**11-442 / 11-642 / 11-742:**
**Search Engines**

# Feature-Based Retrieval

Jamie Callan and Chenyan Xiong
Carnegie Mellon University
callan@cs.cmu.edu

1

# Introduction

**We have covered many different approaches to retrieval**
- Retrieval models:  Vector space, BM25, language models, …
- Representations:  Title, body, url, inlink, …
- Query templates:  Sequential dependency models, …
- Query-independent evidence:  PageRank, url depth, …

**Different approaches can be combined to improve accuracy**
- We have seen combinations that were created manually

**Manually exploring the space of combinations is impractical**
- Too many combinations and too many parameters
- Use machine learning instead…

2

2

## Introduction

**A brief introduction to machine learning**
- Use <u>training data</u>    $X \to Y$          $X$: examples (feature vectors)
                                                                      $Y$: labels (desired values)
  to learn a <u>model</u> $h$ of how labels are assigned to examples
$$Y = h(X; \boldsymbol{w}) \quad \boldsymbol{w}: \text{ feature weights}$$
- For a <u>new (unlabeled)</u> example $x$, predict the label $y$
$$y = h(\boldsymbol{x}; \boldsymbol{w})$$

**Issues that must be addressed to apply this approach to search**
- How are examples encoded as feature vectors?
- What are the examples and labels?

3

**3**

---

## Introduction

**"Learning to rank (LeToR or LTR or L2R)"**
- Feature:  Evidence about how well query $q$ matches document $d$
- <u>Learn</u> a model that combines many types of evidence

**Example features**

| | |
|---|---|
| $f_1$:  $\text{BM25}_{\text{Title}}$ (q, d) | $f_7$:  $\text{VSM}_{\text{Title}}$ (q, d) |
| $f_2$:  $\text{BM25}_{\text{Body}}$ (q, d) | $f_8$:  $\text{VSM}_{\text{Body}}$ (q, d) |
| $f_3$:  $\text{Indri}_{\text{Title}}$ (q, d) | $f_9$:  PageRank (d) |
| $f_4$:  $\text{Indri}_{\text{Body}}$ (q, d) | $f_{10}$:  $\text{Length}_{\text{URL}}$ (d) |
| $f_5$:  $\text{Indri\_SDM}_{\text{Body}}$ (q, d) | $f_{11}$:  $\text{Length}_{\text{Terms}}$ (q) |
| $f_6$:  $\text{RankedBoolean}_{\text{Body}}$ (q, d) |     :      :      : |

4

**4**

2

## Introduction

**"Learning to rank (LeToR or LTR or L2R)"**

- **Feature:** Information about q, d, or the pair (q, d)
  - E.g., $BM25_{Body}$ (q, d) or ContainsCityName (q) or PageRank (d)
- **Feature vector:** A list of features
- **Label:** {Relevant, Not Relevant} or {0, 1, 2, 3, 4} or …
  - Sometimes the label is <u>known</u> (*training*)
  - Sometimes the label is <u>predicted</u> using the model (*testing*)
- **Model:** A method of <u>combining evidence</u> to predict a label

5

**5**

## Introduction

**LeToR is a type of supervised learning**

**Familiar supervised learning tasks**
- Classify an example into a discrete category
  - Find h: $X \rightarrow Y; Y \in \{cat, dog, mouse, ...\}$
  - E.g. Naïve Bayes, SVM
- Map each example to a numeric score
  - Find h: $X \rightarrow Y; Y \in R$
  - E.g. Linear regression, logistic regression

6

**6**

## Introduction

**LeToR is a type of supervised learning**

**Ranking task**
- Produce the best <u>ranking</u> of given documents
- Usually done by finding <u>ranking scores</u> for (q,d)
  - h: $\boldsymbol{x} \to y; y \in R$
    - » $\boldsymbol{x}$: A feature vector that describes how well q matches d
    - » $y$: The predicted label (a ranking score)
  - Note: We only care about the <u>ranking order,</u> not the scores
- Predict scores of $d_1, d_2, d_3, d_4, \ldots$ for query q
- Sort documents by their scores to produce a ranking

## Introduction

**We start with linear models, which are easy**

$$Similarity(q,d) = \begin{matrix} w_1 \times f_1(q,d) + \\ w_2 \times f_2(q,d) + \\ \vdots \\ w_n \times f_n(q,d) + \end{matrix} = \sum w_i \, f_i(q,d) = h(\boldsymbol{x}, \boldsymbol{w})$$

$f_i$: A feature measuring how well q matches d (e.g., $BM25_{Body}$)

$w_i$: A weight indicating the importance of feature $f_i$

h: The model

$\boldsymbol{x}$: A feature vector (i.e., $f_1, f_2, \ldots f_n$)

$\boldsymbol{w}$ A weight vector (i.e., $w_1, w_2, \ldots w_n$)

## Introduction

**If you care about <u>efficiency</u>, you may be worrying about costs**
- $Similarity(q, d) = \sum w_i \times f_i(q, d)$
    - Each $f_i$ is a function that may be expensive
        - » E.g., BM25, Indri, SDM, …
    - There may be many features
        - » E.g., 20-100

