

PV211: Introduction to Information Retrieval

<https://www.fi.muni.cz/~sojka/PV211>

IIR 15-2: Learning to rank Handout version

Petr Sojka, Hinrich Schütze et al.

Faculty of Informatics, Masaryk University, Brno
Center for Information and Language Processing, University of Munich

2023-04-19

(compiled on 2023-04-12 10:39:41)

Overview

- 1 Zone scoring
- 2 Machine-learned scoring
- 3 Ranking SVMs

Take-away today

- Basic idea of learning to rank (LTR): We use machine learning to learn the relevance score (retrieval status value) of a document with respect to a query.
- Zone scoring: a particularly simple instance of LTR
- Machine-learned scoring as a general approach to ranking
- Ranking SVMs

Main idea

- The aim of term weights (e.g., tf-idf) is to measure term salience.
 - The sum of term weights is a measure of the relevance of a document to a query and the basis for ranking.
- Now we view this ranking problem as a machine learning problem – we will learn the weighting or, more generally, the ranking.
 - Term weights can be learned using training examples that have been judged.
- This methodology falls under a general class of approaches known as [machine learned relevance](#) or [learning to rank](#).

Learning weights

Main methodology

- Given a set of **training examples**, each of which is a tuple of: a query q , a document d , a relevance judgment for d on q
 - Simplest case: $R(d, q)$ is either relevant (1) or non-relevant (0)
 - More sophisticated cases: graded relevance judgments
- Learn weights from these examples, so that the learned scores approximate the relevance judgments in the training examples

Binary independence model (BIM)

- Is what BIM does a form of learning to rank?
- Recap BIM:
 - Estimate classifier of probability of relevance on training set
 - Apply to all documents
 - Rank documents according to probability of relevance

Learning to rank vs. Text classification

- Both are machine learning approaches
- Text classification, BIM and relevance feedback (if solved by text classification) are **query-specific**.
 - We need a query-specific training set to learn the ranker.
 - We need to learn a new ranker for each query.
- Learning to rank usually refers to **query-independent** ranking.
- We learn a single classifier.
- We can then rank documents for a query that we don't have any relevance judgments for.

Learning to rank: Exercise

- One approach to learning to rank is to represent each query-document pair as a data point, represented as a vector.
- We have two classes.
 - Class 1: the query is relevant to the document.
 - Class 2: the query is not relevant to the document.
- This is a standard classification problem, except that the data points are query-document pairs (as opposed to documents).
- Documents are ranked according to probability of relevance of corresponding document-query pairs.
- What features/dimensions would you use to represent a query-document pair?

Simple form of learning to rank: Zone scoring

- Given: a collection where documents have three **zones** (a.k.a. **fields**): `author`, `title`, `body`
- Weighted zone scoring requires a separate weight for each zone, e.g. g_1 , g_2 , g_3
- Not all zones are equally important:
e.g. `author` < `title` < `body`
→ $g_1 = 0.2$, $g_2 = 0.3$, $g_3 = 0.5$ (so that they add up to 1)
- Score for a zone = 1 if the query term occurs in that zone, 0 otherwise (**Boolean**)

Example

Query term appears in `title` and `body` only

Document score: $(0.3 \cdot 1) + (0.5 \cdot 1) = 0.8$.

General form of weighted zone scoring

Given query q and document d , weighted zone scoring assigns to the pair (q, d) a score in the interval $[0,1]$ by computing a **linear combination** of document zone scores, where each zone contributes a value.

- Consider a set of documents, which have l zones
- Let $g_1, \dots, g_l \in [0, 1]$, such that $\sum_{i=1}^l g_i = 1$
- For $1 \leq i \leq l$, let s_i be the Boolean score denoting a match (or non-match) between q and the i^{th} zone
 - $s_i = 1$ if a query term occurs in zone i , 0 otherwise

Weighted zone score a.k.a **ranked Boolean retrieval**

Rank documents according to $\sum_{i=1}^l g_i s_i$

Learning weights in weighted zone scoring

- Weighted zone scoring may be viewed as **learning a linear function** of the Boolean match scores contributed by the various zones.
- No free lunch: labor-intensive assembly of user-generated relevance judgments from which to learn the weights
 - Especially in a dynamic collection (such as the Web)
 - Major search engine put considerable resources into creating large training sets for learning to rank.
- Good news: once you have a large enough training set, the problem of learning the weights g_i reduces to a simple optimization problem.

