

## SVM untuk Ranking

Dr. rer. nat. Hendri Murfi

Intelligent Data Analysis (IDA) Group

Departemen Matematika, Universitas Indonesia – Depok 16424

Telp. +62-21-7862719/7863439, Fax. +62-21-7863439, Email. hendri@ui.ac.id

## Model Linear

- Model dasar yang akan digunakan adalah model linear, yaitu model yang merupakan kombinasi linear dari fungsi basis, yaitu:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Dimana  $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$  adalah variabel input, dan  $\mathbf{w} = (w_0, w_1, \dots, w_D)^T$  adalah parameter,  $\boldsymbol{\phi}(\mathbf{x})$  adalah *fungsi basis*,  $M$  adalah jumlah fungsi basis

- Pada banyak metode  $\phi_0(\mathbf{x}) = 1$ , sehingga  $w_0$  berfungsi sebagai bias
- Ada banyak pilihan yang mungkin untuk fungsi basis  $\boldsymbol{\phi}(\mathbf{x})$ , misal fungsi polinomial, fungsi gaussian, fungsi sigmoidal, fungsi pemetaan, dll

# Model Linear

## Kutukan Dimensi

- Model linear memiliki sifat-sifat yang penting baik dari aspek komputasi maupun analitik. Penggunaan model linear dengan pendekatan parametrik pada metode klasik memiliki keterbatasan pada aplikasi praktis disebabkan oleh kutukan dimensi (*curse of dimensionality*)

$$y(x, w) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

Diagram illustrating the growth of parameters in a polynomial model. The equation shows terms up to the third order. Callouts indicate the number of parameters for each term: a blue box with '1' points to the  $w_3 x^3$  term, and a blue box with  $D^3$  points to the  $w_{ijk} x_i x_j x_k$  term.

untuk model berorde M, maka pertumbuhan jumlah parameter w proporsional dengan  $D^M$

3

# Model Linear

## Solusi Praktis

- Untuk menerapkan metode ini pada masalah skala besar, pendekatan umum yang dilakukan adalah membuat fungsi basis dapat beradaptasi dengan data pembelajaran, dan mengeset data pembelajaran tersebut sebagai pusat fungsi basis. Selanjutnya, memilih sebagian dari data pembelajaran tersebut selama proses training (nonparametrik).
- Dasarnya adalah bahwa data real biasanya memiliki sifat mulus, artinya perubahan sedikit pada data input hanya akan memberikan sedikit perubahan pada output
- Penggunaan fungsi kernel sebagai fungsi basis adalah salah satu contoh pendekatan seperti ini yang banyak digunakan saat ini.

4

# Metode Kernel

## Fungsi Kernel: Definisi

- Fungsi kernel adalah suatu fungsi  $k$  yang mana untuk semua vektor input  $\mathbf{x}, \mathbf{z}$  akan memenuhi kondisi

$$k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

dimana  $\phi(\cdot)$  adalah fungsi pemetaan dari ruang input ke ruang fitur

- Dengan kata lain, fungsi kernel adalah fungsi perkalian dalam (*inner product*) pada ruang fitur.

5

# Metode Kernel

## Fungsi Kernel: Contoh

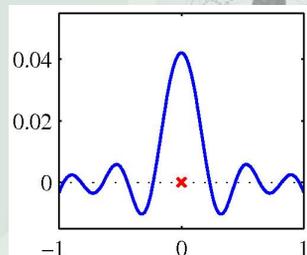
- Salah satu contoh fungsi kernel yang banyak digunakan adalah fungsi Gaussian, yaitu:

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$$

dimana  $\mathbf{x}'$  adalah „inti“ yang biasanya dipilih dari data pembelajaran.

- Misal untuk data  $x=5$  dengan target  $t=0.04$ , maka data tsb dapat digambarkan dengan fungsi basis sbb:

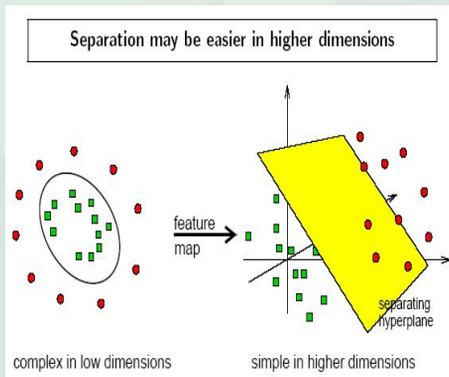
$$y(x) = 0.04 * \exp(-\|x-5\|^2 / 2\sigma^2)$$



6

# Metode Kernel

## Fungsi Kernel: Keuntungan



- Fungsi kernel memungkinkan kita untuk mengimplementasikan suatu model pada ruang dimensi lebih tinggi (ruang fitur) tanpa harus mendefinisikan fungsi pemetaan dari ruang input ke ruang fitur
- Sehingga, untuk kasus yang *non-linearly separable* pada ruang input, diharapkan akan menjadi *linearly separable* pada ruang fitur
- Selanjutnya, kita dapat menggunakan *hyperplane* sebagai *decision boundary* secara efisien

