

Kernel

HPV

1 01,2 2 1,2,3
 3 , 2
 {jgyoung⁰, sjaugh, btzhang }@bi.snu.ac.kr

HPV Risk Classification Using Kernel Based Learning

Je -Gun Joung^{0,1,2} Sirk June Augh² Byoung -Tak Zhang^{1,2,3}

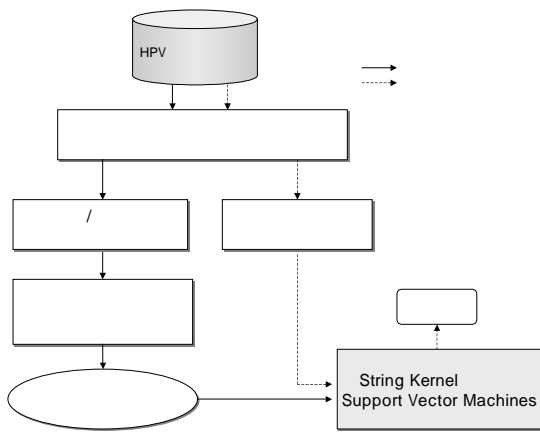
¹Interdisciplinary Program in Bioinformatics, Seoul National University

²Center for Bioinformation Technology, Seoul National University

³School of Computer Science and Engineering, Seoul National University

DNA (human papillomavirus: HPV)
 HPV 가 .
 HPV (kernel)

1. 8kb (human papillomavirus: HPV) HPV 가 .
 가 DNA HPV 85 (string kernels) HPV
 가 (genotype) HPV 가 (map) 가 (feature space)
 120 [1]. HPV 가 .
 (low - risk type) (high - risk type) HPV 16, 18, 2. HPV
 31 HPV , HPV type 1 HPV
 가 HPV HPV HPV
 (text mining) 가 가 [2]. 가 HPV HPV
 (decision tree) HPV 가 (text mining) (multiple sequence alignment)
 HPV 가 가 가 (point)
 가 HPV (window) 가
 (subsequence)



1. HPV

S 가 score(S)

$(S_{Pos}) / P_M(S_{Neg})$ 가

$\log(P_M)$ 가

SVM

HPV

3. Support Vector Machines
Support Vector Machines
1995 Vapnik

[3].
2

(classification)

SVM

(feature space)

(boundary)

$S = \{x_i, y_i\}, i=1, \dots, n$ $\Phi(S) = \{\Phi(x_i), y_i\} = \{z_i, y_i\}, i=1, \dots, n$

SVM
(linear discriminant function)

$f(z) = \langle w \cdot x \rangle + b$ hyperplane

$f(z) = 0$

hyperplane $f(z) = \langle w \cdot x \rangle + b = 0$

$z^* \in \{z_1, \dots, z_n\}$ (Euclidean distance)

hyperplane (margin)

hyperplane (normalization) hyperplane

$1/\|w\|$

SVM 2 (quadratic programming) (maximal margin classifier)

$$\text{maximize } \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle z_i \cdot z_j \rangle,$$

$$\text{subject to } \alpha_i \geq 0 \quad (1 \leq i \leq n), \sum_{i=1}^n \alpha_i y_i = 0.$$

α Lagrange multiplier
 $\alpha_1^*, \dots, \alpha_n^*$ hyperplane
 $f^*(z) = 0$

$$f^*(z) = \sum_{i=1}^n \alpha_i^* y_i \langle z_i \cdot z \rangle + b^*$$

$$b^* = y_s - \sum_{i=1}^n \alpha_i^* y_i \langle z_i \cdot z_s \rangle \text{ for some } \alpha_s^* \neq 0$$

Φ

$$\langle z_i \cdot z_j \rangle = \langle \Phi(x_i) \cdot \Phi(x_j) \rangle,$$

$$K(x_i \cdot x_j) = \langle \Phi(x_i) \cdot \Phi(x_j) \rangle$$

4. HPV

Mismatch

Mismatch

spectrum

[4].

k -spectrum

k - 가

(subsequences)

k - 가

spectrum

(hybridization)

k - 가

가

[5].