**Won't this be incredibly slow?**

**9**

## Introduction

**Where does learning to rank fit in a large search engine?**

**First, query classification to determine the type of query**
- Make decisions about how the query will be evaluated

**Then retrieval is done by a <u>sequence of retrieval models</u>**

**One example (there are many variations)**
- Exact-match Boolean:        <u>Form a set</u> of documents
- Best-match retrieval:        <u>Rank</u> the set, select the top $n_1$
    - E.g., Indri, BM25
- Learned models:        <u>Re-rank</u> the top $n_1$, select the top $n_2$
    - **Note:** LeToR does not <u>find</u> documents, it just <u>scores</u> them
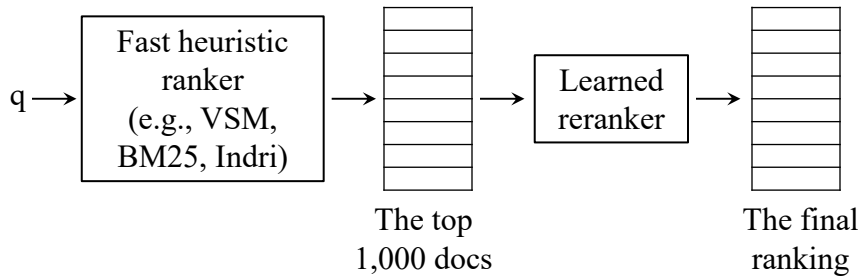
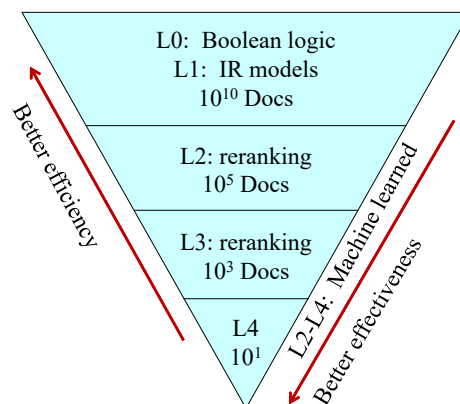**10**

## Introduction

**A standard reranking architecture**



This pipeline can have an arbitrary number of rerankers

**11**

## Introduction



**Each layer ranks, prunes, and returns results**

**Layered evaluation gives control over search costs**

- Simpler models are applied to massive data
  - Efficient
- Sophisticated models are applied to little data
  - Effective

(Pederson, 2010)

**12**

## Introduction

**LeToR methods are described along three main dimensions**
- The document representation (the <u>features</u>)
- The type (or style) of <u>training data</u>
- The machine learning algorithm

**Our focus is the first two (this is not a machine learning class) … but we will discuss some example learning algorithms, too**
- 11-641, *Machine Learning for Text Mining* provides more in-depth coverage of the learning algorithms

## Introduction

**For a query q, model document d as $\vec{x} = \Phi(d, q)$**
- $\Phi$ extracts features for (d, q).  Think of it as a feature factory.
- It produces a feature vector $\vec{x}$ or **x**

**Use training data to learn a model $h(\vec{x})$ or $h(\mathbf{x})$**
- $h(\vec{x})$ generates a real-valued score that is used to rank *d* for *q*

**Note:** This notation makes *q* and *d* <u>implicit</u> in the feature values
- It is consistent with a typical machine learning presentation
- Think only of feature vectors $\vec{x}$ and learned models $h(\vec{x})$
  - Each feature vector $\vec{x}$ describes a (d, q) pair

**Introduction:**
**Example Features**

**An initial set of features**

- $x_{CoordinationMatch}$    Number of query terms that match document d
- $x_{VSM}$    Vector space score for query q and document d
  (or a field in document d)
- $x_{BM25}$    BM25 score for query q and document d
  (BM25 features have been in Bing since ~2006)
- $x_{Indri}$    Indri score for query q and document d
- $x_{PageRank}$    PageRank score for document d
- $x_{Spam}$    Spam score for document d
- $x_{UrlDepth}$    Depth of document d in a website
- …