7

# Metode Kernel

## Fungsi Kernel: Penggunaan

- Secara umum, ada dua cara penggunaan metode kernel pada *machine learning*, yaitu:
  - Penggunaan langsung, yaitu fungsi kernel digunakan sebagai fungsi basis dari model machine learning tersebut, contoh: *radial basis function*
  - Penggunaan tidak langsung melalui *kernel trick*, yaitu merepresentasikan suatu model ke dalam representasi dual yang mengandung *inner product* dari fungsi pemetaan, contoh: **support vector machine**, *kernel linear regression*, *kernel Perceptron*, *kernel PCA*, dll

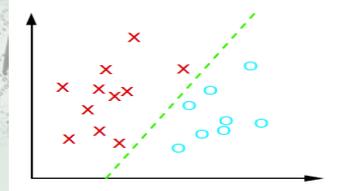
# Support Vector Machine

## Klasifikasi dan Regresi

Diberikan data pembelajaran  $\{\mathbf{x}_n, t_n\}$ ,  $n = 1$  sd  $N$

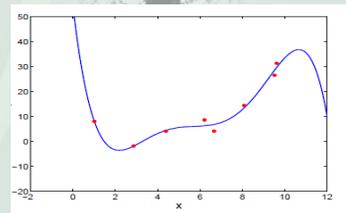
- **Klasifikasi**

- $t_n$  bernilai diskrit (kelas) berhingga tidak terurut (skala nominal)
- Permasalahan: menentukan *decision boundary* yang dapat mengklasifikasi data dengan benar



- **Regresi**

- $t_i$  bernilai kontinu (bilangan real)
- Permasalahan: menentukan fungsi regresi yang dapat memprediksi nilai data dengan benar



9

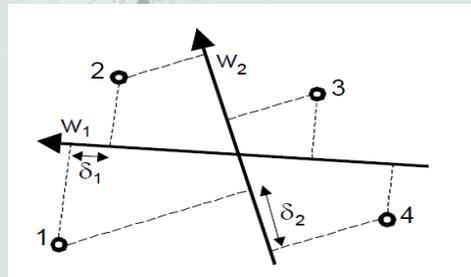
# Support Vector Machine

## Regresi Ordinal

Diberikan data pembelajaran  $\{\mathbf{x}_n, t_n\}$ ,  $n = 1$  sd  $N$

- **Regresi Ordinal**

- $t_n$  bernilai diskrit (kelas) berhingga seperti klasifikasi
- Terdapat urutan diantara elemen  $t_n$  (skala ordinal) seperti regresi
- Permasalahan: menentukan fungsi regresi ordinal yang dapat mengklasifikasikan data dengan benar



10

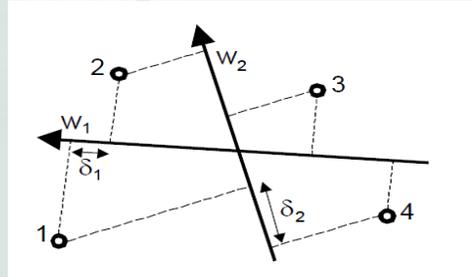
# Support Vector Machine

## Ranking

Diberikan data pembelajaran  $\{x_n, t_n\}, n = 1 \text{ sd } N$

- **Ranking**

- $t_n$  bernilai real yang menyatakan urutan (ranking) dari data
- Permasalahan: menentukan fungsi ranking yang dapat mengurutkan data dengan benar

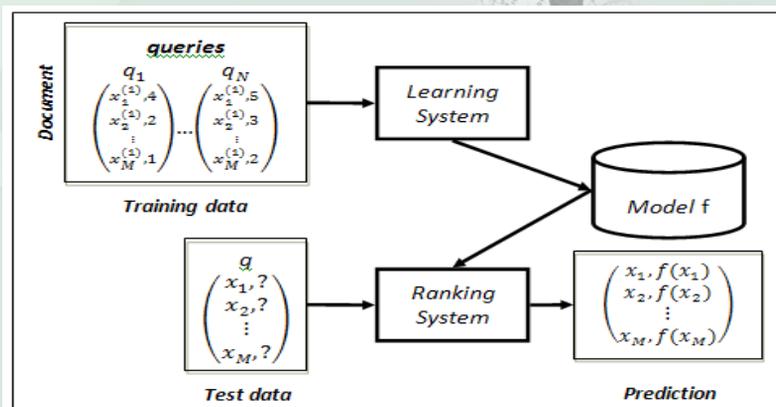


11

# SVM untuk Ranking

## Learning to Rank

- *Learning to Rank* adalah suatu bidang penelitian yang bertujuan membangun fungsi ranking dengan menggunakan machine learning
- Misal  $x_i$  adalah suatu item/dokumen,  $q_j$  adalah suatu query,  $f$  adalah suatu fungsi ranking, maka *Learning to Rank* dapat diilustrasikan sbb:



12

# SVM untuk Ranking

## Bentuk Umum

- Model linear yang akan digunakan sebagai fungsi regresi ordinal memiliki bentuk umum sbb:

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

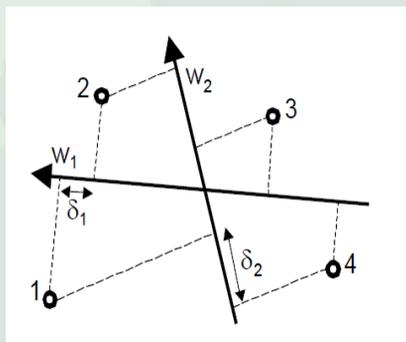
dimana  $\mathbf{x}$  adalah vektor input,  $\mathbf{w}$  adalah parameter bobot, dan  $\phi(\mathbf{x})$  adalah fungsi pemetaan

13

# SVM untuk Ranking

## Formulasi Masalah

- Diberikan data pembelajaran  $\{(\mathbf{x}(q_k, d_n), t)\}$ , dimana  $\mathbf{x}(q_k, d_n) \in R^D$  adalah suatu vektor dimana dimensi/fitur merupakan parameter-parameter kemiripan antara query  $q_k \in R^S, k=1..K$  dan item/dokumen  $d_n \in R^S, n=1..N$ , dan  $t \in \{r_1, r_2, \dots, r_p\}$  adalah suatu skala ordinal.



- Permasalahan adalah menentukan fungsi  $y(\mathbf{x})$  sedemikian sehingga untuk sembarang pasangan data training  $(\mathbf{x}_i, t_i)$  dan  $(\mathbf{x}_j, t_j)$  maka

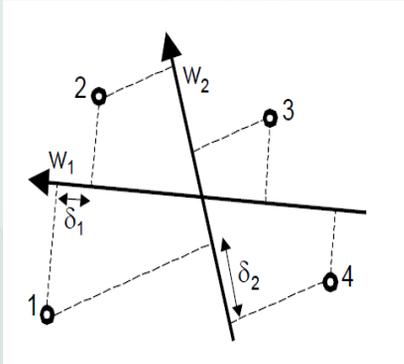
$$y(\mathbf{x}_i) > y(\mathbf{x}_j) \iff t_i > t_j$$

Contoh: Diberikan data training  $(\mathbf{x}_1, 4), (\mathbf{x}_2, 3), (\mathbf{x}_3, 2), (\mathbf{x}_4, 1)$ , maka  $\mathbf{w}_1$  adalah fungsi yang benar

14

# SVM untuk Ranking

## Formulasi Masalah



- Misal  $P_k$  adalah himpunan dari pasangan  $(i,j)$  untuk suatu query  $q_k$ , dimana  $x_i$  memiliki ranking yang lebih tinggi dari  $x_j$ :

$$P_k = \{(i,j) : t_i > t_j\}$$

Contoh: dari contoh sebelumnya, maka  $P_1 = \{(1,2), (1,3), (1,4), (2,3), (2,4), (3,4)\}$

15

# SVM untuk Ranking

## Formulasi Masalah

- Maka penentuan fungsi  $y(\mathbf{x})$  harus memenuhi kondisi berikut ini:

$$\forall (i, j) \in P_1 : w^T \phi(x(q_1, d_i)) \geq w^T \phi(x(q_1, d_j))$$

$$\forall (i, j) \in P_2 : w^T \phi(x(q_2, d_i)) \geq w^T \phi(x(q_2, d_j))$$

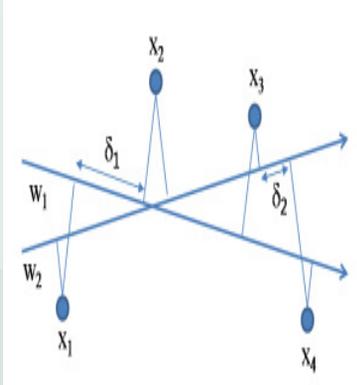
⋮

$$\forall (i, j) \in P_k : w^T \phi(x(q_k, d_i)) \geq w^T \phi(x(q_k, d_j))$$

16

# SVM untuk Ranking

## Formulasi Masalah



- Asumsikan semua data dapat di ranking dengan benar (*linearly rankable*), maka jarak antara proyeksi dua data  $\mathbf{x}(q_k, d_i)$  dan  $\mathbf{x}(q_k, d_j)$  pada  $\mathbf{w}$  adalah:

$$\frac{w^T (\phi(x(q_k, d_i)) - \phi(x(q_k, d_j)))}{\|w\|}$$

- Margin ( $\delta_i$ ) adalah jarak terdekat dari proyeksi dua data training pada vektor bobot  $\mathbf{w}$ . Sehingga, memaksimumkan margin dapat dideskripsikan sbb:

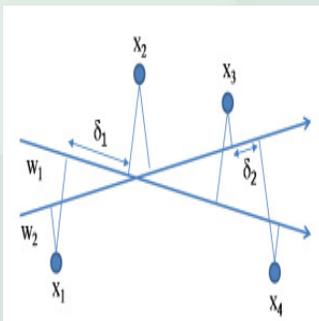
$$\arg \max_w \left\{ \frac{1}{\|w\|} \min_{i,j} [w^T (\phi(x(q_k, d_i)) - \phi(x(q_k, d_j)))] \right\}$$

17

# SVM untuk Ranking

## Formulasi Masalah

- Solusi langsung dari masalah optimasi sebelumnya akan sangat kompleks, sehingga perlu dikonversi ke masalah yang ekuivalen yang lebih mudah diselesaikan



- Salah satu metode adalah menggunakan bentuk kanonik, yaitu:

$$w^T (\phi(x(q_k, d_i)) - \phi(x(q_k, d_j))) = 1$$

untuk data yang terdekat. Selanjutnya, semua data pembelajaran harus memenuhi kondisi berikut ini:

$$w^T (\phi(x(q_k, d_i)) - \phi(x(q_k, d_j))) \geq 1$$

18

# SVM untuk Ranking

## Formulasi Masalah

- Sehingga masalah penentuan fungsi ranking  $y(\mathbf{x})$  dapat ditulis sbb:

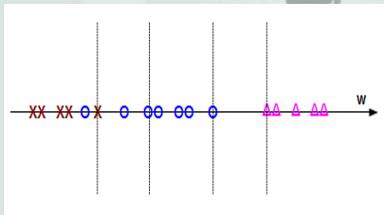
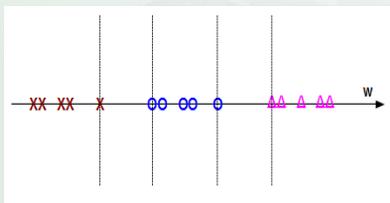
$$\begin{aligned} \arg \min_w & \frac{1}{2} \|w\|^2 \\ \text{s.t} & \quad \forall (i, j) \in P_1 : w^T \phi(x(q_1, d_i)) \geq w^T \phi(x(q_1, d_j)) \\ & \quad \forall (i, j) \in P_2 : w^T \phi(x(q_2, d_i)) \geq w^T \phi(x(q_2, d_j)) \\ & \quad \vdots \\ & \quad \forall (i, j) \in P_k : w^T \phi(x(q_k, d_i)) \geq w^T \phi(x(q_k, d_j)) \end{aligned}$$

- Dengan kata lain, penentuan fungsi ranking  $y(\mathbf{x})$  menjadi masalah pemrograman kuadrat (*quadratic programming*), yaitu meminimumkan suatu fungsi kuadrat dengan kendala pertidaksamaan linear

19

# SVM untuk Ranking

## Soft Margin: Landasan

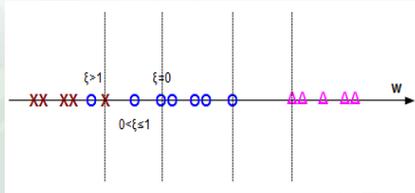


- Diberikan data pembelajaran  $\{\mathbf{x}(q_k, d_n), t\}$
- Pada formulasi sebelumnya, semua data diasumsikan *linearly rankable* pada ruang fitur  $\phi(\mathbf{x})$ .
- Dalam aplikasi praktis, kondisi tersebut sering tidak dapat terpenuhi bahkan setelah data di transformasi ke ruang fitur  $\phi(\mathbf{x})$ .
- Masalah ini diatasi dengan membuat margin lunak (*soft margin*) yang memungkinkan beberapa data pada „urutan yang salah“

20

# SVM untuk Ranking

Soft Margin: Formulasi Masalah



- Untuk merealisasikan *soft margin* ini, diperkenalkan variabel *slack*,  $\xi_{ijk} \geq 0$  untuk masing-masing pasangan data pembelajaran

- Variabel *slack* tersebut bernilai

$$\forall (i, j) \in P_k : w^T [\phi(x(q_k, d_i)) - w^T \phi(x(q_k, d_j))] = 1 - \xi_{ijk}$$

- Sehingga,  $\xi_{ijk} = 0$  adalah pasangan data yang terletak pada *margin*,  $0 < \xi_{ijk} \leq 1$  adalah pasangan data yang terletak didalam margin, dan  $\xi_{ijk} > 1$  adalah pasangan data yang salah ranking

21

# SVM untuk Ranking

Soft Margin: Formulasi Masalah

- Sehingga masalah penentuan fungsi ranking  $y(\mathbf{x})$  dapat ditulis sbb [1]:

$$\arg \min_w \frac{1}{2} \|w\|^2 + C \sum \xi_{ijk}$$

$$s.t \quad \forall (i, j) \in P_1 : w^T [\phi(x(q_1, d_i)) - w^T \phi(x(q_1, d_j))] \geq 1 - \xi_{ijk}$$

$$\forall (i, j) \in P_2 : w^T [\phi(x(q_2, d_i)) - w^T \phi(x(q_2, d_j))] \geq 1 - \xi_{ijk}$$

