

Lecture 16: Supervised Methods, Neural Networks and Learning to Rank

IR Spring 22

Instructor: Prof. Somak Aditya

Slides Courtesy: Learning to Rank Tutorial

Categorization/Classification

Given:

- A representation of a document d
 - Issue: how to represent text documents.
 - Usually some type of high-dimensional space – bag of words, vector space
- A fixed set of classes: $C = \{c_1, c_2, \dots, c_J\}$

Determine:

- The category of d : $\gamma(d) \in C$, where $\gamma(d)$ is a **classification function**
- **We want to build classification functions** ("classifiers").

Classification Examples

- Classify an email as spam / legitimate
- Classify a news article as per its topic (Politics / Sports / Entertainment / ..)

Classification Methods

- Manual Classification

pol.
science
arts

- Rule-based Classification

$q = ?$

$d = \checkmark$

- Given: A document d , A fixed set of classes: $C = \{c_1, c_2, \dots, c_J\}$

- A training set D of documents each with a label in C
- Determine: A learning method or algorithm which will enable us to learn a classifier γ
- For a test document d , we assign it the class $\gamma(d) \in C$

$$\sigma(\theta^T x)$$

Classification Methods

- Supervised learning
 - Naive Bayes (simple, common)
 - k-Nearest Neighbors (simple, powerful)
 - Support-vector machines (generally more powerful)
 - Decision trees → random forests → gradient-boosted decision trees (e.g., xgboost)
 - Neural models
 - ...
- No free lunch: need hand-classified training data
- Many commercial systems still use some of them, but mostly NN/DL dominates.

How to Classify?

$$\sigma(\theta^T x) \rightarrow \theta^T x$$
$$(\theta_1^T \sigma(\theta_1^T x))$$

- Three aspects
 - Input Representation (features – handcrafted, now we learn them)
 - Function representation (hypothesis space)
 - Output classes
- Supervised learning classifiers can use any sort of feature
 - URL, email address, punctuation, capitalization, dictionaries, network features
- In the simplest bag of words view of documents
 - We use only word features

The bag of words representation

$Y(\text{I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.}) = C$

The bag of words representation

$$Y(\text{Table}) = C$$

great	2
love	2
recommen d	1
laugh	1
happy	1
...	...

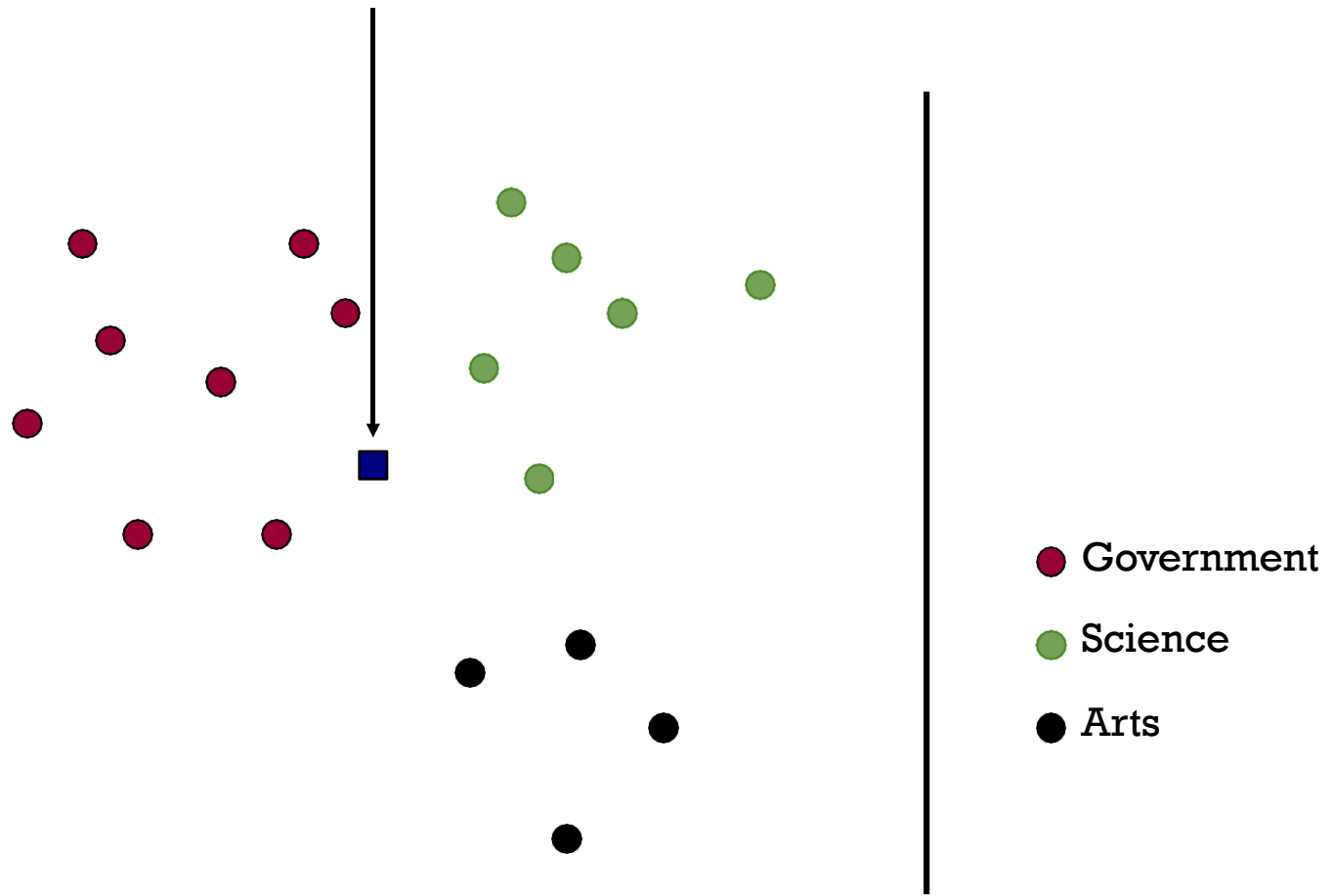
Remember: Vector Space Representation

- Each document is a vector, one component for each term (= word).
- Normalize vectors to unit length.
- High-dimensional vector space:
 - Terms are axes
 - 10,000+ dimensions, or even 100,000+
 - Docs are vectors in this space
- How can we do classification in this space?