**15**

---

**Introduction:**
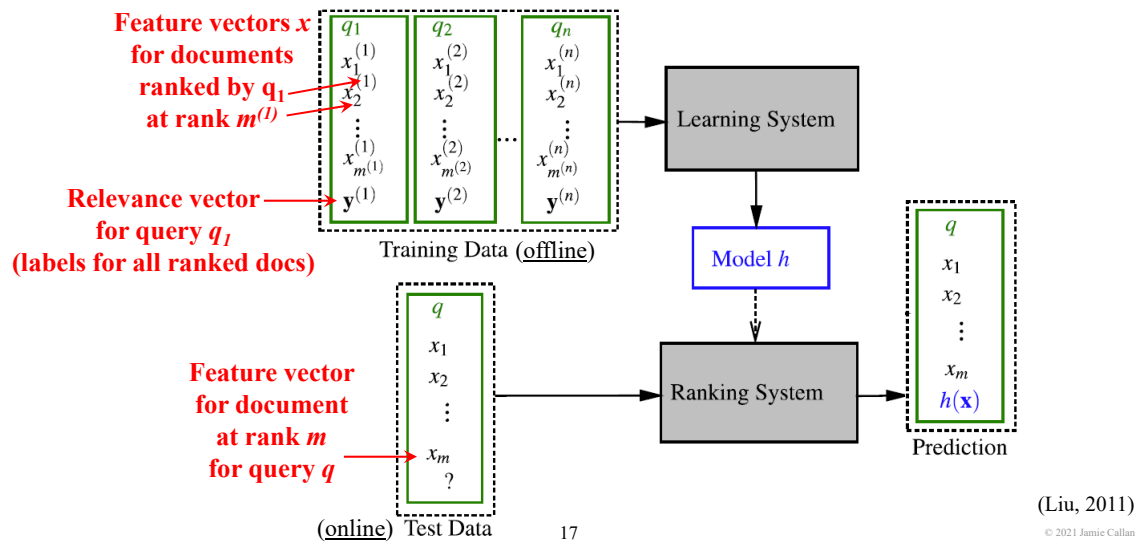**Example Features**

**More features…**

- Content length
- Query term density
- Number of inlinks
- Number of referring domains
- Number of referring IPs
- Time on site
- Pages per session
- Bounce rate
- Website security (https)

**16**

# Introduction:
# Learning to Rank Framework

**Feature vectors $x$ for documents ranked by $q_1$ at rank $m^{(1)}$**

| $q_1$ | $q_2$ | $q_n$ |
|---|---|---|
| $x_1^{(1)}$ | $x_1^{(2)}$ | $x_1^{(n)}$ |
| $x_2^{(1)}$ | $x_2^{(2)}$ | $x_2^{(n)}$ |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $x_{m^{(1)}}^{(1)}$ | $x_{m^{(2)}}^{(2)}$ | $x_{m^{(n)}}^{(n)}$ |
| $\mathbf{y}^{(1)}$ | $\mathbf{y}^{(2)}$ | $\mathbf{y}^{(n)}$ |

**Relevance vector for query $q_1$ (labels for all ranked docs)**

Training Data (<u>offline</u>)

Learning System

Model $h$

**Feature vector for document at rank $m$ for query $q$**

| $q$ |
|---|
| $x_1$ |
| $x_2$ |
| $\vdots$ |
| $x_m$ |
| ? |

Ranking System

| $q$ |
|---|
| $x_1$ |
| $x_2$ |
| $\vdots$ |
| $x_m$ |
| $h(\mathbf{x})$ |

Prediction

(<u>online</u>) Test Data

(Liu, 2011)

17

---

# Introduction

**Machine learning algorithms are designed to <u>predict labels</u>**

**What are the labels for (q, d)?**
- Relevant and not relevant?        Easy to get, hard prediction problem
- Relevance grades (e.g., 0, 1, 2, 3, 4)?   Easy to get, hard prediction problem
- A score in the range [0.0 … 1.0]?      Hard to get, easier prediction problem
- Position in an ideal ranking?        Hard to get, easier prediction problem
- …?

**Methods differ in how they address this issue**

18

*9*

## Introduction

**Today we cover three approaches to learning to rank**

**How are they similar?**
- Each uses a trained model $h$ to estimate the score of **x**
  - **x** describes how well document $d$ matches query $q$
- Documents are ranked by the score $h(\mathbf{x})$

**How are they different?**
- Three different approaches to <u>training the model</u> $h$
  - Different types of training data

**19**

## Outline

Introduction to feature-based methods

**Three approaches to training learning algorithms**
- Pointwise
- Pairwise
- Listwise

Benchmark datasets

Sample results

**20**

## Pointwise Approaches

**Approach: Train a model using <u>individual documents</u>**
- Training data:      $\mathbf{x} \rightarrow$ label        (q, d, relevance)
- Learned model:    $h(\mathbf{x}) \rightarrow$ label     (q, d, ?)
- (Remember, $\mathbf{x} = \Phi(d_i, q)$ )

**The type of label depends on the learning algorithm**
1. Classification (e.g., SVM):       Predict a relevance <u>category</u>
   - $h(\mathbf{x})$ = relevant
2. Regression (e.g., linear regression):   Predict a relevance <u>score</u>
   - $h(\mathbf{x})$ = 0.53

**21**

---

## Pointwise Approaches

**LeToR Procedure**
- **Offline:** Use Cranfield evaluation data to train a model
  - Queries, documents, relevance assessments
- **Online:** Given a new query q
  - Run a ranker (e.g., BM25) to get initial documents
    » E.g., the top 100 returned by BM25
  - For each initial document
    » Generate a feature vector $\mathbf{x_i} = \Phi(d_i, q)$
    » Use the model to generate a reranking label $h(\mathbf{x_i})$ for $d_i$
  - Rerank the initial documents