⋮

$$\forall (i, j) \in P_k : w^T [\phi(x(q_k, d_i)) - w^T \phi(x(q_k, d_j))] \geq 1 - \xi_{ijk}$$

$$\xi_{ijk} \geq 0$$

dimana parameter  $C > 0$  akan mengontrol *trade-off* antara lebar margin (kapabilitas generalisasi) dan jumlah pasangan data training yang tertukar urutannya (kesalahan urutan)

22

# SVM untuk Ranking

Solusi Masalah: Lagrange Multipliers

- Masalah optimisasi diatas tampak ekuivalen dengan masalah optimisasi dari SVM untuk klasifikasi, dimana setiap vektor diganti dengan pasangan vektor jarak  $[\phi(\mathbf{x}(\mathbf{q}_k, \mathbf{d}_i)) - \phi(\mathbf{x}(\mathbf{q}_k, \mathbf{d}_k))]$ . Sehingga penyelesaian masalah optimisasi ini dapat mengadaptasi pendekatan yang digunakan oleh SVM untuk klasifikasi

23

# SVM untuk Ranking

Prediksi Data Baru

- Misal diberikan himpunan data  $\{\mathbf{x}(\mathbf{q}_{new}, \mathbf{d}_1), \mathbf{x}(\mathbf{q}_{new}, \mathbf{d}_2), \dots, \mathbf{x}(\mathbf{q}_{new}, \mathbf{d}_N)\}$ , maka ranking dari data tersebut ditentukan berdasarkan fungsi ranking  $y(\mathbf{x})$ , yaitu:

$$\mathbf{x}_i \propto \mathbf{x}_j \text{ jika } y(\mathbf{x}_i) > y(\mathbf{x}_j)$$

dimana  $\propto$  didefinisikan „memiliki ranking yang lebih tinggi“

24

# SVM untuk Ranking

Algoritma & Perangkat Lunak

- Ada beberapa algoritma dan perangkat lunak yang telah dikembangkan untuk memecahkan masalah optimasi pada SVM untuk Ranking, antara lain:
  - SVM<sup>light</sup> [1][2] (<http://svmlight.joachims.org>)
  - SVM<sup>rank</sup> [3] (<http://svmlight.joachims.org>)

25

# SVM untuk Ranking

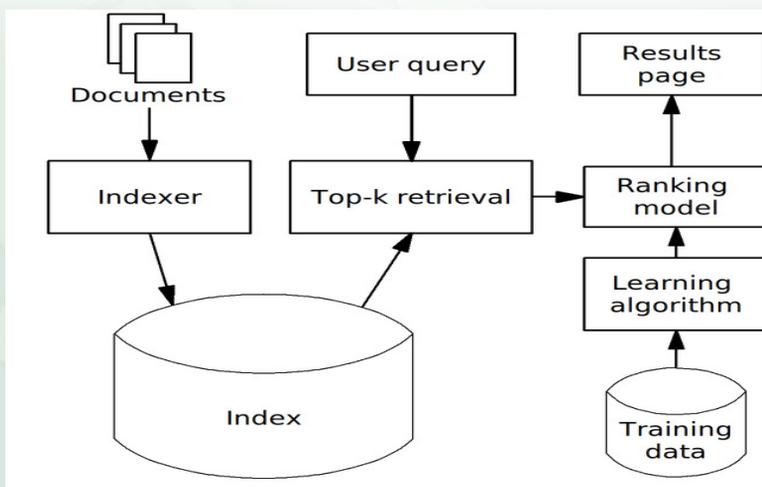
Aplikasi

- **Information Retrieval** untuk meranking dokumen/web yang relevan dengan suatu query yang diberikan oleh pengguna
- **Machine Translation** untuk meranking sekumpulan hipotesis translasi
- **Computational Biology** untuk meranking kandidat struktur 3-D dari suatu protein
- **Recommender System** untuk meranking artikel berita terkait dengan berita yang sedang dibaca oleh pengguna

26

# SVM untuk Ranking

Learning to Rank for Information Retrieval



Contoh arsitektur dari suatu *machine-learned search engine* (sumber: www.wikipedia.org)

# SVM untuk Ranking

Data Training

- Salah satu data training yang sering digunakan adalah LETOR, yaitu koleksi data yang diorganisasi oleh Microsoft yang menyimpan penentuan relevansi web terhadap query yang diberikan oleh pengguna pada mesin pencari Microsoft Bing (<http://www.bing.com>)
- Pada contoh kasus ini akan digunakan data LETOR 4.0 MQ2008 (<http://research.microsoft.com/users/LETOR/>)