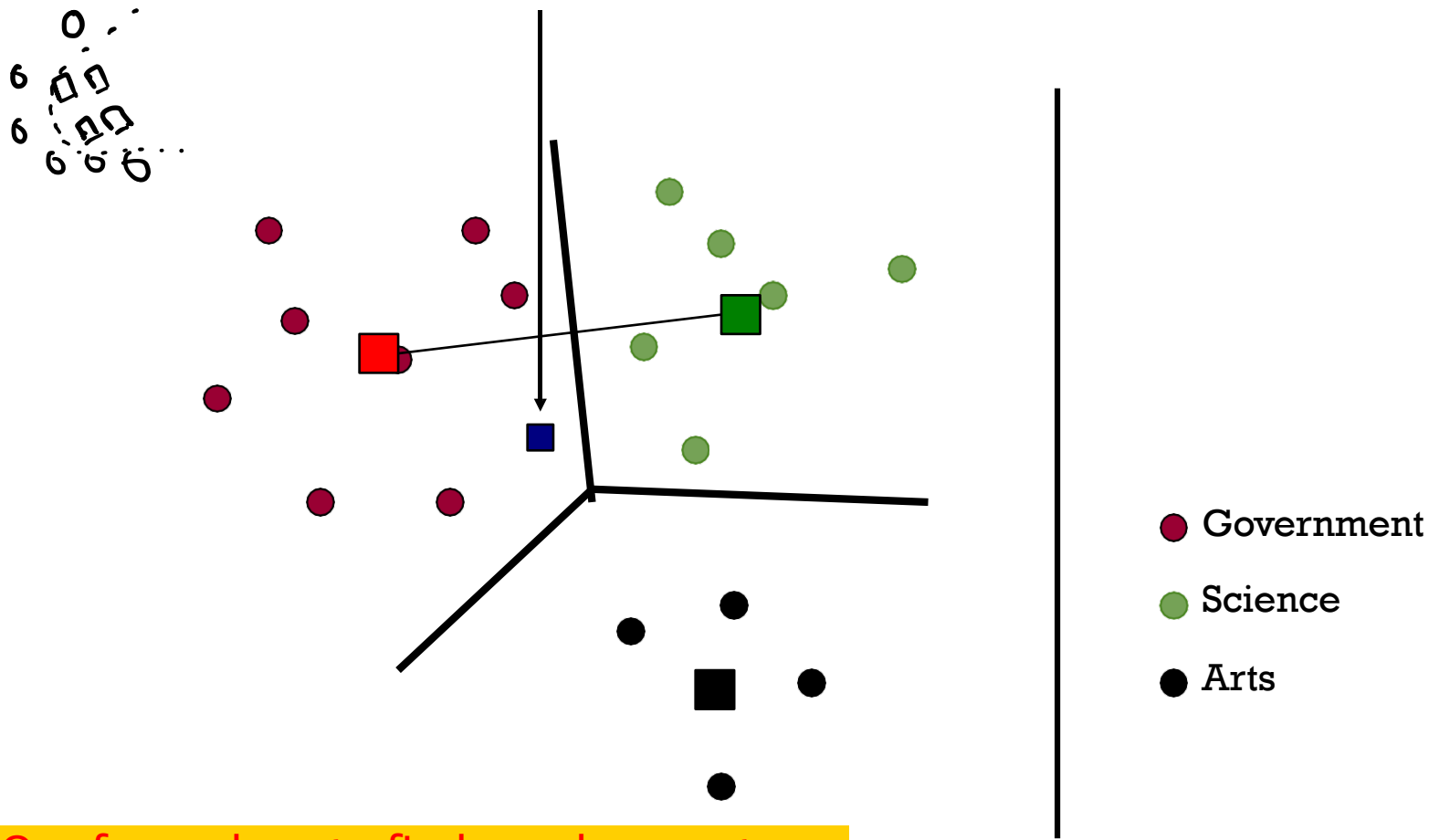
Classification Using Vector Spaces

- In vector space classification, training set corresponds to a labeled set of points (equivalently, vectors)
- **Premise 1:** Documents in the same class form a contiguous region of space
- **Premise 2:** Documents from different classes don't overlap (much)
- Learning a classifier: build surfaces to delineate classes in the space

Test Document of what class?



Test Document = Government



Our focus: how to find good separators

Another supervised learning problem: Regression

- Instead of mapping instances (e.g., documents) to discrete labels/classes ...
- Map instances to a **continuous (real) value**

$$x^2 + y^2 \neq 1 \Rightarrow$$

- E.g.,
 - Given the features (e.g., height, food habit, occupation, ~~weight~~) of many individuals, **predict the weight of a new person** for whom the other features are known
 - Given the bag of words representation of an email, predict the probability of it being spam

$$f(x) = y \quad y \in \mathbb{R} \quad 0/1 \quad 0.8$$

SUPERVISED METHODS FOR RETRIEVAL

Machine learning for IR ranking?