**22**

## Pointwise Approaches: Classification

**Classification (e.g., SVM):    Predict a relevance <u>category</u>**

$x_1 = \Phi(d_1, q), \quad y_1 = $ NRel   (not relevant)           $doc_1$

$x_2 = \Phi(d_2, q), \quad y_2 = $ HRel   (highly relevant)       $doc_2$

**Loss function:** Typically, 0 if the label is correct, otherwise 1

- The goal is to minimize the number of misclassified examples

**Training data is (relatively) easy to get**

- Use the Cranfield methodology
- Hire relevance assessors to assign documents to categories
  - E.g., { not relevant, relevant, highly relevant }

**23**

---

## Pointwise Approaches: Problems with Classification

**Categories are ordinal (ordered), which most algorithms ignore**

- Some errors are more serious than others
- Consider two trained models…

| $C_1$ | Label | Predicted | |
|---|---|---|---|
| $d_1$ | HRel | Rel | ✗ |
| $d_2$ | Rel | HRel | ✗ |
| $d_3$ | NRel | NRel | ✓ |

**Loss:  2**

| $C_2$ | Label | Predicted | |
|---|---|---|---|
| $d_1$ | HRel | NRel | ✗ |
| $d_2$ | Rel | Rel | ✓ |
| $d_3$ | NRel | NRel | ✓ |

**Loss: 1**

  - $C_1$ is the better ranker, but $C_2$ has lower loss
- What is the relative importance of different types of errors?
  - There isn't good theory to guide us

**24**

Page 12

## Pointwise Approaches:
## Regression

**Regression (e.g., linear regression): Predict a relevance <u>score</u>**

$$x_1 = \Phi(d_1, q), \quad y_1 = -1 \qquad doc_1$$
$$x_2 = \Phi(d_2, q), \quad y_2 = 3 \qquad doc_2$$

**Loss function:** Squared error
- The goal is to avoid large errors

**Training data is difficult to get**
- What is the desired score of $d_{21}$ for q?
  (see next slide)

**25**

## Pointwise Approaches:
## Problems with Regression

**Predicting numeric labels is difficult because <u>numeric labels are arbitrary</u>**
- During <u>training</u>, different mappings from relevance categories to numeric values cause the learning algorithm to learn different models
  - { NRel, Rel, HRel } → { 0, 1, 2 }      <u>Equal</u> differences among categories
  - { NRel, Rel, HRel } → { 1, 4, 9 }      <u>Prefer</u> HRel correct
  - { NRel, Rel, HRel } → { 1, 16, 81 }   <u>Must</u> get HRel correct
- What is the relative importance of different types of errors?
  - What are the right numbers to assign for Nrel, Rel, and HRel?
  - There isn't good theory to guide us
  - Thus, pointwise training data is inherently noisy

**26**

## Pointwise Approaches: Problems

**Pointwise methods focus on making good predictions for <u>individual documents</u>**
- E.g., minimize the number of misclassified examples
- E.g., minimize the sizes of classification errors

**Search metrics care about <u>the order</u> in <u>a set</u> of documents**
- E.g., NDCG@k, MAP@k

**Pointwise methods tend to be less accurate than other methods**
- Noise is inherent in the training data
- The loss functions don't model the task well
- The model focuses on the wrong goal (predict scores, not ranks)

27

**27**

---

## Outline

Introduction to feature-based methods

**Three approaches to training learning algorithms**
- Pointwise
- **Pairwise**
- Listwise

Benchmark datasets

Sample results

28

**28**

**Pairwise Approaches**

**Approach:  Train a model using <u>pairs of documents</u>**
- Training data:        prefer $(x_1, x_2)$
- Learned model:        $h(x_1) > h(x_2)$

**What is the pair value?**
- Binary assessments $\{x_1 > x_2, x_1 < x_2\}$

**Loss function**
- 0 if a pair is ordered correctly, otherwise 1

**Minimize the number of misclassified document pairs**
- Focus on <u>preference</u>, not specific scores or labels

**29**

---

**Pairwise Approaches**

**How is training data produced?**

| Relevance Data | Training Data$_q$ |
|---|---|
| $d_1$, relevant | $x_1 = \Phi(d_1, q) > x_2 = \Phi(d_2, q)$ |
| $d_2$, not relevant | $x_2 < x_1$ |
| $d_3$, not relevant | $x_1 > x_3$ |
| $d_4$, relevant | $x_3 < x_1$ |
| $d_5$, not relevant | :    : |
| $d_6$, relevant | $x_4 > x_2$ |
| $d_7$, not relevant | $x_2 < x_4$ |
| :   :   : | :    : |

**Each preference is one example (one instance of training data)**

**30**

## Pairwise Approaches

**How is preference data used to train a linear classifier?**

- **Training data:** prefer $(d_1, d_2)$
- Extract feature vectors: $\quad \vec{x_1} = \Phi(d_1, q) \qquad \vec{x_2} = \Phi(d_2, q)$
- Assume a linear model: $\quad h(\vec{x_i}) = w^T \vec{x_i}$

$$\text{prefer } (d_i, d_j) \;\rightarrow\; h(\vec{x_i}) > h(\vec{x_j})$$
$$\rightarrow\; w^T \vec{x_i} > w^T \vec{x_j}$$
$$\rightarrow\; w^T \vec{x_i} - w^T \vec{x_j} > 0$$
$$\rightarrow\; w^T (\vec{x_i} - \vec{x_j}) > 0$$
$$\rightarrow\; h(\vec{x_i} - \vec{x_j}) > 0$$

|       | $x_i$ | $x_j$ | $x_i$-$x_j$ |
|-------|-------|-------|-------------|
| $f_1$ | 0.67  | 0.59  | 0.08        |
| $f_2$ | 0.01  | 0.02  | -0.01       |
| $f_3$ | 0.23  | 0.19  | 0.04        |
| $f_4$ | 0.15  | 0.13  | 0.02        |
| $f_5$ | 0.80  | 0.82  | -0.02       |
|       | : :   | : :   | : :         |