Dokumen							
1	qid:10264	1:0.138801	2:0.000000	3:0.000000		46:0.966667	#docid = Gx004-93-7097963
1	qid:10264	1:0.000000	2:0.000000	3:0.000000		46:0.000000	#docid = Gx010-40-4497720
2	qid:10264	1:0.072555	2:0.666667	3:1.000000	...	46:1.000000	#docid = Gx016-32-14546147
0	qid:10264	1:0.078864	2:0.333333	3:0.000000		46:0.266667	#docid = Gx020-25-8391882
1	qid:10264	1:0.031546	2:0.000000	3:0.000000		46:0.100000	#docid = Gx025-94-0531672
2	qid:10264	1:1.000000	2:0.000000	3:0.500000		46:0.533333	#docid = Gx026-03-13004845
2	qid:10264	1:0.056782	2:1.000000	3:1.000000		46:0.000000	#docid = Gx048-02-13747475
2	qid:10264	1:0.198738	2:0.666667	3:1.000000	...	46:0.000000	#docid = Gx268-53-13016636
0	qid:10266	1:0.022052	2:0.000000	3:0.000000		46:0.052758	#docid = Gx000-11-6487904
0	qid:10266	1:0.037392	2:0.000000	3:0.000000		46:0.013189	#docid = Gx000-22-1713857

Relevansi
Fitur
Dokumen ID

# SVM untuk Ranking

Data Training: Fitur

- Masing-masing data training terdiri dari 46 fitur yang merupakan tingkat kemiripan (similarity) antara query dan web

No	Deskripsi	No	Deskripsi	No	Deskripsi
1	TF of body	17	DL of anchor	33	BM25 of URL
2	TF of anchor	18	DL of title	34	LMIR.ABS of URL
3	TF of title	19	DL of the URL	35	LMIR.DIR of URL
4	TF of URL	20	DL of document	36	LMIR.IM of URL
5	TF of document	21	BM25 of body	37	BM25 of document
6	IDF of body	22	LMIR.ABS of body	38	LMIR.ABS of document
7	IDF of anchor	23	LMIR.DIR of body	39	LMIR.DIR of document
8	IDF of title	24	LMIR.IM of body	40	LMIR.IM of document
9	IDF of URL	25	BM25 of anchor	41	PageRank
10	IDF of document	26	LMIR.ABS of anchor	42	Inlink number
11	TF * IDF of body	27	LMIR.DIR of anchor	43	Outlink number
12	TF * IDF of anchor	28	LMIR.IM of anchor	44	Number of slash in URL
13	TF * IDF of title	29	BM25 of title	45	Length of URL
14	TF * IDF of URL	30	LMIR.ABS of title	46	Number of child page
15	TF * IDF of document	31	LMIR.DIR of title		
16	DL of body	32	LMIR.IM of title		

29

# SVM untuk Ranking

Data Training: Fitur

Diberikan suatu himpunan dokumen  $\{d_1, d_2, \dots, d_N\}$  dan suatu query  $q$

- Term Frequency (TF)

$TF(q, d_i)$  = jumlah kata-kata (*terms*) dari query  $q$  pada dokumen  $d_i$

- Inverse Document Frequency (IDF)

$$IDF(q, d_i) = \log(N/DF(q, d_i))$$

dimana  $DF(q, d_i)$  = jumlah dokumen yang mengandung kata-kata (*terms*) dari query  $q$

- Term Frequency \* Inverse Document Frequency (TFIDF)

$$TFIDF(q, d_i) = TF(q, d_i) * IDF(q, d_i)$$

30

# SVM untuk Ranking

Satuan Ukuran Akurasi: P@n

- *Precision at Position n (P@n):*

$$P @ n = \frac{\text{jumlah web yang relevan pada posisi } n}{n}$$

Contoh:

Misal dokumen yang relevan terhadap query  $q_1$  adalah  $\{d_1, d_2, d_3, d_4\}$

Jika ranking yang diberikan oleh suatu fungsi ranking adalah  $\{d_1, d_4, d_3, d_2, d_5\}$ , maka:

- $P@1 = 1/1, P@2 = 2/2, P@3 = 3/3, P@4 = 3/4, P@5 = 4/5, P@6 = 4/6$

31

# SVM untuk Ranking

Satuan Ukuran Akurasi: MAP

- *Mean Average Precision (MAP):*

Misal himpunan dokumen yang relevan terhadap query  $q_j \in Q$  adalah  $\{d_1, d_2, \dots, d_{m_j}\}$ , dan  $R_{jk}$  adalah himpunan hasil retrieval yang diranking dari atas sampai ke dokumen  $d_k$ , maka:

$$MAP = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} P @ R_{jk}$$

dimana  $P@R_{jk}$  adalah Precision dari  $R_{jk}$ , yaitu

$$P @ R_{jk} = \frac{\text{jumlah web yang relevan pada } R_{jk}}{|R_{jk}|}$$

32

# SVM untuk Ranking

Satuan Ukuran Akurasi: MAP

Contoh:

Misal dokumen yang relevan terhadap query  $q_1$  adalah  $\{d_1, d_2, d_3, d_4\}$