- We've looked at methods for ranking documents in IR
 - Cosine similarity, inverse document frequency, BM25, proximity, pivoted document length normalization, Pagerank, ...
- We've discussed supervised learning problems - classification and regression
- Can we use *machine learning* to rank the documents displayed in search results?
 - Sounds like a good idea
 - Known as "machine-learned relevance" or "learning to rank"
 - Actively researched and used by Web search engines

$$I: q, d, o: \frac{o}{1}$$

Simple example: Using classification for ad hoc IR

$$I, O, [\text{objective}]$$

$$\downarrow \vec{x}$$

Collect a training corpus of (q, d, r) triples

- Relevance r is here binary (but may be multiclass, with 3-7 values)
- Query-Document pair is represented by a feature vector
- Train a machine learning model to predict the class r of a document-query pair

$$\begin{array}{ccc}
 I \rightarrow \vec{x} & |V| & |V| \\
 q, d \quad ? & \text{vsm} & 2|V| \\
 & & \text{brut } \vee \text{ ceaser} \\
 & & \text{Doc : } \dots \dots \dots \\
 \underline{f(x)} = y & f(q, d_1) = 0.8 & \text{Label: } 1 \\
 & & f(q, d_2) = 0.8
 \end{array}$$

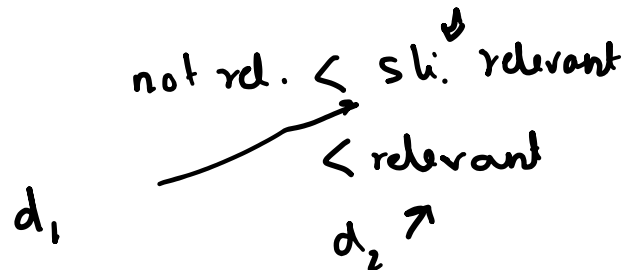
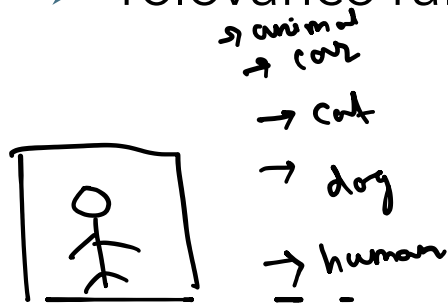
“Learning to rank”

- Classification probably isn't the right way to think about approaching ad hoc IR:
 - Classification problems: Map to an unordered set of classes
 - Regression problems: Map to a real value
 - Ordinal regression (or “ranking”) problems: Map to an *ordered* set of classes

“Learning to rank”

- Assume a number of categories **C** of relevance exist
 - These are totally ordered: $c_1 < c_2 < \dots < c_J$
 - This is the ordinal regression setup
- Assume training data is available consisting of document query pairs (d, q) represented as feature vectors x_i with
 - relevance ranking c_i

$$0, 1, 2$$
$$0 < 1 < 2$$



LEARNING TO RANK

Learning to rank (L2R)

Definition

- "... the task to automatically construct a ranking model using training data, such that the model can sort new objects according to their degrees of relevance, preference, or importance." - Liu [2009]

L2R models represent

- a rankable item e.g., a document, given
 - some context, e.g., a user-issued query
- as a numerical vector $\vec{x} \in \mathbb{R}^n$.

The ranking model $f: \vec{x} \rightarrow \mathbb{R}$ is **trained** to map the vector to a real-valued score such that relevant items are scored higher

Approaches

$$\begin{matrix} \textcircled{3} & \textcircled{1} & \textcircled{2} \\ f_{\theta}(x) & = & y \\ \theta = \arg \min_{\theta} J(\theta) & \rightarrow & \textcircled{4} \end{matrix}$$

$$q, d, y_{q,d} = 0.7$$

$$0/1$$

$$0/1/2$$

Based on training objectives [Liu 2009]:

- **Pointwise approach**: Relevance label $y_{q,d}$ is a number
 - Supervision: binary or graded human judgments or implicit user feedback (e.g., CTR).
 - Classification/regression to predict $y_{q,d}$, given $\vec{x}_{q,d}$.
- **Pairwise approach**: pairwise preference between documents for a query ($d_i \succ_q d_j$) as label.

q, d_1, d_2

1	0
0	1

 - Supervision: pairwise preference
 - Task: Given $\{q, d_i, d_j\}$, predict 1 (if d_i is preferred) or 0 otherwise.
- **Listwise approach**: optimize for rank-based metric, such as NDCG
 - difficult because these metrics are often not differentiable w.r.t. model parameters.

