# **Ranking and Learning**

293S UCSB, Tao Yang, 2023 Partially based on Manning, Raghavan, and Schütze's text book.

## **Table of Content**

- Weighted scoring for ranking
  - Ranking Features
- Learning to rank:
  - A simple example
  - Generalization
  - Type of learning-to-rank methods
- Learning to ranking as classification
  - Iearning-to-rank strategies
- Training process and evaluation

## **Aspects of Ranking Marching User Intent**

#### Relevance

- Documents need to be relevant to a user query.
- Authoritativeness.
  - High quality content is normally preferred since users rely on trustful information to learn or make a decision.
- Freshness.
  - Latest information is desired for time-sensitive queries.
- Preference
  - Personal or geographical preference can impact the choices

## **Weighted Scoring**

- Scoring with weighted features
  - Consider each document/query is a vector of features
  - Dot-product similarity of query and document vectors
- Example:
  - A simple weighted scoring method: use a linear combination of subscores:
    - E.g.,

Score = 0.6\*< <u>Title score></u> + 0.3\*<<u>Abstract score></u> + 0.1\*<<u>Body score</u>>

#### Example with binary subscores

Query term "ucsb admission" appears in title, and "ucsb" appears in body. Document score:  $(0.6 \cdot 2) + (0.1 \cdot 1) = 1.3$ .

## Simple Model of Ranking with Similarity [Croft, Metzler, Strohman's textbook slides]

#### Document features are topical or quality-based



## Simple Model of Ranking with Similarity [Croft, Metzler, Strohman's textbook slides]



### **Aspects of Ranking Marching User Intent**



## **Ranking Features used in Web Search**

- Modern systems especially on the Web use a great number of features:
  - Major web search engines use "hundreds" of such features – and they keep refinement
    - Text features: Query word frequency, Highlighted on page.
    - Document features: URL length, URL contains "~", Page length, Page freshness

- Categories of ranking signals
  - Query-dependent
  - Query-independent

## **Ranking Signals: Query Dependent**

- Text score
  - Document text.
    - -Text frequency: TFIDF, BM25
    - -Text proxmity:
      - Closeness of keywords that appear in a document
      - Sum of 1/distance<sup>2</sup>(w<sub>1</sub>,w<sub>2</sub>) for all keyword pairs
      - Query word span window
  - Anchor text
  - URL text
    - http://www.microsoft.com/en-us/download/

## **Ranking Signals: Query Dependent**

- Historical queries that yield document clicks
  - www.marriott.com for mariott, marriot
- Query classification and preference
  - Local, commerical products, news, image, video
  - Geo-location
- Link citation from documents that match the same query
  - # citations from documents relevant to a query
  - Hub authority analysis

## **Ranking Signals: Query independent**

- Document specific:
  - Link analysis: Page Rank
    - #incoming links to a URL
  - Quality of documents:
    - Spam analysis
  - Page classification and properties
    - Geo location
    - Country/language classification
    - Homepage/personal page classification
  - Freshness

#### Site specific

- Site quality:
  - Well-known sites
- Site classification: e.g. Country classification

## **Table of Content**

- Weighted scoring for ranking
  - Ranking Features
- Learning to rank:
  - A simple example
  - Generalization
  - Type of learning-to-rank methods
- Learning to ranking as classification
  - Iearning-to-rank strategies
  - Convert ranking as SVM based classification
- Training process and evaluation



## **Machine learning for ranking**

- How do we combine these signals into a good ranker?
  - How to derive weights if linear combination is used
  - What are other machine-learned models?
- Learning to rank
  - Learning from examples (called training data)



# Learning weights: Methodology

- •Given a set of training examples,
  - each contains (query q, document d, relevance score r).
  - r is relevance judgment for d on q
    - Simplest scheme
      - relevant (1) or nonrelevant (0)
    - More sophisticated: graded relevance judgments
      - 1 (bad), 2 (Fair), 3 (Good), 4 (Excellent), 5 (Perfect)

 Learn weights from these examples, so that the learned scores approximate the relevance judgments in the training examples