Learning weights in weighted zone scoring: Simple case

- Let documents have two zones: title, body
- The weighted zone scoring formula we saw before:

$$\sum_{i=1}^I g_i s_i$$

- Given q, d , $s_T(d, q) = 1$ if a query term occurs in title, 0 otherwise; $s_B(d, q) = 1$ if a query term occurs in body, 0 otherwise
- We compute a score between 0 and 1 for each (d, q) pair using $s_T(d, q)$ and $s_B(d, q)$ by using a constant $g \in [0, 1]$:

$$\text{score}(d, q) = g \cdot s_T(d, q) + (1 - g) \cdot s_B(d, q)$$

Learning weights: determine g from training examples

Example

Φ_j	d_j	q_j	s_T	s_B	$r(d_j, q_j)$
Φ_1	37	linux	1	1	Relevant
Φ_2	37	penguin	0	1	Nonrelevant
Φ_3	238	system	0	1	Relevant
Φ_4	238	penguin	0	0	Nonrelevant
Φ_5	1741	kernel	1	1	Relevant
Φ_6	2094	driver	0	1	Relevant
Φ_7	3194	driver	1	0	Nonrelevant

- Training examples: triples of the form $\Phi_j = (d_j, q_j, r(d_j, q_j))$
- A given training document d_j and a given training query q_j are assessed by a human who decides $r(d_j, q_j)$ (either relevant or nonrelevant)

Learning weights: determine g from training examples

Example

Example	DocID	Query	s_T	s_B	Judgment
Φ_1	37	linux	1	1	Relevant
Φ_2	37	penguin	0	1	Nonrelevant
Φ_3	238	system	0	1	Relevant
Φ_4	238	penguin	0	0	Nonrelevant
Φ_5	1741	kernel	1	1	Relevant
Φ_6	2094	driver	0	1	Relevant
Φ_7	3194	driver	1	0	Nonrelevant

- For each training example Φ_j we have Boolean values $s_T(d_j, q_j)$ and $s_B(d_j, q_j)$ that we use to compute a score:

$$\text{score}(d_j, q_j) = g \cdot s_T(d_j, q_j) + (1 - g) \cdot s_B(d_j, q_j)$$

Learning weights

- We compare this score $score(d_j, q_j)$ with the human relevance judgment for the same document-query pair (d_j, q_j) .
- We define the error of the scoring function with weight g as

$$\epsilon(g, \Phi_j) = (r(d_j, q_j) - score(d_j, q_j))^2$$

- Then, the total error of a set of training examples is given by

$$\sum_j \epsilon(g, \Phi_j)$$

- The problem of learning the constant g from the given training examples then reduces to picking the value of g that minimizes the total error.

Minimizing the total error ϵ : Example (1)

Training examples

Example	DocID	Query	s_T	s_B	Judgment
Φ_1	37	linux	1	1	1 (relevant)
Φ_2	37	penguin	0	1	0 (nonrelevant)
Φ_3	238	system	0	1	1 (relevant)
Φ_4	238	penguin	0	0	0 (nonrelevant)
Φ_5	1741	kernel	1	1	1 (relevant)
Φ_6	2094	driver	0	1	1 (relevant)
Φ_7	3194	driver	1	0	0 (nonrelevant)

- Compute score:

$$score(d_j, q_j) = g \cdot s_T(d_j, q_j) + (1 - g) \cdot s_B(d_j, q_j)$$
- Compute total error: $\sum_j \epsilon(g, \Phi_j)$, where

$$\epsilon(g, \Phi_j) = (r(d_j, q_j) - score(d_j, q_j))^2$$
- Pick the value of g that minimizes the total error

Minimizing the total error ϵ : Example (2)

- Compute score $score(d_j, q_j)$

$$score(d_1, q_1) = g \cdot 1 + (1 - g) \cdot 1 = g + 1 - g = 1$$

$$score(d_2, q_2) = g \cdot 0 + (1 - g) \cdot 1 = 0 + 1 - g = 1 - g$$

$$score(d_3, q_3) = g \cdot 0 + (1 - g) \cdot 1 = 0 + 1 - g = 1 - g$$

$$score(d_4, q_4) = g \cdot 0 + (1 - g) \cdot 0 = 0 + 0 = 0$$

$$score(d_5, q_5) = g \cdot 1 + (1 - g) \cdot 1 = g + 1 - g = 1$$

$$score(d_6, q_6) = g \cdot 0 + (1 - g) \cdot 1 = 0 + 1 - g = 1 - g$$

$$score(d_7, q_7) = g \cdot 1 + (1 - g) \cdot 0 = g + 0 = g$$

- Compute total error $\sum_j \epsilon(g, \Phi_j)$

$$(1-1)^2 + (0-1+g)^2 + (1-1+g)^2 + (0-0)^2 + (1-1)^2 + (1-1+g)^2 + (0-g)^2 = 0 + (-1+g)^2 + g^2 + 0 + 0 + g^2 + g^2 = 1 - 2g + 4g^2$$

- Pick the value of g that minimizes the total error

Setting derivative to 0, gives you a minimum of $g = \frac{1}{4}$.