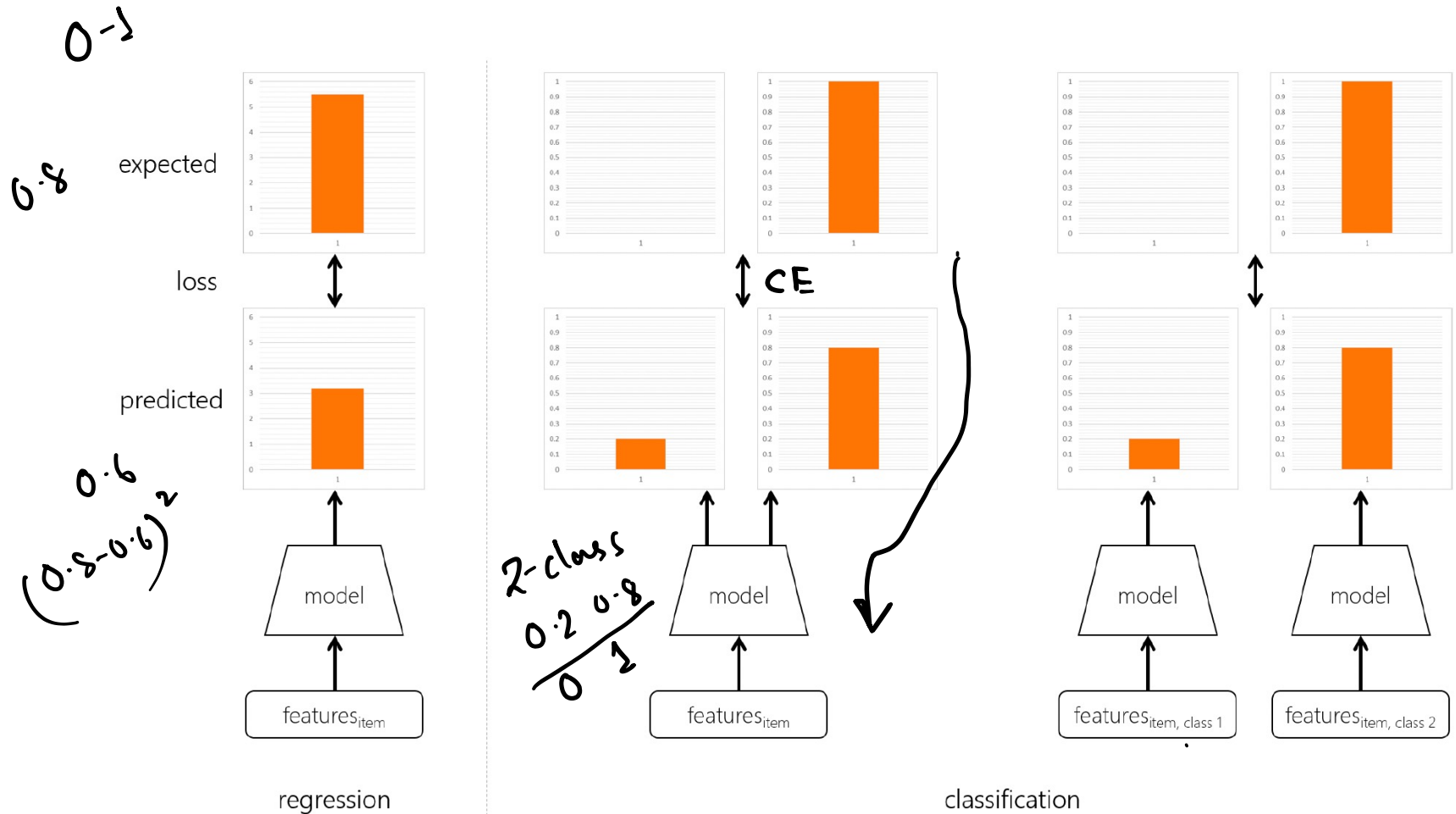
Features

Traditional L2R models employ hand-crafted features

They can often be categorized as:

- Query-independent or static features (e.g., incoming link count and document length)
- Query-dependent or dynamic features (e.g., BM25)
- Query-level features (e.g., query length)

Short Intro: Cross Entropy + NN

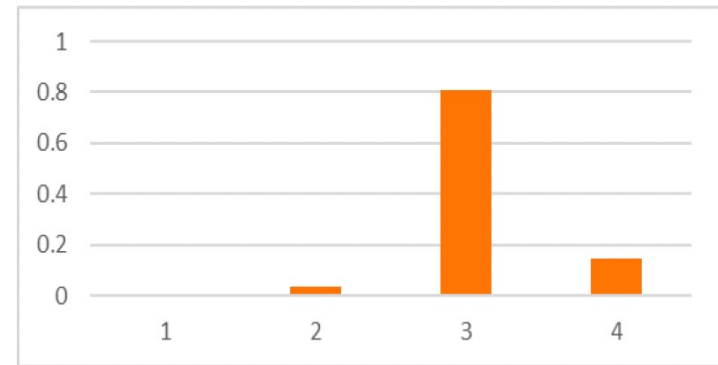
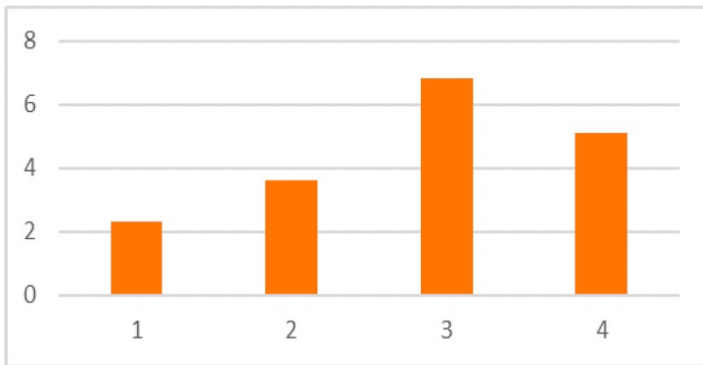


Short Intro : What is softmax

- The softmax function is popularly used to normalize the neural network output scores across all the classes