**feature vectors**

- Now it is a standard two-class learning problem: $\quad h(\vec{x_j} - \vec{x_i}) < 0 \qquad h(\vec{x_i} - \vec{x_j}) > 0$
  - Use your favorite algorithm (e.g., regression, SVM)

31

© 2021 Jamie Callan

31

---

## Example: Ranking SVM

**Goal:** Minimize the number of misclassified document pairs

$$\text{minimize:} \quad \tfrac{1}{2} \vec{w} \cdot \vec{w} \quad + \quad C \sum \xi_{i,j,k}$$

    **complexity**    **accuracy (on training data)**

subject to:

$$\forall i \; \forall j \; \forall k: \; \xi_{i,j,k} \geq 0$$
$$\forall (d_i, d_j) \in r_1 : \vec{w} \cdot \vec{x_i} > \vec{w} \cdot \vec{x_j} + 1 - \xi_{i,j,k}$$
$$\dots$$
$$\forall (d_i, d_j) \in r_n : \vec{w} \cdot \vec{x_i} > \vec{w} \cdot \vec{x_j} + 1 - \xi_{i,j,k}$$

↑

**Preference pairs as constraints**

*C:* **Controls model complexity (small *C*) and training error (large *C*)**
*ξ :* **Slack variable**
*r_n:* *k* **training pairs from the** *n*th **ranking (query)**

(Croft, et al., 2010)

32

© 2021 Jamie Callan

32

---

Page 16

*16*

## Example: Ranking SVM

**Properties of Ranking SVM**
- Generalizes well
- Kernels can be used to improve accuracy
    - But linear kernels often work well
- Inherits desirable properties of SVM
    - Many open source tools
    - Much research effort on optimization → fast training
    - Theoretical guarantees (in learning)

33

**33**

## Pairwise Approaches

**LeToR Procedure**
- **Offline:** Use Cranfield evaluation data to train a model
    - Queries, documents, relevance assessments
- **Online:** Given a new query q
    - Run a ranker (e.g., BM25) to get initial documents
    - For each initial document
        » Generate a feature vector $\mathbf{x_i} = \Phi(d_i, q)$
        » Use the model to generate a reranking label $h(\mathbf{x_i})$
    - Rerank the initial documents

34

**34**

## Pairwise Approaches

**There have been many pairwise approaches**
- RankNet (Burges, et al., 2005)
- Frank (Tsai, et al., 2007)
- RankBoost (Freund, et al., 2003)
- Ranking SVM (Herbrich, et al., 2000; Joachims, 2002)
- MHR
- IR-SVM
- …

(Liu, 2011)

© 2021 Jamie Callan

**35**

## Problems with Pairwise Approaches

**Queries with an <u>even balance</u> of relevant and non-relevant documents dominate the training data**
- Number of pairs = $|R| \times |NR|$
  - $q_1$: 5 relevant, 5 not relevant: 5 R × 5 NR = 25 pairs
  - $q_2$: 2 relevant, 8 not relevant: 2 R × 8 NR = 16 pairs
  - $q_3$: 9 relevant, 1 not relevant: 9 R × 1 NR = 9 pairs

**The pairwise approach is more sensitive to noisy labels**
- One noisy label generates many training instances

**Usually the ranking position is ignored**
- As with the pointwise approaches

© 2021 Jamie Callan

**36**

## Outline

Introduction to feature-based methods

**Three approaches to training learning algorithms**
- Pointwise
- Pairwise
- Listwise

**Benchmark datasets**

**Sample results**

© 2021 Jamie Callan

## Listwise Approaches

**Approach:  Train a model using <u>a list of documents</u>**
- Training data:  $x_1 > x_2 > \ldots > x_n$
- Learned model:  $h(x_1) > h(x_2) > \ldots$

**What is the value of a particular ranking?**
- Some metric over the ranking
  - E.g., NDCG@n, with n=1, 3, 5, 10, …
  - E.g., MAP@n

**The goal is to maximize the value of the metric**
- Directly align the model with the ranking target

© 2021 Jamie Callan

# Listwise Approaches

**How is training data produced?**

| **Binary Relevance** | **Multi-Valued Relevance** | |
|---|---|---|
| $d_1$, relevant | $d_4$, perfect | |
| $d_4$, relevant | $d_1$, excellent | **Order documents** |
| $d_6$, relevant | $d_6$, very good | **by their relevance** |
| $d_2$, not relevant | $d_2$, poor | **scores** |
| $d_3$, not relevant | $d_3$, poor | **(best to worst)** |
| $d_5$, not relevant | $d_5$, poor | |
| $d_7$, not relevant | $d_7$, poor | |
| : : : | : : : | |