Jika ranking yang diberikan oleh suatu fungsi ranking adalah  $\{d_1, d_4, d_3, d_2, d_5\}$ , maka:

- $P@R_{11} = 1/1$ ,  $P@_{21} = 4/5$ ,  $P@R_{31} = 3/3$ ,  $P@R_{41} = 2/2$
- Sehingga *Average Precision* (AP) adalah  $(1+4/5+1+1)/4 = 0.95$
- Dan MAP  $0.95/1 = 0.95$ , karena hanya ada satu query

33

# SVM untuk Ranking

Satuan Ukuran Akurasi: NDCG@n

- *Normalized Discount Cumulative Gain at Position n (NDCG@n)*

Misal  $R(j,m)$  adalah nilai relevansi dari suatu dokumen urutan ke- $m$  untuk query  $q_j$ , maka:

$$NDCG @ n = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_n \sum_{m=1}^n \frac{2^{R(j,m)} - 1}{\log(1+m)}$$

dimana  $Z_n$  adalah faktor normalisasi yang dihitung untuk membuat sedemikian rupa sehingga ranking sempurna pada posisi  $n$  akan bernilai 1

34

# SVM untuk Ranking

Satuan Ukuran Akurasi: NDCG@n

Contoh:

Diberikan nilai relevansi beberapa URL untuk suatu query sbb:

	URL	Relevansi
$d_1$	<a href="http://abc.go.com/">http://abc.go.com/</a>	Perfect: 5
$d_2$	<a href="http://abcnews.go.com/sections/">http://abcnews.go.com/sections/</a>	Excellent: 4
$d_3$	<a href="http://www.abc.net.au/">http://www.abc.net.au/</a>	Excellent: 4
$d_4$	<a href="http://abcnews.go.com">http://abcnews.go.com</a>	Excellent: 4
$d_5$	<a href="http://www.abcteach.com/">http://www.abcteach.com/</a>	Fair: 2

Tentukan nilai  $NDCG@5$  untuk urutan  $\{d_1, d_5, d_2, d_3, d_4\}$  ?

35

# SVM untuk Ranking

Satuan Ukuran Akurasi: NDCG@n

– Gain dan Cumulative Gain

	URL	Gain	CG
#1	<a href="http://abc.go.com/">http://abc.go.com/</a>	Perfect: $31=2^5-1$	31
#2	<a href="http://www.abcteach.com/">http://www.abcteach.com/</a>	Fair: $3=2^2-1$	$34 = 31 + 3$
#3	<a href="http://abcnews.go.com/sections/">http://abcnews.go.com/sections/</a>	Excellent: $15=2^4-1$	$49 = 31 + 3 + 15$
#4	<a href="http://www.abc.net.au/">http://www.abc.net.au/</a>	Excellent: 15	$64 = 31 + 3 + 15 + 15$
#5	<a href="http://abcnews.go.com/">http://abcnews.go.com/</a>	Excellent: 15	$79 = 31 + 3 + 15 + 15 + 15$

36

# SVM untuk Ranking

Satuan Ukuran Akurasi: NDCG@n

- Discounting Factor:  $\log(2) / \log(1 + \text{rank})$

	URL	Gain	DCG
#1	<a href="http://abc.go.com/">http://abc.go.com/</a>	Perfect: $31=2^5-1$	$31 = 31 \times 1$
#2	<a href="http://www.abcteach.com/">http://www.abcteach.com/</a>	Fair: $3=2^2-1$	$32.9 = 31 + 3 \times 0.63$
#3	<a href="http://abcnews.go.com/sections/">http://abcnews.go.com/sections/</a>	Excellent: $15=2^4-1$	$40.4 = 32.9 + 15 \times 0.50$
#4	<a href="http://www.abc.net.au/">http://www.abc.net.au/</a>	Excellent: 15	$46.9 = 40.4 + 15 \times 0.43$
#5	<a href="http://abcnews.go.com/">http://abcnews.go.com/</a>	Excellent: 15	$52.7 = 46.9 + 15 \times 0.39$

37

# SVM untuk Ranking

Satuan Ukuran Akurasi: NDCG@n

- Normalized factor: nilai maksimum yang mungkin untuk suatu urutan tertentu

	URL	Gain	DCG	Max DCG
#1	<a href="http://abc.go.com/">http://abc.go.com/</a>	Perfect: $31=2^5-1$	31	$31 = 31 \times 1$
#2	<a href="http://www.abcteach.com/">http://www.abcteach.com/</a>	Fair: $3=2^2-1$	32.9	$40.5 = 31 + 15 \times 0.63$
#3	<a href="http://abcnews.go.com/sections/">http://abcnews.go.com/sections/</a>	Excellent: $15=2^4-1$	40.4	$48.0 = 40.5 + 15 \times 0.50$
#4	<a href="http://www.abc.net.au/">http://www.abc.net.au/</a>	Excellent: 15	46.9	$54.5 = 48.0 + 15 \times 0.43$
#5	<a href="http://abcnews.go.com/">http://abcnews.go.com/</a>	Excellent: 15	52.7	$60.4 = 54.5 + 15 \times 0.39$