### Simple example of learning-to-rank

- Each doc has two zones: <u>Title</u> and <u>Body</u>
- For a chosen  $w \in [0,1]$ , score for doc d on query q

$$score(d,q) = w \cdot s_T(d,q) + (1-w)s_B(d,q)$$
  
where:

- $s_T(d, q) \in \{0, 1\}$  is a Boolean denoting whether q matches the <u>Title</u> and
- $s_B(d, q) \in \{0, 1\}$  is a Boolean denoting whether q matches the <u>Body</u>

## **Examples of Training Data**

Example	DocID	Query	$s_T$	$s_B$	Judgment
$\Phi_1$	37	linux	1	1	Relevant
$\Phi_2$	37	penguin	0	1	Non-relevant
$\Phi_3$	238	system	0	1	Relevant
$\Phi_4$	238	penguin	0	0	Non-relevant
$\Phi_5$	1741	kernel	1	1	Relevant
$\Phi_6$	2094	driver	0	1	Relevant
$\Phi_7$	3191	driver	1	0	Non-relevant

From these 7 examples, learn the best value of w.



- For each example Φ<sub>t</sub> we can compute the score based or score(d<sub>t</sub>, q<sub>t</sub>) = w ⋅ s<sub>T</sub>(d<sub>t</sub>, q<sub>t</sub>) + (1 − w)s<sub>B</sub>(d<sub>t</sub>, q<sub>t</sub>).
- We quantify Relevant as 1 and Non-relevant as 0
- Would like the choice of w to be such that the computed scores are as close to these 1/0 judgments as possible
  - Denote by  $r(d_t, q_t)$  the judgment for training instance  $\Phi_t$
- Then minimize total squared regression error  $\sum_{\Phi_t} (r(d_t, q_t) - score(d_t, q_t))^2$

## **Optimize the selection of weights**

- There are 4 kinds of training examples
- Thus only four possible values for score
  - And only 8 possible values for error (relevant vs irrelevant)
- Let n<sub>01r</sub> be the number of training examples for which *title score* 0, *body score* 1, judgment = *Relevant*.
- Similarly define  $n_{00r}$ ,  $n_{10r}$ ,  $n_{11r}$ ,  $n_{00i}$ ,  $n_{01i}$ ,  $n_{10i}$ ,  $n_{11i}$

Error:

$$[1 - (1 - \omega)]^2 n_{01r} + [0 - (1 - \omega)]^2 n_{01i}$$

Add up contributions from various cases to get total error

 $(n_{01r} + n_{10i})w^2 + (n_{10r} + n_{01i})(1 - w)^2 + n_{00r} + n_{11i}$ 

• Now differentiate with respect to *w* to get optimal value of *w* as:

$$\frac{n_{10r} + n_{01i}}{n_{10r} + n_{10i} + n_{01r} + n_{01i}}.$$

### **Learning-based Web Search**

- Given features  $x_1, x_2, ..., x_M$  for each document, learn a ranking function  $f(x_1, x_2, ..., x_m)$  that minimizes the loss function *L* under a query
- f\*=min L( f(x<sub>1</sub>,x<sub>2</sub>,...,x<sub>M</sub> ), GroundTruth)
- Some related issues
  - The functional space
    - linear/non-linear? continuous? Derivative?
  - The search strategy
  - The loss function

## **Table of Content**

- Weighted scoring for ranking
  - Ranking Features
- Learning to rank:
  - A simple example
  - Generalization
- Learning to ranking as classification
  - Convert ranking as SVM based classification
- Training process and evaluation

## Relationship to Classification Problem: An example

- Collect a training corpus of (q, d, r) triples
  - Relevance r is still binary for now
  - Document is represented by a feature vector
    - $\mathbf{x} = (\alpha, \omega)$   $\alpha$  is cosine similarity,  $\omega$  is minimum query window size
      - ω is the shortest text span that includes all query words (Query term proximity in the document)
- Train a machine learning model to predict the class r of a document-query pair

example	docID	query	cosine score	ω	judgment
$\Phi_1$	37	linux operating system	0.032	3	relevant
$\Phi_2$	37	penguin logo	0.02	4	nonrelevant
$\Phi_3$	238	operating system	0.043	2	relevant
$\Phi_4$	238	runtime environment	0.004	2	nonrelevant
$\Phi_5$	1741	kernel layer	0.022	3	relevant
$\Phi_6$	2094	device driver	0.03	2	relevant
$\Phi_7$	3191	device driver	0.027	5	nonrelevant