Weight g that minimizes error in the general case



$$g = \frac{n_{10r} + n_{01n}}{n_{10r} + n_{10n} + n_{01r} + n_{01n}}$$

- $n_{...}$ are the counts of rows of the training set table with the corresponding properties:

n_{10r} $s_T = 1$ $s_B = 0$ document relevant

n_{10n} $s_T = 1$ $s_B = 0$ document nonrelevant

n_{01r} $s_T = 0$ $s_B = 1$ document relevant

n_{01n} $s_T = 0$ $s_B = 1$ document nonrelevant

- Derivation: see book
- Note that we ignore documents that have 0 match scores for both zones or 1 match scores for both zones – the value of g does not change their final score.

Exercise: Compute g that minimizes the error

	DocID	Query	s_T	s_B	Judgment
Φ_1	37	linux	0	0	Relevant
Φ_2	37	penguin	1	1	Nonrelevant
Φ_3	238	system	1	0	Relevant
Φ_4	238	penguin	1	1	Nonrelevant
Φ_5	238	redmond	0	1	Nonrelevant
Φ_6	1741	kernel	0	0	Relevant
Φ_7	2094	driver	1	0	Relevant
Φ_8	3194	driver	0	1	Nonrelevant
Φ_9	3194	redmond	0	0	Nonrelevant

More general setup of machine learned scoring

- So far, we have considered a case where we combined match scores (Boolean indicators of relevance).
- Now consider more general factors that go beyond Boolean functions of query term presence in document zones.

Two examples of typical features

- The vector space cosine similarity between query and document (denoted α)
- The minimum window width within which the query terms lie (denoted ω)
 - Query term proximity is often indicative of topical relevance.
- Thus, we have one feature that captures overall query-document similarity and one features that captures proximity of query terms in the document.

Learning to rank setup for these two features

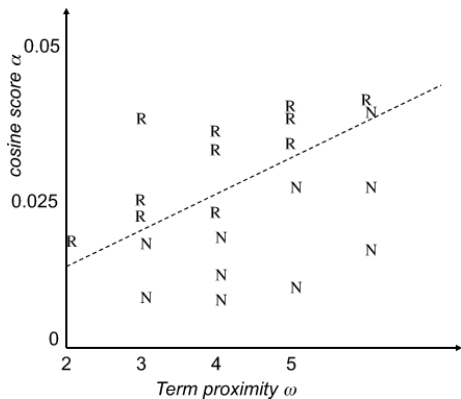
Example

Example	DocID	Query	α	ω	Judgment
Φ_1	37	linux	0.032	3	relevant
Φ_2	37	penguin	0.02	4	nonrelevant
Φ_3	238	operating system	0.043	2	relevant
Φ_4	238	runtime	0.004	2	nonrelevant
Φ_5	1741	kernel layer	0.022	3	relevant
Φ_6	2094	device driver	0.03	2	relevant
Φ_7	3191	device driver	0.027	5	nonrelevant

α is the cosine score. ω is the window width.

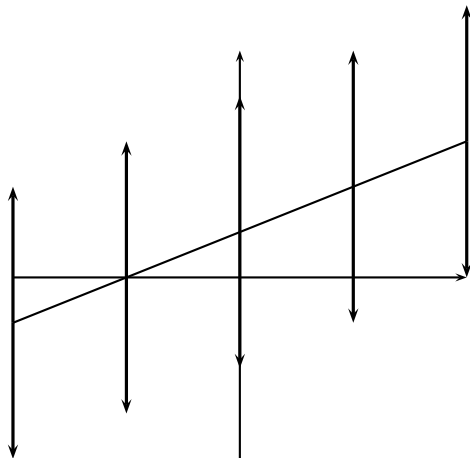
This is exactly the same setup as for zone scoring except we now have more complex features that capture whether a document is relevant to a query.

Graphic representation of the training set



This should look familiar.