38

# SVM untuk Ranking

Satuan Ukuran Akurasi: NDCG@n

– Normalized Discounting Cumulative Gain

	URL	Gain	DCG	Max DCG	NDCG
#1	<a href="http://abc.go.com/">http://abc.go.com/</a>	Perfect: $31=2^5-1$	31	31	$1 = 31/31$
#2	<a href="http://www.abcteach.com/">http://www.abcteach.com/</a>	Fair: $3=2^2-1$	32.9	40.5	$0.81=32.9/40.5$
#3	<a href="http://abcnews.go.com/sections/">http://abcnews.go.com/sections/</a>	Excellent: $15=2^4-1$	40.4	48.0	$0.84=40.4/48.0$
#4	<a href="http://www.abc.net.au/">http://www.abc.net.au/</a>	Excellent: 15	46.9	54.5	$0.86=46.9/54.5$
#5	<a href="http://abcnews.go.com/">http://abcnews.go.com/</a>	Excellent: 15	52.7	60.4	$0.87=52.7/60.4$

39

# SVM untuk Ranking

Desain Eksperimen

- Eksperimen menggunakan 5-fold Cross Validation dengan statistik data sbb:

FOLD	TRAINING	VALIDATION	TEST
Fold 1	{S1, S2, S3}	S4	S5
Fold 2	{S2, S3, S4}	S5	S1
Fold 3	{S3, S4, S5}	S1	S2
Fold 4	{S4, S5, S1}	S2	S3
Fold 5	{S5, S1, S2}	S3	S4

FOLD	TRAINING	VALIDATION	TEST
Fold 1	9630	2707	2874
Fold 2	9404	2874	2933
Fold 3	8643	2933	3635
Fold 4	8514	3635	3062
Fold 5	9442	3062	2707

40

# SVM untuk Ranking

## Hasil Eksperimen

- Tingkat akurasi metode SVM untuk Ranking dengan menggunakan perangkat lunak SVM<sup>Rank</sup>

Performance on testing set											
Para	NDCG@1	NDCG@2	NDCG@3	NDCG@4	NDCG@5	NDCG@6	NDCG@7	NDCG@8	NDCG@9	NDCG@10	MeanNDCG
10	0.3611	0.3777	0.4088	0.4248	0.4469	0.4659	0.4738	0.4215	0.2078	0.2117	0.4577
2	0.3036	0.3450	0.3761	0.4018	0.4224	0.4395	0.4406	0.4230	0.1718	0.1738	0.4296
1	0.3461	0.3710	0.4082	0.4354	0.4452	0.4596	0.4664	0.4489	0.2433	0.2494	0.4686
0.2	0.4161	0.4682	0.4909	0.5078	0.5290	0.5391	0.5428	0.4955	0.2860	0.2892	0.5442
0.5	0.3864	0.4305	0.4589	0.4845	0.5042	0.5215	0.5290	0.4933	0.2107	0.2155	0.5159
Avg	0.36266	0.39848	0.42858	0.45086	0.46954	0.48512	0.49052	0.45644	0.22392	0.22792	0.4832

Para	P@1	P@2	P@3	P@4	P@5	P@6	P@7	P@8	P@9	P@10	MAP
10	0.4167	0.3942	0.3868	0.3654	0.3436	0.3216	0.2949	0.2748	0.2564	0.2423	0.4502
2	0.3631	0.3471	0.3397	0.3264	0.3121	0.2972	0.2748	0.2540	0.2385	0.2229	0.4213
1	0.4013	0.3758	0.3694	0.3487	0.3223	0.3015	0.2830	0.2675	0.2484	0.2357	0.4529
0.2	0.4968	0.4904	0.4544	0.4299	0.4089	0.3864	0.3576	0.3368	0.3171	0.2981	0.5284
0.5	0.4586	0.4268	0.4013	0.3774	0.3503	0.3259	0.3003	0.2779	0.2633	0.2465	0.4950
Avg	0.4273	0.40686	0.39032	0.36956	0.34744	0.32652	0.30212	0.2822	0.26474	0.2491	0.46956

sumber: <http://research.microsoft.com/en-us/um/beijing/projects/letor/letor4baseline.aspx>

## Referensi

- (1) T. Joachims. *Optimizing Searching Engines Using Clickthrough Data*. Proceeding of the ACM Conference on Knowledge Discovery and Data Mining, 2002
- (2) T. Joachim. *Making Large-Scale SVM Learning Practical*. In B. Schoelkopf, C. J. C. Burges, and A. J. Smola (Eds), *Advances in Kernel Methods – Support Vector Learning*, pp. 169-184, MIT Press, 1999
- (3) T. Joachims. *Training Linear SVMs in Linear Time*. Proceeding of the ACM Conference on Knowledge Discovery and Data Mining, 2006