$f_{\theta}(x) = y$
 $g(\theta^T x) = \begin{bmatrix} 0.7 & 0.5 \\ 0.8 & 0.2 \end{bmatrix}$

$$p(z_i) = \frac{e^{z_i}}{\sum_z e^z}$$

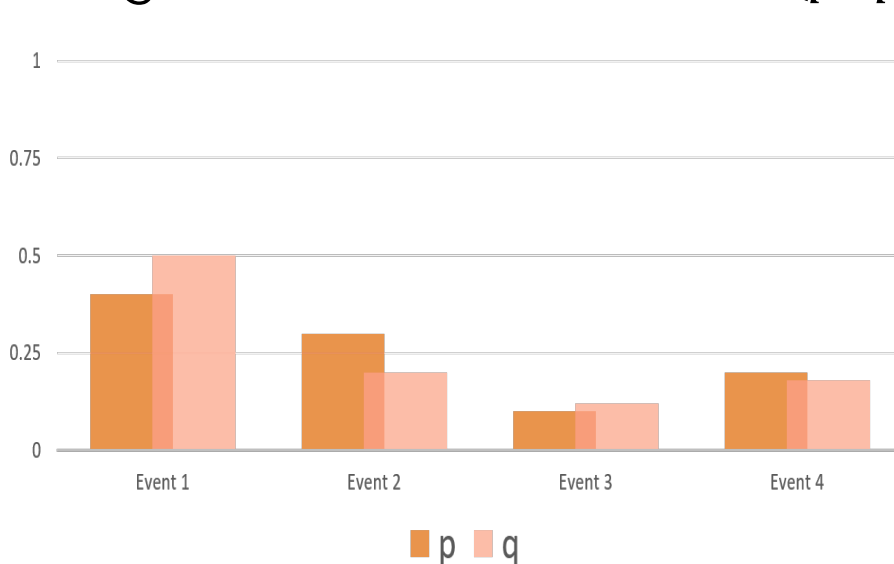


Short Intro: Cross-entropy

- The cross entropy between two probability distributions p and q over a discrete set of events is given by,

$$CE(p, q) = - \sum_i p_i \log(q_i)$$

- Single-label classification: $CE(p, q) = -\log(q_{correct})$



Handwritten calculation for single-label classification:

	not	sl	rd
q, d	0	1	0
	0.4	0.3	0.5

Calculation:

$$-0.4 \log 0.4 +$$

$$-1 \cdot \log 0.3 +$$

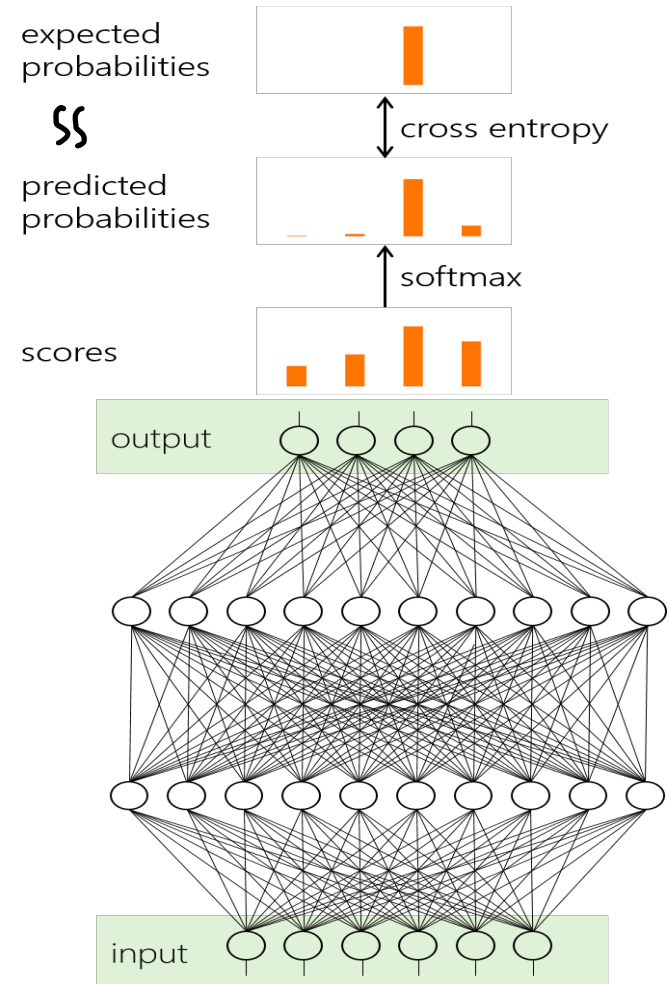
$$-0.5 \log 0.5$$

Short Intro: CE with softmax

- Cross entropy with softmax is a popular loss function for classification

$$\mathcal{L}_{CE} = -\log\left(\frac{e^{z_{correct}}}{\sum_{z \in Z} e^z}\right)$$

$\mathcal{T}(\theta)$



L2R Loss Functions

PointWise Loss

PairWise Loss

ListWise Approaches

Pointwise Loss

Regression-based or classification-based approaches are popular

- Regression loss
 - Given $\langle q, d \rangle$ predict the value of $y_{q,d}$
 - E.g., square loss for binary or categorical labels,

$$L_{Squared} = \|y_{q,d} - f(\vec{x}_{q,d})\|^2$$

- where, $y_{q,d}$ is (generally) the actual value of the label

Pointwise Loss

Supervision X, Y_{GT}
← Training $X_{tr}, Y_{GT, tr}$
Test X_{test}

Regression-based or classification-based approaches are popular

- Classification loss

- Given $\langle q, d \rangle$ predict the value of $y_{q,d}$
- E.g., Cross-Entropy with Softmax over categorical labels Y ,

$$L_{CE}(q, d, y_{q,d}) = -\log(p(y_{q,d}|q, d)) = -\log\left(\frac{e^{s_{y_{q,d}}}}{\sum_{y \in Y} e^{s_y}}\right)$$

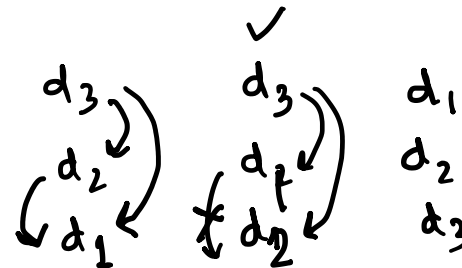
- where, $s_{y_{q,d}}$ is model's score for label $y_{q,d}$

Pairwise Loss

$$s_i = 0.7, s_j = 0.5$$

$$1 + 0.5 - 0.7 = 0.8$$

q



$$s_i = q$$

Minimizes the average number of **inversions** in ranking

➤ i.e., $d_i \succ_q d_j$ but d_j is ranked higher than d_i

$$d_1, d_2, d_3 \quad \text{sim}(d_1, d_3) > \text{sim}(d_1, d_2)$$

• For $\langle q, d_i \rangle$ and $\langle q, d_j \rangle$ Feature vectors: \vec{x}_i and \vec{x}_j

• Model scores: $s_i = f(\vec{x}_i)$ and $s_j = f(\vec{x}_j)$

• Say, d_i is more relevant. \rightarrow
 $s_i > s_j$

c.