Spectrum

spectrum

$$\Phi_k(x) = (\phi_\alpha(x))_{\alpha \in A^k}$$

$\phi_\alpha(x)$ x k -mer (k - 가) α

가 α 20

(amino acid) k -mer A^k

x_i, x_j k -spectrum:

$$K_k(x_i \cdot x_j) = \langle \Phi_k(x_i) \cdot \Phi_k(x_j) \rangle$$

k -

mer

k -spectrum

가

Mismatch

k -mer

m

mismatch

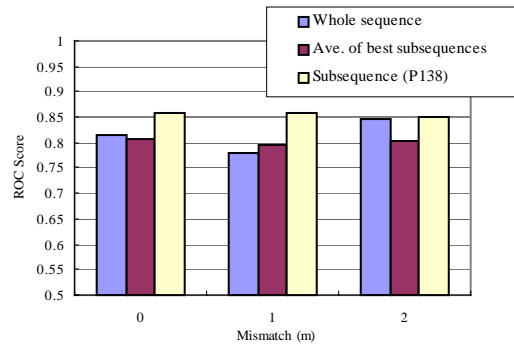
(k, m) -

mismatch

k -mer α 가

a_1

$a_2 \dots a_k$ α (k, m)
 β β
 $\phi_\beta(\alpha) = P(b_1 | a_1)P(b_2 | a_2) \dots P(b_k | a_k)$
 \cdot a b
 $P(b | a)$ a b
 (substitution)
 PAM[12] BLOSUM[6]
 β α mismatch
 $\Phi_{(k,m)}(\alpha) = (\phi_\beta(\alpha))_{\beta \in A^k}$
 $\Phi_{(k,m)}(x) = \sum_{k\text{-mers } \alpha \text{ in } x}^x \Phi_{(k,m)}(\alpha)$



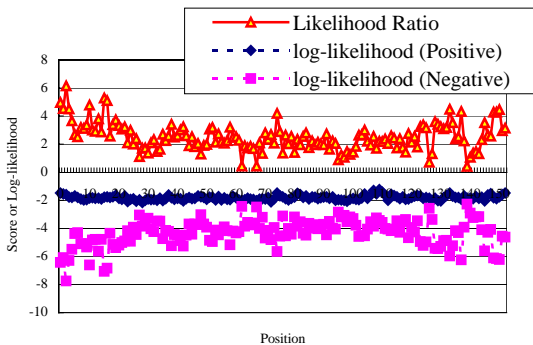
3. likelihood

mismatch

$$K_{(k,m)}(x_i \cdot x_j) = \langle \Phi_{(k,m)}(x_i) \cdot \Phi_{(k,m)}(x_j) \rangle$$

5.

8 E6
 80
 E6
 가



2. E6

2 E6

가

log - likelihood

가

log - likelihood

가

SVM

3

ROC (receiver - operating characteristic)

가

6.

HPV

HPV

SVM

가

가

IMT-

2000

- zur Hausen, H., Papillomaviruses causing cancer: evasion from host - cell control in early events in carcinogenesis, *Journal of National Cancer Inst .*, Vol. 92, pp.690 - 698, 2000.
- , , , , pp. 148 - 160, 2002.
- Vapnik, V. N., *Statistical learning theory*, Springer, 1998.
- Leslie, C., Eskin, E., Weston, J. and Noble, W., Mismatch string kernels for SVM protein classification, *NIPS 2002*. (to appear)
- Peer, I. and Shamir, R., Spectrum alignment: efficient resequencing by hybridization, *In ISMB*, pp. 260 - 268, AAAI Press, 2000.
- Henikoff, S. and Henikoff, J. G., Amino acid substitution matrices from protein blocks, *PNAS*, Vol. 89, pp. 10915-10919, 1992.