# Learning to Rank

srihari@buffalo.edu
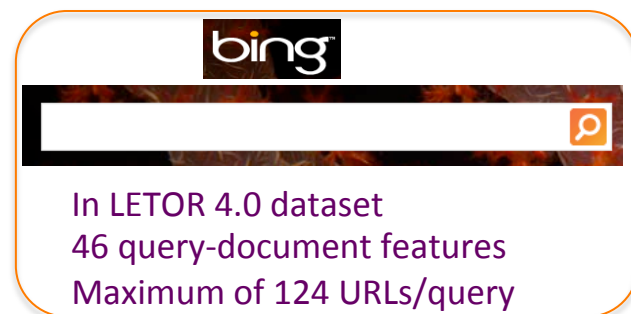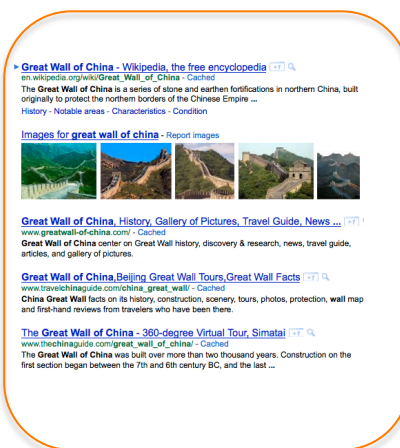
# Learning to Rank Problem

- LeToR

- Multiple Inputs

- Target Value
  - $t$ is discrete (eg, $1,2..6$) in training set but a continuous value in $[1,6]$ is learnt and used to rank objects

# Regression with multiple inputs: LeToR

**Input ($x_i$):**

($d$ Features of Query-URL pair)

- Log frequency of query in anchor text
- Query word in color on page
- # of images on page
- # of (out) links on page
- PageRank of page
- URL length
- URL contains "~"
- Page length
- TF/IDF

($d > 200$)

In LETOR 4.0 dataset
46 query-document features
Maximum of 124 URLs/query

Yahoo! data set has $d = 700$

**Target ($t$):**
Relevance Value
(0,1,2):higher value is better match

**Regression returns**

continuous value, $y$

Allows fine-grained ranking of URLs

# Query-URL Features

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

| feature id | feature description | stream | comments |
|---|---|---|---|
| **Feature List of Microsoft Learning to Rank Datasets** | | | |
| 1 | | body | |
| 2 | | anchor | |
| 3 | covered query term number | title | |
| 4 | | url | |
| 5 | | whole document | |
| 6 | | body | |
| 7 | | anchor | |
| 8 | covered query term ratio | title | |
| 9 | | url | |
| 10 | | whole document | |
| 11 | | body | |
| 12 | | anchor | |
| 13 | stream length | title | |
| 14 | | url | |
| 15 | | whole document | |
| 16 | | body | |
| 17 | | anchor | |
| 18 | IDF(Inverse document frequency) | title | |
| 19 | | url | |
| 20 | | whole document | |
| 21 | | body | |
| 22 | | anchor | |
| 23 | sum of term frequency | title | |
| 24 | | url | |
| 25 | | whole document | |
| 26 | | body | |
| 27 | | anchor | |
| 28 | min of term frequency | title | |
| 29 | | url | |
| 30 | | whole document | |
| 31 | | body | |
| 32 | | anchor | |
| 33 | max of term frequency | title | |
| 34 | | url | |
| 35 | | whole document | |
| 36 | | body | |
| 37 | | anchor | |
| 38 | mean of term frequency | title | |
| 39 | | url | |
| 40 | | whole document | |
| 41 | | body | |
| 42 | | anchor | |
| 43 | variance of term frequency | title | |
| 44 | | url | |
| 45 | | whole document | |
| 46 | | body | |
| 47 | sum of stream length | anchor | |
| 48 | normalized term | title | |
| 49 | frequency | url | |
| 50 | | whole document | |

raw frequency: $tf(t,d) = f_{t,d}$ , log-normalized: $1 + \log(f_{t,d})$

$idf(t,D) = \frac{\log N}{|\{d \in D: t \in d\}|}$ N=|D|. no. of docs in corpus, $|\{d \in D: t \in d\}|$: no. of docs where term $t$ appears

# Feature Statistics

- Most of 46 features are normalized as continuous values from 0 to 1, exception some features are all 0s'.

| Feature | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| Min | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Max | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |
| Mean | 0.254 | 0.1598 | 0.1392 | 0.2158 | 0.1322 | 0.1614 | 0 | 0 | 0 | 0 | 0 | 0.2841 | 0.1382 | 0.2109 | 0.1218 |

| 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 0.2879 | 0.1533 | 0.2258 | 0.3057 | 0.3332 | 0.1534 | 0.5473 | 0.5592 | 0.5453 | 0.5622 | 0.1675 | 0.1377 | 0.1249 | 0.126 | 0.2109 | 0.1705 |

| 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| 0.1694 | 0.1603 | 0.1275 | 0.0762 | 0.0762 | 0.0728 | 0.5479 | 0.5574 | 0.5502 | 0.5673 | 0.4321 | 0.3361 | 0 | 0.1065 | 0.1211 |

# Returning to LeToR Problem

- Try:
- Several Basis Functions
- Quadratic Regularization
- Express results as $E_{RMS}$
  - rather than as squared error $E(\mathrm{w}*)$ or as Error Rate with thresholded results

$$E_{RMS} = \sqrt{2E(\mathrm{w}*)/N}$$