• Pairwise loss generally has the following form [Chen et al., 2009],

$$L_{\text{pairwise}} = \phi(s_i - s_j)$$

where, ϕ can be,

• Hinge function $\phi(z) = \max(0; 1 - z)$

• Logistic function $\phi(z) = \log(1 + e^{-z})$

$$\max(0; 1 - \frac{0.7 - 0.5}{s_i - s_j}) \quad 1.2$$

RankNet

$$p(d_i|q) \quad p(d_j|q)$$

$$p(R=1|d_i,q) \quad p(R=1|d_j,q)$$

RankNet [Burges et al., 2005] is a **pairwise** loss function—popular choice for training neural L2R models and also an industry favourite [Burges, 2015]

Predicted probabilities: $p_{ij} = p(s_i > s_j) \equiv \frac{e^{\gamma \cdot s_i}}{e^{\gamma \cdot s_i} + e^{\gamma \cdot s_j}} = \frac{1}{1 + e^{-\gamma(s_i - s_j)}}$

and $p_{ji} \equiv \frac{1}{1 + e^{-\gamma(s_j - s_i)}}$

Desired probabilities: $\bar{p}_{ij} = 1$ and $\bar{p}_{ji} = 0$

Computing cross-entropy between \bar{p} and p ,

$$\mathcal{L}_{RankNet} = -\bar{p}_{ij} \log(p_{ij}) - \bar{p}_{ji} \log(p_{ji}) \quad (9)$$

$$= -\log(p_{ij}) \quad (10)$$

$$= \log(1 + e^{-\gamma(s_i - s_j)}) \quad (11)$$

CE with softmax over Documents

An alternative loss function assumes a single relevant document d^+ and compares it against the full collection D

Probability of retrieving d^+ for q is given by the softmax function,

$$p(d^+|q) = \frac{e^{\gamma \cdot s(q, d^+)}}{\sum_{d \in D} e^{\gamma \cdot s(q, d)}} \quad \text{1M x 64} \quad (12)$$

The cross entropy loss is then given by,

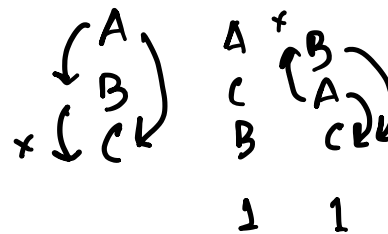
$$\mathcal{L}_{\text{CE}}(q, d^+, D) = -\log(p(d^+|q)) \quad (13)$$

$$= -\log\left(\frac{e^{\gamma \cdot s(q, d^+)}}{\sum_{d \in D} e^{\gamma \cdot s(q, d)}}\right) \quad (14)$$

CE vs RankNet

- If we consider only a pair of relevant and non-relevant documents in the denominator, CE reduces to RankNet
- Computing the denominator is prohibitively expensive -- L2R models typically consider few negative candidates
- Large body of work in NLP to deal with similar issue that may be relevant to future L2R models
 - Importance sampling, negative sampling

ListWise

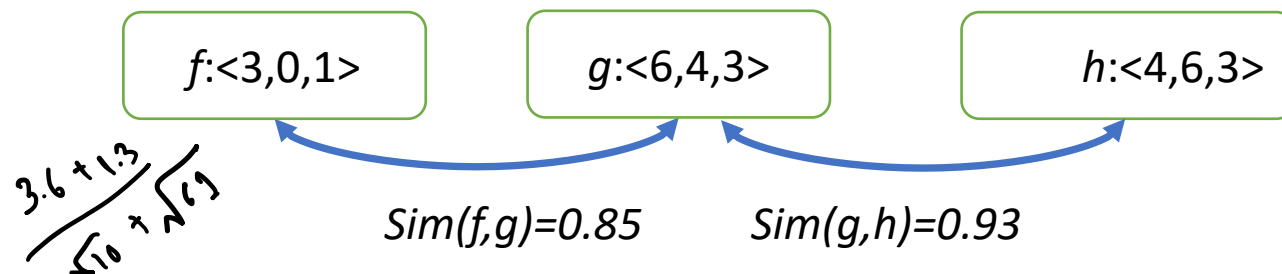


A Simple Example:

- function f : $f(A)=3, f(B)=0, f(C)=1$ $ACB \leftarrow$
- function h : $h(A)=4, h(B)=6, h(C)=3$ $BAC \leftarrow$
- **ground truth g : $g(A)=6, g(B)=4, g(C)=3$ $ABC \leftarrow$**

Question: which function is closer to ground truth?

- Based on pointwise similarity: $\text{sim}(f,g) < \text{sim}(g,h)$.
- Based on pairwise similarity: $\text{sim}(f,g) = \text{sim}(g,h)$
- Based on cosine similarity between score vectors?



- Acc. To position-wise discount f should be closer to g .