#### Window-based text span score

- Query text span in a document is the minimum length of word interval that covers all query words
- Example document:

Fred's tropical fish shop is the best place to find tropical fish at low price

- Span for query "tropical fish": 2
- Span for query "Fred's fish shop": 4

## **Using classification for deciding relevance**



## Summary: Ranking vs. Classification

- Classification
  - Well studied: Bayesian, Neural network, Decision tree, SVM, Boosting, ...
  - Training data: points: Positive: x1, x2, x3, Negative: x4, x5



#### Ranking

- Two ways to transform ranking problem to classification:
  - Assign a document to a class (relevant/nonrelevant)

Or assign to multiple classes such as perfect, excellent, good, fair, bad)

2. Classify the relationship of two documents in answering a query

## Classification Algorithms in Machine Learning

- Bayes
- Decision trees
- SVM (Supporting Vector Machines)
- Learning ensembles with many classifiers
  - Random Forest
  - Boosting regression trees
- Neural networks
  - Simple linear classifier (perceptron )
  - Deep multi-layer neural networks

- E.g. Convolution neural network

Training learns parameters involved in a network. A popular learning algorithm: SGD

Given a feature vector  $\mathbf{x}=(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ 

$$o(x_1,...,x_n) = \begin{cases} 1 \text{ if } w_0 + w_1 x_1 + ... + w_n x_n > 0\\ -1 \text{ otherwise} \end{cases}$$

Sometimes we will use simpler vector notation :



n+1 weight parameters are learned during training

UMN CS 5751 Machine Learning



## Strategies for "learning to rank"

- Point-wise learning
  - Given a query-document pair, predict a score (e.g. relevancy score)
    - Map f(x) to one of relevance vaules 0,1,2...
- Pair-wise learning
  - the input is a pair of results for a query, and the classification target is the relevance ordering relationship between them
  - Correct Order: f(x<sub>1</sub>) >f (x<sub>2</sub>) if x<sub>1</sub> is more relevant than x<sub>2</sub>
  - Otherwise incorrect.
- List-wise learning
  - Directly optimize the ranking metric (e.g. NDCG) for each query with a list of ranked results

#### **Point-wise learning: Example**

- Goal is to learn a threshold to separate each rank
- Assume 3 relevance levels: 1, 2, 3



## Pair-wise Learning

- A ranking should correctly classify the order of documents based on their relevance score:
  - Assume query q has matched documents ordered as x1, x2, x3, x4, x5

- Correct order

(x1, x2), (x1, x3), (x1, x4), (x1, x5), (x2, x3), (x2, x4) ...

- Other orders are incorrect

Convert ranking into binary classification

#### Learning to Rank – Example of Loss Functions

Given 3 documents for Query  $q: d_1, d_2, d_3$ , a ranker  $f_{\theta}$ , define loss *L*:

#### **Pointwise:**

 $L(f_{\theta}, q, d_1, d_2, d_3) = L(f_{\theta}, q, d_1) + L(f_{\theta}, q, d_2) + L(f_{\theta}, q, d_3)$ 

#### **Pairwise:**

 $L\left(f_{\theta}, q, d_{1}, d_{2}, d_{3}\right) = L\left(f_{\theta}, q, d_{1}, d_{2}\right) + L\left(f_{\theta}, q, d_{1}, d_{3}\right) + L\left(f_{\theta}, q, d_{2}, d_{3}\right)$ 

#### **Listwise:** $L(f_{\theta}, q, d_1, d_2, d_3) = L(f_{\theta}, q, d_1, d_2, d_3)$

Pretrained Transformers for Text Ranking: BERT and Beyond, 2021, Andrew Yates, Rodrigo Nogueira, and Jimmy Lin

#### **Types of Losses: Application Example**

# $\operatorname{Rel}(d_1) > \operatorname{Rel}(d_2) > \operatorname{Rel}(d_3), f_{\theta} = a \text{ neural network that outputs a probability}$



Pretrained Transformers for Text Ranking: BERT and Beyond, 2021, Andrew Yates, Rodrigo Nogueira, and Jimmy Lin