In this case: LTR approach learns a linear classifier in 2D



- A linear classifier in 2D is a line described by the equation $w_1 d_1 + w_2 d_2 = \theta$
- Example for a 2D linear classifier
- Points $(d_1 \ d_2)$ with $w_1 d_1 + w_2 d_2 \geq \theta$ are in the class c .
- Points $(d_1 \ d_2)$ with $w_1 d_1 + w_2 d_2 < \theta$ are in the complement class \bar{c} .

Learning to rank setup for two features

- Again, two classes: relevant = 1 and nonrelevant = 0
- We now seek a scoring function that combines the values of the features to generate a value that is (close to) 0 or 1.
- We wish this function to be in agreement with our set of training examples as much as possible.
- A linear classifier is defined by an equation of the form:

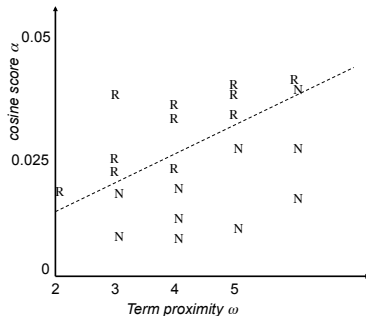
$$\text{Score}(d, q) = \text{Score}(\alpha, \omega) = a\alpha + b\omega + c,$$

where we learn the coefficients a, b, c from training data.

- Regression vs. classification
 - We have only covered binary classification so far.
 - We can also cast the problem as a regression problem.
 - This is what we did for zone scoring just now.

Linear classification in this case

- We pick a threshold θ .
- If $Score(\alpha, \omega) > \theta$, we declare the document **relevant**, otherwise we declare it **nonrelevant**.
- As before, all points that satisfy $Score(\alpha, \omega) = \theta$ form a line (dashed here) → linear classifier that separates relevant from nonrelevant instances.



Summary

- The problem of making a binary relevant/nonrelevant judgment is cast as a classification or regression problem, based on a training set of query-document pairs and associated relevance judgments.
- In the example: The classifier corresponds to a line $Score(\alpha, \omega) = \theta$ in the α - ω plane.
- In principle, any method learning a linear classifier (including least squares regression) can be used to find this line.
- Big advantage of learning to rank: we can avoid hand-tuning scoring functions and simply learn them from training data.
- Bottleneck of learning to rank: maintaining a representative set of training examples whose relevance assessments must be made by humans.

Learning to rank for more than two features

- The approach can be readily generalized to a large number of features.
- In addition to cosine similarity and query term window, there are lots of other indicators of relevance: PageRank-style measures, document age, zone contributions, document length, etc.
- If these measures can be calculated for a training document collection with relevance judgments, any number of such measures can be used to machine-learn a classifier.

LTR features used by Microsoft Research (1)

- Zones: body, anchor, title, url, whole document
- Features derived from standard IR models: query term number, query term ratio, length, idf, sum of term frequency, min of term frequency, max of term frequency, mean of term frequency, variance of term frequency, sum of length normalized term frequency, min of length normalized term frequency, max of length normalized term frequency, mean of length normalized term frequency, variance of length normalized term frequency, sum of tf-idf, min of tf-idf, max of tf-idf, mean of tf-idf, variance of tf-idf, boolean model, BM25

LTR features used by Microsoft Research (2)

- Language model features: LMIR.ABS, LMIR.DIR, LMIR.JM
- Web-specific features: number of slashes in url, length of url, inlink number, outlink number, PageRank, SiteRank
- Spam features: QualityScore
- Usage-based features: query-url click count, url click count, url dwell time
- See link in resources for more information

Shortcoming of our LTR approach so far

- Approaching IR ranking like we have done so far is not necessarily the right way to think about the problem.
- Statisticians normally first divide problems into **classification** problems (where a categorical variable is predicted) versus **regression** problems (where a real number is predicted).
- In between is the specialized field of **ordinal regression** where a ranking is predicted.
- Machine learning for ad hoc retrieval is most properly thought of as an ordinal regression problem, where the goal is to rank a set of documents for a query, given training data of the same sort.
- Next up: **ranking SVMs**, a machine learning method that learns an ordering directly.

Exercise

Example

Example	DocID	Query	Cosine	ω	Judgment
Φ_1	37	linux	0.03	3	relevant
Φ_2	37	penguin	0.04	5	nonrelevant
Φ_3	238	operating system	0.04	2	relevant
Φ_4	238	runtime	0.02	3	nonrelevant

Give parameters a, b, c of a line $a\alpha + b\omega + c$ that separates relevant from nonrelevant.