Permutation Probability Distribution

AB C	A	B	C
A	0/1	0/1	0/1

- Question:
 - How to represent a ranked list?

A - B

A - C

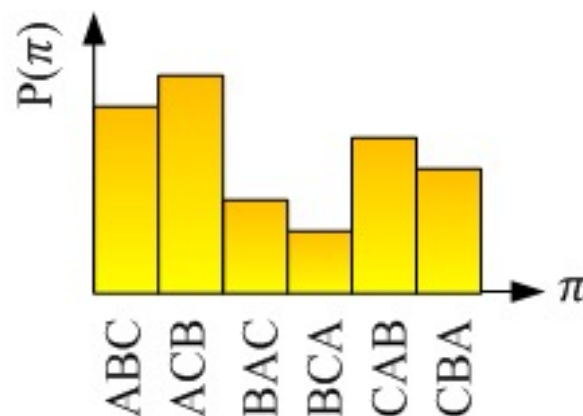
B - C

Solution

- Ranked list \leftrightarrow Permutation probability distribution
- More informative representation for ranked list: permutation and ranked list has 1-1 correspondence.

1 0 0 0 0 0

$f: f(A)=3, f(B)=0, f(C)=1;$
 Ranking by f : ABC



Luce Model: Defining Permutation Probability

- Probability of permutation π is defined as

$$P_s(\pi) = \prod_{j=1}^n \frac{\varphi(s_{\pi(j)})}{\sum_{k=j}^n \varphi(s_{\pi(k)})}$$

$P(ABC) > P(ACB)$

- Example:

$$P_f(ABC) = \frac{\varphi(f(A))}{\varphi(f(A)) + \varphi(f(B)) + \varphi(f(C))} \cdot \frac{\varphi(f(B))}{\varphi(f(B)) + \varphi(f(C))} \cdot \frac{\varphi(f(C))}{\varphi(f(C))}$$

P(A ranked No.1)

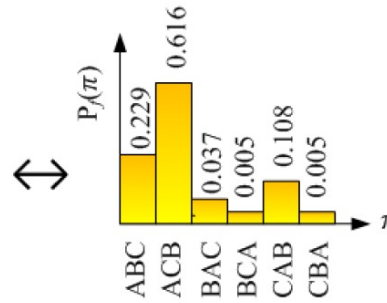
P(B ranked No.2 | A ranked No.1)
 = P(B ranked No.1) / (1 - P(A ranked No.1))

P(C ranked No.3 | A ranked No.1, B ranked No.2)

Distance between Ranked Lists

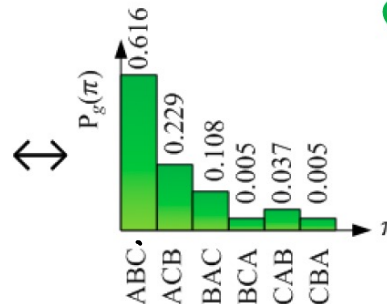
$\varphi = \exp$

$f: f(A) = 3, f(B)=0, f(C)=1;$
Ranking by f : ABC



Using **KL-divergence**
to measure difference
between distributions

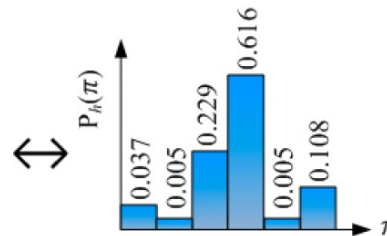
\rightarrow $g: g(A) = 6, g(B)=4, g(C)=3;$
Ranking by g : ABC



Closer!

$$dis(f,g) = 0.46$$

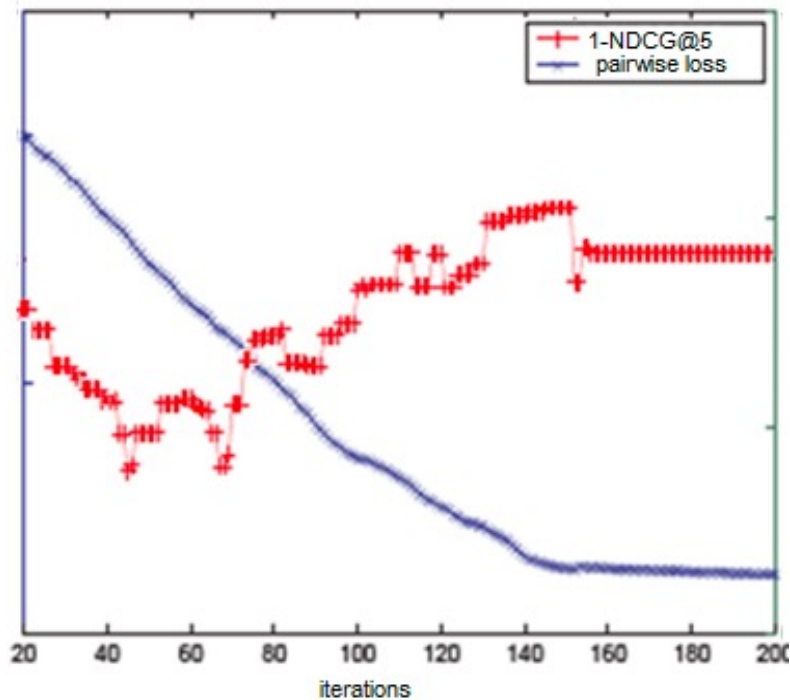
$h: h(A) = 4, h(B)=6, h(C)=3;$
Ranking by h : ACB