**The training algorithm generates and scores rankings**

---

# Listwise Approaches:
# Challenges

**Directly optimizing some metrics is hard**
- Popular metrics (e.g., NDCG@n) are not continuous, nor convex
  - Very hard for optimization algorithms

**Two common strategies in listwise approaches:**
1. Find another metric that is intuitive and easy to optimize.
   - E.g. likelihood of 'best' rankings in training data
2. Directly optimize ranking evaluation metrics, with approximation

**Listwise Approach: ListMLE**
*List*wise *M*aximum *L*ikelihood *E*stimation

**General Idea**
- Construct the probability of a ranking p($\mathbf{x_1}$, $\mathbf{x_2}$, ... , $\mathbf{x_n}$)
- Find the model $h(x)$ that maximizes the probability (likelihood) of best rankings in the training data
- Use $h(x)$ in ranking

**The key step is constructing p($\mathbf{x_1}$, $\mathbf{x_2}$, ... , $\mathbf{x_n}$)**
- It is impossible to define $p(x_1, x_2, ..., x_n)$ directly
    - The sample space is all possible ranks: $n!$
- ListMLE defines a generative process with much smaller space
    - With the help of independence assumptions

(Xia, et al., 2008)

41

**41**

---

**Listwise Approach: ListMLE**
*List*wise *M*aximum *L*ikelihood *E*stimation

**ListMLE's <u>generative definition</u> of the ranking process**
- The ranking is assembled iteratively from candidate documents
- Given a set of unranked candidate documents $S_i = \{d_1, ..., d_n\}$
    - Pick the best candidate $d_i$ from $S_i$ to appear at rank i
    $$p(d_i|S_i;w) = \frac{\exp(w^T x_i)}{\sum_{x_j \in S_i} \exp(w^T x_j)}$$
    - Remove $d_i$ from the set of unranked documents: $S_{i+1} = S_i \setminus d_i$
- Repeat until all candidate documents are ranked
- **Assumption:** $p(d_i|S_i;w)$ is independent of other documents

**Thus, the learning task is to learn $\exp(w^T x)$**

(Xia, et al., 2008)

42

**42**

**Listwise Approaches:**
**ListMLE**

**The likelihood of ranking $r_k$**

$$l(r_k; w) = p(x_1, x_2, \ldots, x_n; w) = \prod_{i=1}^{n} p(d_i | S_i; w)$$

**Use Maximum Likelihood Estimation to find parameters**

$$w^* = \operatorname{argmax}_w l(r_k; w)$$

- Solved by typical gradient methods

**For new query, documents are ranked by:**

$$h(x) = \exp(w^T x)$$

**Note:** Different learning process, but familiar type of learned model

(Xia, et al., 2008)

43

© 2021 Jamie Callan

**43**

---

**Listwise Approaches:**
**ListMLE**

**LeToR Procedure**
- **Offline:** Use Cranfield evaluation data to train a model
  - Queries, documents, relevance assessments
- **Online:** Given a new query q
  - Run a ranker (e.g., BM25) to get initial documents
  - For each initial document
    » Generate a feature vector $\mathbf{x_i} = \Phi(d_i, q)$
    » Use the model to generate a reranking label $h(\mathbf{x_i})$
  - Rerank the initial documents

44

© 2021 Jamie Callan

**44**

Page 22

22

## Three Approaches to Training: Summary

**Pointwise**
- Training data is a document category or score
- Accurate category or score $\neq$ accurate ranking
- Position information ignored

**Pairwise**
- Training data is a preference among a pair of documents
- Accurate preference $\neq$ accurate ranking
- Position information ignored

**Listwise**
- Training data is a ranking of documents
- Hard to directly optimize ranking metrics

**45**

## Three Approaches to Training: Summary

**Pointwise is the weakest of the three approaches**

**Pairwise and listwise are about equally useful**
- Pairwise has an imperfect learning target, but is easy to achieve
  - Minimizes pairwise errors, but we want the best ranking
  - A simpler learning problem with theoretical guarantees
- Listwise has a perfect learning target, but is harder to achieve
  - Learning target is exactly the same as what we want
  - A harder learning problem

**Which do you prefer?**

**46**

# Outline

Introduction to feature-based methods

Three approaches to training learning algorithms

- Pointwise
- Pairwise
- Listwise

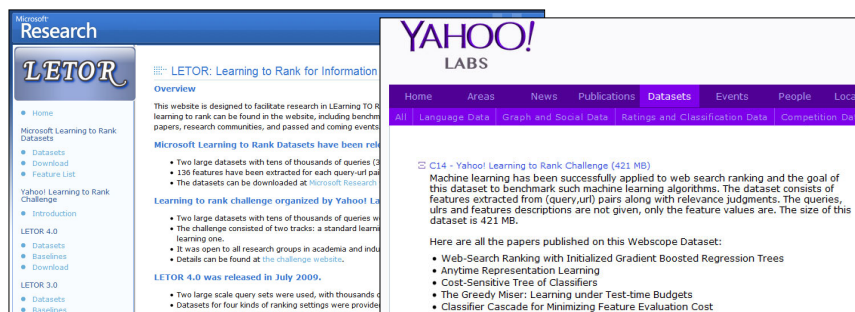**Benchmark datasets**

Sample results

**47**

---

# Benchmark Datasets

**Three popular LeToR benchmark datasets**

1. The LeToR collection (Microsoft) (2007)
2. Yahoo! Challenge datasets (2010)
3. Microsoft Learning to Rank datasets (2010)