Basic setup for ranking SVMs

- As before we begin with a set of judged query-document pairs.
- But we do not represent them as query-document-judgment triples.
- Instead, we ask judges, for each training query q , to **order the documents** that were returned by the search engine with respect to relevance to the query.
- We again construct a vector of features $\psi_j = \psi(d_j, q)$ for each document-query pair – exactly as we did before.
- For two documents d_i and d_j , we then form the vector of feature differences:

$$\Phi(d_i, d_j, q) = \psi(d_i, q) - \psi(d_j, q)$$

Training a ranking SVM

- Vector of feature differences: $\Phi(d_i, d_j, q) = \psi(d_i, q) - \psi(d_j, q)$
- By hypothesis, one of d_i and d_j has been judged more relevant.
- Notation: We write $d_i \prec d_j$ for “ d_i precedes d_j in the results ordering”.
- If d_i is judged more relevant than d_j , then we will assign the vector $\Phi(d_i, d_j, q)$ the class $y_{ijq} = +1$; otherwise -1 .
- This gives us a training set of pairs of vectors and “precedence indicators”. Each of the vectors is computed as the difference of two document-query vectors.
- We can then train an SVM on this training set with the goal of obtaining a classifier that returns

$$\vec{w}^T \Phi(d_i, d_j, q) > 0 \quad \text{iff} \quad d_i \prec d_j$$

Advantages of Ranking SVMs vs. Classification/regression

- Documents can be evaluated **relative** to other candidate documents for the same query, rather than having to be mapped to a **global scale** of goodness.
- This often is an easier problem to solve since just a ranking is required rather than an absolute measure of relevance.
- Especially germane in web search, where the ranking at the very top of the results list is exceedingly important.

Why simple ranking SVMs don't work that well

- Ranking SVMs treat all ranking violations alike.
 - But some violations are minor problems, e.g., getting the order of two relevant documents wrong.
 - Other violations are big problems, e.g., ranking a nonrelevant document ahead of a relevant document.
- Some queries have many relevant documents, others few.
 - Depending on the training regime, too much emphasis may be put on queries with many relevant documents.
- In most IR settings, getting the order of the top documents right is key.
 - In the simple setting we have described, top and bottom ranks will not be treated differently.
- → Learning-to-rank frameworks actually used in IR are more complicated than what we have presented here.

Example for superior performance of LTR

SVM algorithm that directly optimizes MAP (as opposed to ranking).

Proposed by: Yue, Finley, Radlinski, Joachims, ACM SIGIR 2002.
Performance compared to state-of-the-art models: cosine, tf-idf, BM25, language models (Dirichlet and Jelinek-Mercer)

Model	TREC 9		TREC 10	
	MAP	W/L	MAP	W/L
SVM_{map}^{Δ}	0.242	–	0.236	–
Best Func.	0.204	39/11 **	0.181	37/13 **
2nd Best	0.199	38/12 **	0.174	43/7 **
3rd Best	0.188	34/16 **	0.174	38/12 **

Learning-to-rank clearly better than non-machine-learning approaches

Optimizing scaling/representation of features

- Both of the methods that we've seen treat the features as given and do not attempt to modify the basic representation of the document-query pairs.
- Much of traditional IR weighting involves **nonlinear** scaling of basic measurements (such as log-weighting of term frequency, or idf).
- At the present time, machine learning is very good at producing optimal weights for features in a linear combination, but it is not good at coming up with good nonlinear scalings of basic measurements.
- This area remains the domain of human feature engineering.

Assessment of learning to rank

- The idea of learning to rank is old.
 - Early work by Norbert Fuhr and William S. Cooper
- But it is only very recently that sufficient machine learning knowledge, training document collections, and computational power have come together to make this method practical and exciting.
- While skilled humans can do a very good job at defining ranking functions by hand, hand tuning is difficult, and it has to be done again for each new document collection and class of users.
- The more features are used in ranking, the more difficult it is to manually integrate them into one ranking function.
- Web search engines use a large number of features → web search engines need some form of learning to rank.

Exercise

Write down the training set from the last exercise as a training set for a ranking SVM.

Take-away today

- Basic idea of learning to rank (LTR): We use machine learning to learn the relevance score (retrieval status value) of a document with respect to a query.
- Zone scoring: a particularly simple instance of LTR
- Machine-learned scoring as a general approach to ranking
- Ranking SVMs

Resources

- Chapter 15-2 of IIR
- Resources at <https://www.fi.muni.cz/~sojka/PV211/> and <http://cislmu.org>, materials in MU IS and FI MU library
 - References to ranking SVM results
 - Microsoft learning to rank datasets