# Modified example for multi-class mapping with pair-wise learning

- Collect a training corpus of (q, d, r) triples
  - Relevance label r has 4 values
    - Perfect, Relevant, Weak, Nonrelevant
- Train a machine learning model to predict the class r of a document-query pair

example	docID	query	cosine score	ω	judgment
$\Phi_1$	37	linux operating system	0.032	3	Perfect
$\Phi_2$	37	penguin logo	0.02	4	Nonrelevant
$\Phi_3$	238	operating system	0.043	2	Relevant
$\Phi_4$	238	runtime environment	0.004	2	Weak
$\Phi_5$	1741	kernel layer	0.022	3	Relevant
$\Phi_6$	2094	device driver	0.03	2	Perfect
$\Phi_7$	3191	device driver	0.027	5	Nonrelevant

#### The Ranking SVM : Pairwise Learning [Herbrich et al. 1999, 2000; Joachims et al. KDD 2002]

- Aim is to classify training instance pairs as
  - correctly ranked
  - or incorrectly ranked
- This turns an ordinal regression problem back into a binary classification problem
- We want a ranking function *f* such that *c<sub>i</sub>* is ranked before *c<sub>k</sub>* :

 $c_i < c_k \text{ iff } f(\psi_i) > f(\psi_k)$ 

• Suppose that *f* is a linear function

 $f(\Psi_i) = \mathbf{w} \cdot \Psi_i$ 

• Thus

 $c_i < c_k \text{ iff } w(\psi_i - \psi_k) > 0$ 

## How many training examples formed for Ranking SVM?

example	docID	query	cosine score	ω	judgment
Φ <sub>1</sub>	37	linux operating system	0.032	3	 Perfect
$\Phi_2$	37	penguin logo	0.02	4	Nonrelevant
$\Phi_3$	238	operating system	0.043	2	Relevant
$\Phi_4$	238	runtime environment	0.004	2	Weak
$\Phi_5$	1741	kernel layer	0.022	3	Relevant
$\Phi_6$	2094	device driver	0.03	2	Perfect
$\Phi_7$	3191	device driver	0.027	5	Nonrelevant

#### 1 training case formed:

Query: device driver Order: Doc 2094, 3192

How to derive (a, b, c) based on the training examples?  $Score(d, q) = a\alpha + b\omega + c$   $Score(Doc 2094, q) \ge 1 + Score(Doc 3192, q)$  $0.03a + 2b + c \ge 1 + 0.027a + 5b + c$ 

## **Classification vs. Regression**

#### Classification

- Data in the form (x,y), where x is input vector. y is a category label
- Goal is to find indicator function estimation f.  $(0 \text{ if } v = f(\mathbf{x}))$

• Loss: L(y, ...)

$$y, f(\mathbf{x}) = \begin{cases} 0 & \text{if } y = f(\mathbf{x}) \\ 1 & \text{if } y \neq f(\mathbf{x}) \end{cases}$$

#### Regression

- Data in the form (x,y), where x is input vector. y is real-valued output
- Goal is to find function estimation f.
- Example loss:  $L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$



### **Classification vs. regression for learning to rank**

#### Regression

- Find relative rank scores. E.g. Score = af<sub>1</sub>+bf<sub>2</sub>, what is weight a and b?
- Not just classification labels.
- Classification isn't the best model for rank score learning:
  - Classification: Map to an unordered set of classes
  - Regression: Map to a real value
- This regression formulation gives extra power:
  - Relations between relevance levels are modeled
  - Fine grain scoring from highly relevant to irrelevant
  - Not an absolute scale of goodness

### **Popular Benchmarks for Relevance Evaluation**

#### ClueWeb 09 (Web pages) TREC Robust04 (News articles) MS MARCO Dev. Set ( Question/answers based on web pages) MS MARCO Passage/Document Ranking TREC Deep Learning 2019-2021 based on MS MARCO