**48**

# Benchmark Dataset #1:
# The LeToR Collections

**Created by Microsoft using existing publicly-available data**

- **The .gov corpus**
  - 1M .gov web documents from 2002
  - 350 queries (topic, homepage, named page)
  - Top 1000 documents per query returned by BM25
  - 64 features
- **The OHSUMED corpus**
  - 350K PubMed abstracts from 1987-1991
  - 106 queries (informational)
  - All judged documents
  - 40 features

(http://research.microsoft.com/en-us/projects/mslr/feature.aspx)

(Qin, et al., 2010)

49

49

---

# Benchmark Dataset #1:
# The LeToR Collections

**Field-specific features (title, anchor, url, whole) for ($d$, $q$)**

- Tf:  Sum over all query terms:  $\sum_{q_{i \in q}} tf_{q_i}$

- Idf:  Sum over all query terms:  $\sum_{q_{i \in q}} idf_{q_i}$

- Tf × Idf:  Sum over all query terms:  $\sum_{q_{i \in q}} tf_{q_i} \times idf_{q_i}$

- Field length

- Retrieval model scores
  - Boolean, VSM, BM25, $LM_{Abs}$, $LM_{Dirichlet}$, $LM_{JM}$
- Hyperlink-based features, HITS, Topic-specific PageRank, …

(Qin, et al., 2010)

50

50

## Benchmark Dataset #1:
## The LeToR Collection

**Document-level features for (*d, q*)**

- PageRank, number of inlinks
- Url: Number of '/', length
- Number of child pages

**Note that *q* is irrelevant to these features**

- These features prefer certain types of pages

(Qin, et al., 2010)

51

**51**


## Benchmark Dataset #2:
## The Yahoo Challenge! Datasets

**Created by Yahoo using proprietary data**

- **Set 1**
  - 710K feature vectors
  - 30K queries
  - 519 features (not described)
  - Relevance scale with 5 values
- **Set 2**
  - 173K feature vectors
  - 6K queries
  - 596 features (not described)
  - Relevance scale with 5 values

(Liu, 2011)

52

**52**

**Benchmark Dataset #3:**
**Microsoft Learning to Rank Dataset**

**Created by Microsoft using proprietary data**

- 3.7M web documents
- 30K queries
- 136 features
- Relevance scale with 5 values

**Example data**

```
0    qid:1    1:3 2:0 3:2 4:2 ... 135:0 136:0
2    qid:1    1:3 2:3 3:0 4:0 ... 135:0 136:0
```

relevance    query    features
              id

(http://research.microsoft.com/en-us/projects/mslr/)

53

53

---

**Benchmark Dataset #3:**
**Microsoft Learning to Rank Dataset**

**Field-specific features (title, anchor, body, url, whole) for ($d$, $q$)**

- Covered query term number, covered query term ratio
    - Terms in query q that appear in d
- Field length, field length normalized in various ways
- Idf
- Tf:  Min, Max, Sum, Mean, Variance
    - "apple pie recipe":  Min ($tf_{apple}$ $tf_{pie}$ $tf_{recipe}$)
- Tf × idf: Min, Max, Sum, Mean, Variance
- Retrieval model scores:  Boolean, VSM, BM25, $LM_{Dirichlet}$, $LM_{JM}$

(http://research.microsoft.com/en-us/projects/mslr/feature.aspx)

54

54

## Benchmark Dataset #3:
## Microsoft Learning to Rank Dataset

**Document-level features for (*d, q*)**

- Url: Number of '/', length
- Number of inlinks and outlinks
- PageRank, SiteRank, url click count, url dwell time
- Two quality scores
- Query-url click count

**Note that *q* is irrelevant to these features**

- These features prefer certain types of pages

(http://research.microsoft.com/en-us/projects/mslr/feature.aspx)

