持つ特徴量同士の距離を徐々に近づけ、最終

的に類似したデータ同士のペアが作成される。このデータ作成手法によって、教師データの段階的な学習難易度上昇と学習効率の向上を期待できると考えた。

提案手法の有効性を確認するために評価実験を行った。評価実験では提案手法を用いて作成した教師データと無作為な順序で作成した教師データを使用してペアワイズ学習を行い、順序推定モデルを作成した。学習に用いたデータ数が少ない区間での精度変化を確認するため、学習に用いるデータ数は教師データ全体に対して1割から10割までの範囲を1割間隔で推移させて行った。作成したそれぞれのモデルを用いてテストデータについての順序推定を行い、NDCG@5を用いて推定結果の評価を行った。

評価実験の結果、学習率を 10^{-4} に設定した場合に、使用した教師データの割合が2割から10割の場合で、提案手法を用いて作成したモデルの精度が、無作為なデータ順序で作成したモデルの精度と比較して有意に上回った。この結果から、提案手法によって学習効率の向上を実現でき、カリキュラム学習における学習難易度の定義方法が適切だったといえる。

一方で、学習率を 10^{-5} に設定して同様の実験を行った場合には、使用した教師データの割合が5割から10割の場合で、提案手法を用いて作成したモデルの精度が、無作為なデータ順序で作成したモデルの精度と比較して有意な上回ることはなかった。この結果から、学習率を 10^{-5} に下げたために、低い学習難易度のデータについての学習が十分にされないまま高い学習難易度のデータを学習してしまい、モデルの推定精度を高められなかったと考える。これらの結果から、教師データの段階的な難易度上昇に加えて、各難易度段階においてモデルを十分に学習させることで、提案手法による学習効率向上効果をより期待できると考えた。

今後の展望として、学習難易度の段階的な上昇と能動学習手法を組み合わせることを行い、さらなる学習効率向上を目指したい。また、学習過程でのモデルの学習状況の確認や、それに応じた学習率の動的な変化を組み合わせることによって、カリキュラム学習の効果を高めたい。また、提案手法をリストワイズ学習に応用して学習効率向上を実現し、順序推定モデルの精度を高めたい。

謝辞

本研究を進めるにあたり、指導教員である鈴木准教授にはたくさんのご助言、ご指導をいただきました。留学先でも研究を進めやすいようにゼミの時間を調整していただいたり、ポジティブな言葉をかけていただいたことで安心して留学先でも各活動に励むことができました。本当にお世話になりました、感謝しております。事務補佐員の佐野さん、井尾さんには学外発表の手続きを始め、各方面でお世話になりました。そのおかげで研究をはじめとした各活動に集中することができ、学校生活を快適に過ごすことができました。ありがとうございました。友人や家族には日頃から相談に乗ってもらったり生活を支えていただきました。ありがとうございました。最後に鈴木研究室の皆さん、研究に行き詰まった時には自分のことのように一緒に考えてくださったり、知識の深いみなさんはとても頼りになり、安心して研究活動に取り組むことができました。ありがとうございました。皆様の助けなしには本論文作成に至ることはできませんでした。心より感謝申し上げます。

参考文献

- [1] Tie-Yan Liu, et al. Learning to rank for information retrieval. *Foundations and Trends® in Information Retrieval*, Vol. 3, No. 3, pp. 225–331, 2009.
- [2] Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *international conference on machine learning*, pp. 1311–1320. Pmlr, 2017.
- [3] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*, pp. 129–136, 2007.
- [4] David Cohn, Les Atlas, and Richard Ladner. Improving generalization with active learning. *Machine learning*, Vol. 15, pp. 201–221, 1994.
- [5] John A Hartigan, Manchek A Wong, et al. A k-means clustering algorithm. *Applied statistics*, Vol. 28, No. 1, pp. 100–108, 1979.
- [6] Junjie Liang, Jinlong Hu, Shoubin Dong, and Vasant Honavar. Top-n-rank: A scalable list-wise ranking method for recommender systems. In *2018 IEEE International Conference on Big Data (Big Data)*, pp. 1052–1058. IEEE, 2018.
- [7] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48, 2009.
- [8] Jeffrey L Elman. Learning and development in neural networks: The importance of starting small. *Cognition*, Vol. 48, No. 1, pp. 71–99, 1993.
- [9] Wojciech Zaremba and Ilya Sutskever. Learning to execute. *arXiv preprint arXiv:1410.4615*, 2014.
- [10] Aristidis Likas, Nikos Vlassis, and Jakob J Verbeek. The global k-means clustering algorithm. *Pattern recognition*, Vol. 36, No. 2, pp. 451–461, 2003.
- [11] Tao Qin and Tie-Yan Liu. Introducing LETOR 4.0 datasets. *arXiv preprint arXiv:1306.2597*, 2013.

発表リスト

- [1] エルゲン瑛夏, 鈴木優『ランキング学習における順序発生確率の不確実性に基づいた能動学習手法の提案』, 東海関西データベースワークショップ 2023, 2023
- [2] エルゲン瑛夏, 鈴木優『List-wise 手法における順位に基づいた多段階分割と学習データの上位遷移による能動学習』, 第 16 回データ工学と情報マネジメントに関するフォーラム (DEIM2024), 2024
- [3] エルゲン瑛夏, 鈴木優『ペアワイズ手法におけるクラスタリングを用いたペアデータの段階的難易度上昇による学習効率向上手法』, 第 17 回データ工学と情報マネジメントに関するフォーラム (DEIM2025), 2025