Dataset	Domain	# Query	# Doc	Quer y Lengt h	Doc Lengt h	# judgement s per query	Graded relevance
ClueWeb09	Web	150	50M	1-5	857	90	yes
Robust04	News	250	0.5M	1-4	479	70	yes
MS MARCO passages - Dev	Q&A, Web	6980	8.8M	2-15	57	1	no
TREC DL 19		43				95	yes
TREC DL 20		54				66	yes
MSMARCO Documents - dev		5193	3.2M	2-15	1131	1	no

## **Evaluation Metrics for Relevance**

- Given a query and a ranked list of documents, how to measure relevance? Each document is marked as relevant or nonrelevant.
- Compute the mean values for the following metrics for all queries. 1 is the best and 0 is the worst.
  - MRR (mean reciprocal rank) at Position P

$$\mathbf{RR}(R,q) = \frac{1}{\mathrm{rank}_i}$$
Inverse of the rank of first relevant document up to position P

- Mean Precision@P: %relevant results up to Position P
- Mean Recall@P: %relevant results appeared up to Position
   P out of ALL known relevant results in the data collection.

#### **Relevance metric for multi-level judgement labels**

- Each doc is judged in multiple levels. E.g. bad, good, excellent, perfect.
- nDCG: the total discounted cumulative gain (DCG) up to position p scaled by the ideal DCG by perfect ranking

$$DCG_p = rel_1 + \sum_{i=2}^{p} \frac{rel_i}{\log_2 i}$$

or 
$$DCG_p = \sum_{i=1}^{p} \frac{2^{rel_i} - 1}{\log(1+i)}$$



Compute the mean values for all queries. 1 is the best and 0 is the worst

## Training and Evaluation: How to Evaluate Accuracy with Training Data



- The accuracy/error estimates on the training data: good or bad?
  - Not good
  - Because new data will probably not be **exactly** the same as the training data!
  - The algorithms do well on the training data may overfit, may not do well for future data

## **Evaluation with Independent Test Data**

 Estimation with independent test data is used when we have plenty of data and there is a natural way to forming training and test data.



• For example: reported experiments for which the classifiers were trained on data from 2017 and tested on data from 2018.

## Hold-out Method

• The hold-out method splits the data into training data and test data (usually 2/3 for train, 1/3 for test). Then we build a classifier using the train data and test it using the test data.



 The hold-out method is usually used when we have a sufficient large dataset for training and testing separately

## **Classification with two training/test datasets**

- A labeled dataset is divided into two sets
  - Training set is used to form a classifier that fits data
  - Test set is used to report classification errors with no bias
  - Test metric:
    - Binary classification. Accuracy is the percentage of cases that the derived classifier prdicts correctly.
- How to compute the error with more than 2 classes?
  - For example, 3 Classes: class 1, class 2, class 3.
  - Sqaured error sum
    - Sum (predicted class value target value)^2
    - Normalized by dividing the number of cases
  - Another way: Measure # of cases classified correctly for Class 1, and # of cases classifed correctly for Case 2 etc. Then compute average, or weighted average.

# Divide a dataset into 3 sets: Training set, validation set, and test set

- For more advanced setting, a labeled dataset is divided into 3 sets
  - Training set is used to form a tree under some parameters (e.g. when to stop tree growing)
  - Validation set is used to assess the accuracy of the derived classifier, and then readjust training parameters, and reassess again for the best validation performance
  - Test set is used to report accuracy/error of the final classifier with no bias

## **Classification: Train, Validation, Test Split**



The test data can't be used for parameter tuning!

## Making the Most of Available Data



- Difficult to obtain training/testing data
- Importance of more data
  - Generally, the larger the training data the better the classifier (but returns diminish).
  - The larger the test data the more accurate the error estimate.
  - *Can we use all data* to build the final classifier.

#### **k-Fold Cross-Validation**

- Select a subset for training and another subset for testing without overlapping.
  - data is split into k subsets of equal size; select one testing
- Repeat above process for k times
  - each subset in turn is used for testing and the remainder for training or training/validation
- The estimates are averaged to <u>Classifier</u> yield an overall estimate.





- Weighted scoring for ranking
  - Example: linear combination
  - Ranking features for web search
- Learning to rank: A simple example
  - Generalization to a general machine learning problem
- Learning to ranking as classification
  - Point-wise, pair-wise, & list-wise learning
  - Classification vs. regression for ranking
- How to train and evaluate