

We will discuss string kernels
today:

Text Classification using String Kernels

Huma Lodhi

HUMA@CS.RHUL

Craig Saunders

CRAIG@CS.RHUL

John Shawe-Taylor

JOHN@CS.RHUL

Nello Cristianini

NELLO@CS.RHUL

Chris Watkins

CHRISW@CS.RHUL

*Department of Computer Science,
Royal Holloway, University of London,
Egham, Surrey TW20 0EX, UK*

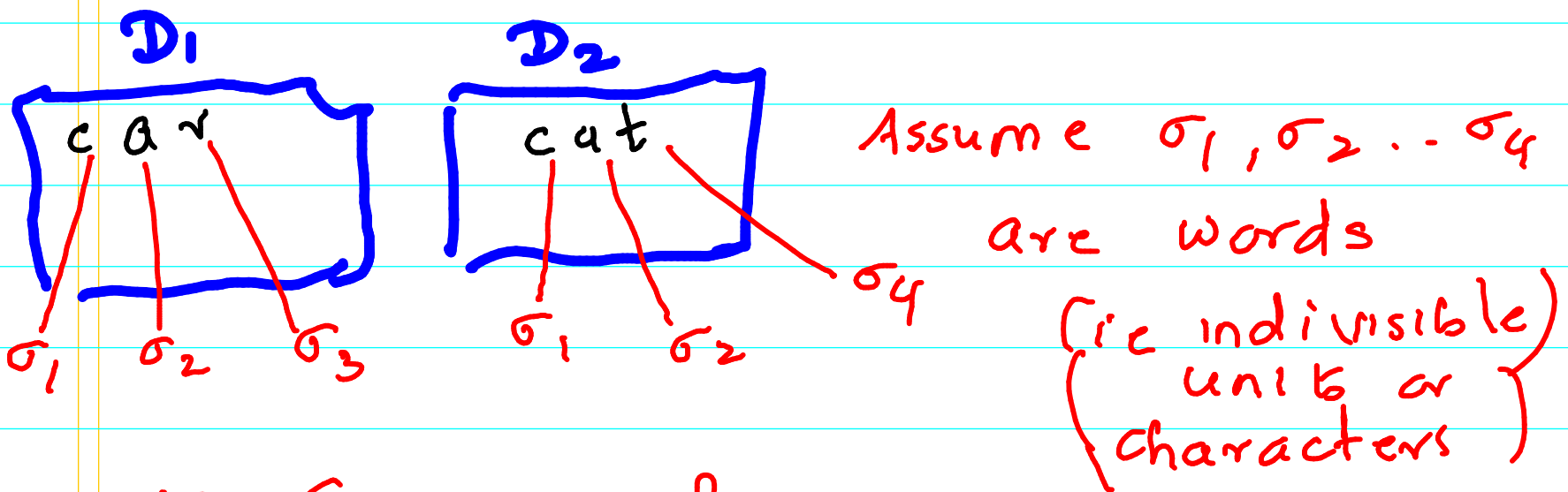
Editor: Bernhard Schölkopf

Definition 1 (String subsequence kernel- SSK) Let Σ be a finite alphabet. A string is a finite sequence of characters from Σ , including the empty sequence. For strings s, t , we denote by $|s|$ the length of the string $s = s_1 \dots s_{|s|}$, and by st the string obtained by concatenating the strings s and t . The string $s[i : j]$ is the substring $s_i \dots s_j$ of s . We say that u is a subsequence of s , if there exist indices $\mathbf{i} = (i_1, \dots, i_{|u|})$, with $1 \leq i_1 < \dots < i_{|u|} \leq |s|$, such that $u_j = s_{i_j}$, for $j = 1, \dots, |u|$, or $u = s[\mathbf{i}]$ for short. The length $l(\mathbf{i})$ of the subsequence in s is $i_{|u|} - i_1 + 1$. We denote by Σ^n the set of all finite strings of length n , and by Σ^* the set of all strings

$$\Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n. \quad (1)$$

We now define feature spaces $F_n = \mathbb{R}^{\Sigma^n}$. The feature mapping ϕ for a string s is given by defining the u coordinate $\phi_u(s)$ for each $u \in \Sigma^n$. We define

$$\phi_u(s) = \sum_{\mathbf{i}: u=s[\mathbf{i}]} \lambda^{l(\mathbf{i})}, \quad (2)$$



$\Sigma = \{\sigma_1, \dots, \sigma_4\}$ is the vocabulary

You could generate sequences of characters from Σ

Eg: $S_1 = c - a$ $S_5 = a - a$
 $S_2 = c - t$
 $S_3 = r - t$
 $S_4 = c - c$

Goal: ① Define feature ϕ so as to measure similarity between D_1 & D_2 in the space of sequences. $S \in \Sigma^*$

② want to compute $\phi^T \phi$ efficiently

③ Go about ① & ② so that "gaps" in the matchings are allowed but with penalty. (special case of no gaps is with infinite penalty)

Defn:

Substring: $S = c_a_t$

Substrings
(contiguous)
of S

$S[1:3] = c_a_t$
 $S[2:3] = a_t$
 $S[2:2] = a$

Subsequence: $S = r_o_l_l_e_r$

subsequences
(non contiguous)
of S

$S[1, 3, 6] = r_l_r$
 $S[2, 4] = o_l$

subsequence (any non-zero entry indicates subsequence)

	c-a	c-t	a-t	b-a	b-t	c-r	a-r	b-r
$\lambda \in (0, 1]$ $\phi(\text{cat})$	λ^2	λ^3	λ^2	0	0	0	0	0
$\phi(\text{car})$	λ^2	0	0	0	0	λ^3	λ^2	0
$\phi(\text{bat})$	0	0	λ^2	λ^2	λ^3	0	0	0
$\phi(\text{bar})$	0	0	0	λ^2	0	0	λ^2	λ^3

Eg: string subsequence



$$\phi(\text{catca}) \mid \overset{\text{c_a}}{\lambda^2 + \lambda^2 + \lambda^2 \times \lambda^3} \\ = 2\lambda^2 + \lambda^5$$

In this setting, each match carries a penalty of λ & each gap carries a penalty of λ .

H/W → To handle "no gap" situation you might need to separate the two penalties & set the gap penalty to ∞

→ Why penalize "matches"?

• of course, if you separate the two penalties, match penalty could be set to 1

• Nevertheless, you often want "short" feature/string matches to be rewarded more than "long" feature/string matches

→ Reason multiplicative scheme for penalty is used is because $\phi^T \phi$ becomes efficiently computable.

H/W: How about additive or o/r forms of penalty or reward?

Let us now formally
define ϕ :

$\phi_u(s)$ for each $u \in \Sigma^n$. We define

$$\phi_u(s) = \sum_{\mathbf{i}: u=s[\mathbf{i}]}$$

match index array.

eg: $u = c_a$

$s = a b r c t q a c a$

$\Rightarrow i = [4, 7], [4, 9], [8, 9]$

$$\neq \lambda^3 + \lambda^5 + \lambda$$

Definition of string kernel

$$K_n(s, t) = \sum_{u \in \Sigma^n} \langle \phi_u(s) \cdot \phi_u(t) \rangle = \sum_{u \in \Sigma^n} \sum_{\mathbf{i}: u=s[\mathbf{i}]} \lambda^{l(\mathbf{i})} \sum_{\mathbf{j}: u=t[\mathbf{j}]} \lambda^{l(\mathbf{j})}$$
$$= \sum_{u \in \Sigma^n} \sum_{\mathbf{i}: u=s[\mathbf{i}]} \sum_{\mathbf{j}: u=t[\mathbf{j}]} \lambda^{l(\mathbf{i})+l(\mathbf{j})}$$

The two i/p documents

Enumerating the u 's is a formidable task!!!

Can we obtain $K_n(s, t)$ from $K_r(p, q)$ for $\left\{ \begin{array}{l} p, q \text{ substring of } s, t \\ r < n \end{array} \right.$ using recursion?

Obviously:

$$K_n(sx, t) = K_n(s, t) + \dots$$

Vinod's recursion

$$K_n(sx, t) = K_n(s, t) + \sum_{i=1}^{|s|} K'_{n-1}(s, i, t) \lambda^{|s|-i}$$

where

$$K'_n(s, i, t) = \sum_{j=1}^{i-1} K'_{n-1}(s, j, t) \lambda^{i-j+1}$$

We will come back to Vinod's recursion after looking at that of the authors:

Base cases:

$$K_n(s, t) = 0 \quad \text{if } n = 0$$

$$= 0 \quad \text{if } |s| \text{ or } |t| < n$$

$$K'_i(s, t) = \sum_{u \in \Sigma^i} \sum_{\mathbf{i}: u=s[\mathbf{i}]} \sum_{\mathbf{j}: u=t[\mathbf{j}]} \lambda^{|s|+|t|-i_1-j_1+2},$$

$$i = 1, \dots, n-1,$$

Like us, the authors also realize the need for an auxiliary $k'_i(s, t)$

$k'_i(s, t)$ is almost similar to $K_n(s, t)$ except that penalty is not based on $l(i)$ or $l(j)$ but is based on the distances from ends of s & t to the first matches i_1 in i & j_1 in j

Let us write recursion for

$$K_i'(s, t):$$

Base case:

$$K_i'(s, t) = 0 \begin{cases} \text{if } i = 0 \\ \text{if } |s| \text{ or } |t| < i \end{cases}$$

$$K_i'(s, t) = \lambda K_i'(s, t) + \sum_{\substack{j: t_j = x \\ |t| - j + 2}} K_{i-1}'(s, t_{1..j-1})$$

factoring in
case where x
has no match
in any $u \in \Sigma$

to account for
 x that matched
 t_j (neither
 x nor t_j are
in argument
for $K_{i-1}'(s, t_{1..j-1})$)

H/w: Write $K_n(sx, t)$
in terms of $K_n(s, t)$
& $K'_i(s, t)$

$$K'_0(s, t) = 1, \text{ for all } s, t,$$

$$K'_i(s, t) = 0, \text{ if } \min(|s|, |t|) < i,$$

$$K_i(s, t) = 0, \text{ if } \min(|s|, |t|) < i,$$

$$K'_i(sx, t) = \lambda K'_i(s, t) + \sum_{j:t_j=x} K'_{i-1}(s, t[1:j-1]) \lambda^{|t|-j+2},$$
$$i = 1, \dots, n-1,$$

$$K_n(sx, t) = K_n(s, t) + \sum_{j:t_j=x} K'_{n-1}(s, t[1:j-1]) \lambda^2.$$