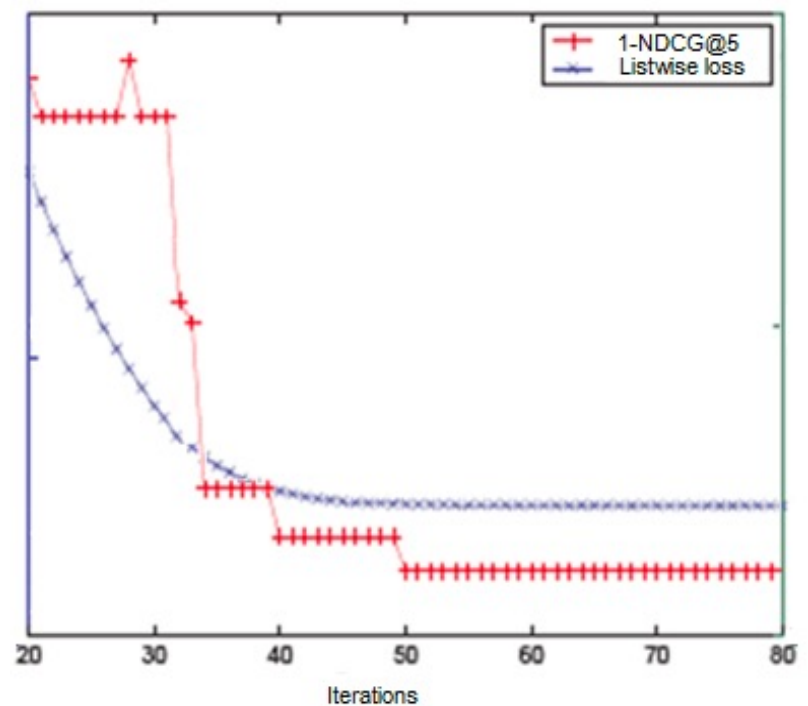
$$dis(g,h) = 2.56$$

Experimental Results (ListNet)

(Z. Cao, T. Qin, T. Liu, et al. ICML 2007)



Pairwise (RankNet)



ListWise (ListNet)

Training Performance on TD2003 Dataset

ListNet vs ListMLE

ListNet [Cao et al., 2007]

- Compute the **probability distribution over all possible permutations** based on model score and ground-truth labels. The loss is then given by the K-L divergence between these two distributions.
- This is computationally very costly, computing permutations of only the top-K items makes it slightly less prohibitive

ideal MAP

ListMLE [Xia et al. 2008]

- Compute the probability of the ideal permutation based on the ground truth. However, with categorical labels more than one permutation is possible which makes this difficult.

Supervision/Annotations

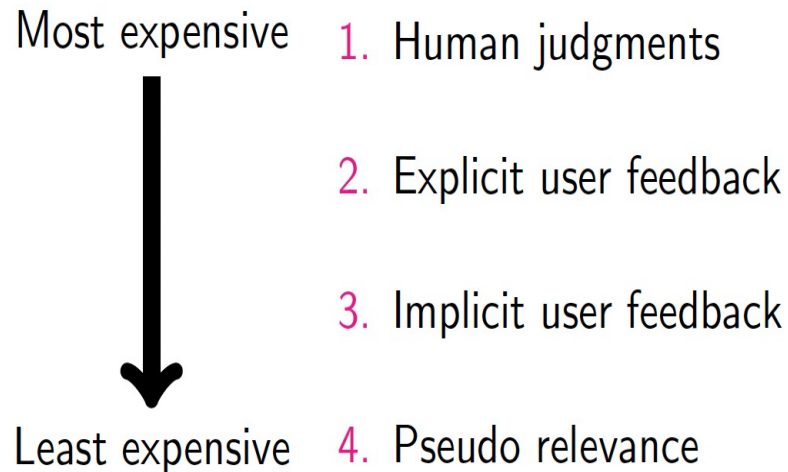
Different Levels of Supervision

Data requirements for training an offline L2R system

- Query/document pairs that encode an **ideal** ranking given a particular query.
- **Ideal ranking?** Relevance, preference, importance [Liu, 2009], novelty & diversity [Clarke et al., 2008].
- What about personalization? Triples of user, query and document.
- **Related to evaluation.** Pairs also used to compute popular offline evaluation measures.
- **Graded or binary.** "documents may be relevant to a different degree"
- **Absolute or relative?** Zheng et al. [2007]

Satisfying Data-Hungry Models

There are different ways to obtain query/document pairs.



Human Judgements

Human judges determine the relevance of a document for a given query.

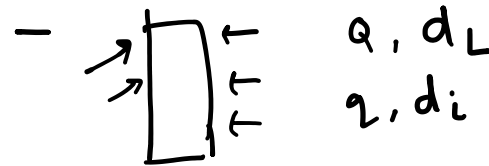
How to determine candidate query/document pairs?

- Obtaining human judgments is **expensive**.
- List of queries: sample of incoming traffic or manually curated.
- Use an existing rankers to obtain rankings and pool the outputs [Sparck Jones and van Rijsbergen, 1976].
- Trade-off between **number of queries (shallow)** and **judgments** (deep) [Yilmaz and Robertson, 2009].

Explicit User Feedback

- When presenting results to the user, ask the user to explicitly judge the documents.
- Unfortunately, users are only rarely willing to give explicit feedback [Joachims et al., 1997].

Extracting pairs from click-through data (training)



Extract implicit judgments from search engine interactions by users.

- Assumption: user clicks \Rightarrow relevance (or, preference).
- Virtually **unlimited** data at very low cost, but interpretation is more difficult.
- **Presentation bias**: users are more likely to click higher-ranked links.
 - How to deal with **presentation bias**? Joachims [2003] suggest to interleave different rankers and record preference.
- Chains of queries (i.e., search sessions) can be identified within logs and more fine-grained user preference can be extracted [Radlinski and Joachims, 2005].