**55**

---

## Outline

**Introduction to feature-based methods**

**Three approaches to training learning algorithms**

- Pointwise
- Pairwise
- Listwise

**Benchmark datasets**

**Sample results**

**56**

# Sample Experimental Results

**LeToR Benchmark (.gov dataset, topic distillation (TD) queries)**

| | Algorithm | NDCG@ | | | P@ | | | MAP |
|---|---|---|---|---|---|---|---|---|
| | | 1 | 3 | 10 | 1 | 3 | 10 | |
| **Point** | Regression | 0.320 | 0.307 | 0.326 | 0.320 | 0.260 | 0.178 | 0.241 |
| **Pair** | RankSVM | 0.320 | **0.344** | 0.346 | 0.320 | **0.293** | 0.188 | 0.263 |
| | RankBoost | 0.280 | 0.325 | 0.312 | 0.280 | 0.280 | 0.170 | 0.227 |
| | FRank | 0.300 | 0.267 | 0.269 | 0.300 | 0.233 | 0.152 | 0.203 |
| **List** | ListNet | **0.400** | 0.337 | **0.348** | **0.400** | **0.293** | **0.200** | **0.275** |
| | AdaRank | 0.260 | 0.307 | 0.306 | 0.260 | 0.260 | 0.158 | 0.228 |
| | SVM$^{Map}$ | 0.320 | 0.320 | 0.328 | 0.320 | 0.253 | 0.170 | 0.245 |

**RankSVM is similar to ListNet except at rank 1**

(Liu, 2011)

---

# Sample Experimental Results

**LeToR Benchmark (.gov dataset, named page (NP) queries)**

| | Algorithm | NDCG@ | | | P@ | | | MAP |
|---|---|---|---|---|---|---|---|---|
| | | 1 | 3 | 10 | 1 | 3 | 10 | |
| **Point** | Regression | 0.447 | 0.614 | 0.665 | 0.447 | 0.220 | 0.081 | 0.564 |
| **Pair** | RankSVM | 0.580 | 0.765 | 0.800 | 0.580 | **0.271** | 0.092 | 0.696 |
| | RankBoost | **0.600** | 0.764 | **0.807** | **0.600** | 0.269 | **0.094** | **0.707** |
| | FRank | 0.540 | 0.726 | 0.776 | 0.540 | 0.253 | 0.090 | 0.664 |
| **List** | ListNet | 0.567 | 0.758 | 0.801 | 0.567 | 0.267 | 0.092 | 0.690 |
| | AdaRank | 0.580 | 0.729 | 0.764 | 0.580 | 0.251 | 0.086 | 0.678 |
| | SVM$^{Map}$ | 0.560 | **0.767** | 0.798 | 0.560 | 0.269 | 0.089 | 0.687 |

**RankSVM, and RankBoost are best**

(Liu, 2011)

## Lessons Learned

**Observations about the effectiveness of different algorithms**
- Many learning algorithms perform relatively well
- Relative effectiveness: Listwise ≈ Pairwise > Pointwise
    - As expected

**Many ML algorithms work with pointwise & pairwise LeToR**
- Easy to develop, reasonably effective

**Listwise algorithms may be more effective**
- But, there are fewer off-the-shelf solutions
- Still an open research topic

## Lessons Learned

**Much of the LeToR literature uses lots of training data**
- Research is driven by web companies that have a lot of data
- But … <u>you</u> may not have a lot of data
    - Their conclusions may not apply to your situation

## Lessons Learned

**Observation from academic research**

- 100-200 labeled queries can support 50-100 features
  - Surprising compared to other classification/regression tasks, which need a higher ratio of examples to features
- The theory behind LeToR is still an open research topic

**Use large numbers of features cautiously**

- Start with a small set of high-quality features, then grow it
- Design features carefully

**61**

---

## Lessons Learned

**Much research is driven by machine learning researchers …their focus is on <u>machine learning</u> algorithms**

**The features used in most systems are surprisingly simple**

- Simple statistics (e.g., tf, idf, tf × idf, field length)
- Obvious variations of existing ranking algorithms
- A few page quality metrics

**Better understanding of <u>search</u> can produce better features      … and better search accuracy**

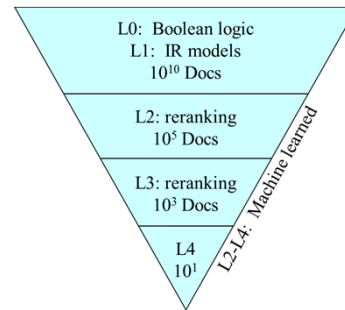- A nice opportunity for future improvement

(Liu, 2011)

**62**

## Outline

**Introduction to feature-based methods**

**Three approaches to training learning algorithms**

- Pointwise
- Pairwise
- Listwise

**Benchmark datasets**

**Sample results**



L0: Boolean logic
L1: IR models
$10^{10}$ Docs

L2: reranking
$10^5$ Docs

L3: reranking
$10^3$ Docs

L4
$10^1$

L2-L4: Machine learned

**63**

---

## For More Information

- B. Croft, D. Metzler, and T. Strohman. *Search Engines: Information Retrieval in Practice.* Addison Wesdley. 2010.
- T.-Y. Liu. *Learning to Rank for Information Retrieval.* Springer. 2011.
- T. Qin, T.-Y. Liu, J. Xu, and H. Li. LETOR: A Benchmark Collection for Research on Learning to Rank for Information Retrieval. *Information Retrieval Journal.* 2010.
- F. Xia, T.-Y. Liu, J. Wang, W. Zhang, and H. Li. Listwise approach to learning to rank – Theory and algorithm. In *Proceedings of the 25th International Conference on Machine Learning.* 2008.

**64**