

---

# COMPUTATIONAL MODELS FOR DATA ANALYSIS

## Kernel Methods

**Alessandra Giordani**

Department of information and communication technology

University of Trento

Email: [moschitti@dit.unitn.it](mailto:moschitti@dit.unitn.it)



# Linear Classifier

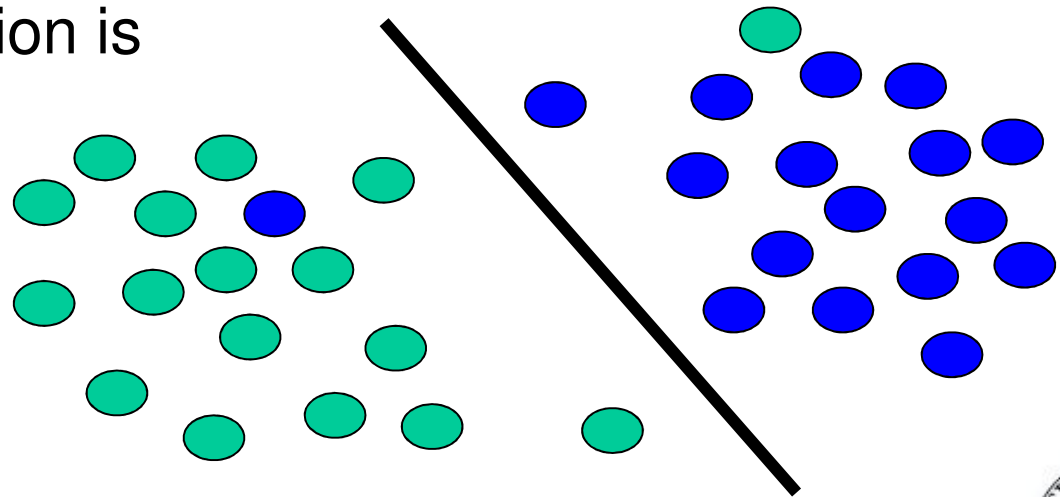
---

- The equation of a hyperplane is

$$f(\vec{X}) = \vec{X} \cdot \vec{W} + b = 0, \quad \vec{X}, \vec{W} \in \mathfrak{R}^n, b \in \mathfrak{R}$$

- $\vec{X}$  is the vector representing the classifying example
- $\vec{W}$  is the gradient of the hyperplane
- The classification function is

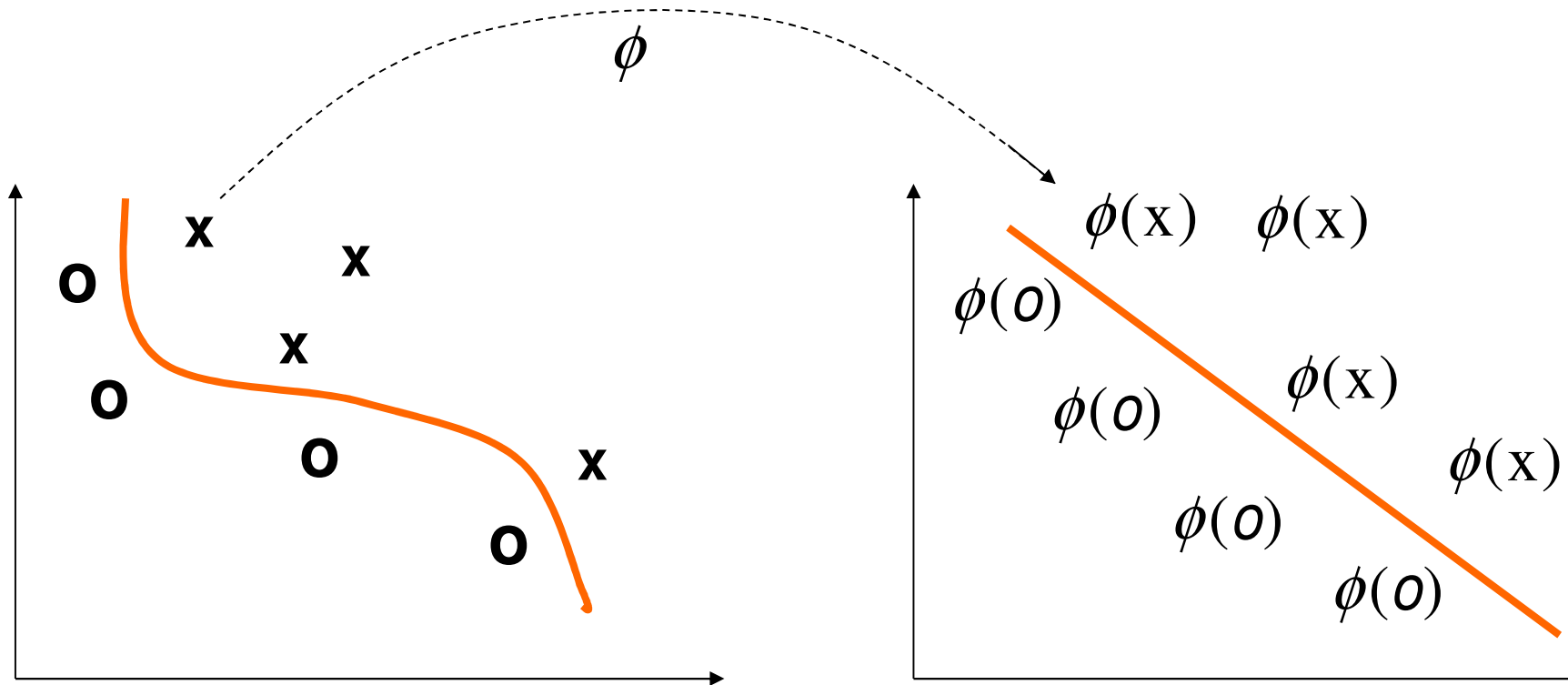
$$h(x) = \text{sign}(f(x))$$



# The main idea of Kernel Functions

---

- Mapping vectors in a space where they are linearly separable  $\vec{X} \rightarrow \phi(\vec{X})$



# A mapping example

---

- Given two masses  $m_1$  and  $m_2$ , one is constrained
- Apply a force  $f_a$  to the mass  $m_1$
- Experiments
  - Features  $m_1$ ,  $m_2$  and  $f_a$
- We want to learn a classifier that tells when a mass  $m_1$  will get far away from  $m_2$
- If we consider the Gravitational Newton Law

$$f(m_1, m_2, r) = C \frac{m_1 m_2}{r^2}$$

- we need to find when  $f(m_1, m_2, r) < f_a$



## A mapping example (2)

---

$$\vec{X} = (x_1, \dots, x_n) \rightarrow \phi(\vec{X}) = (\phi_1(\vec{X}), \dots, \phi_n(\vec{X}))$$

- The gravitational law is not linear so we need to change space

$$(f_a, m_1, m_2, r) \rightarrow (k, x, y, z) = (\ln f_a, \ln m_1, \ln m_2, \ln r)$$

- As

$$\ln f(m_1, m_2, r) = \ln C + \ln m_1 + \ln m_2 - 2 \ln r = c + x + y - 2z$$

- We need the hyperplane

$$\ln f_a - \ln m_1 - \ln m_2 + 2 \ln r - \ln C = 0$$

$(\ln m_1, \ln m_2, -2 \ln r) \cdot (x, y, z) - \ln f_a + \ln C = 0$ , we can decide without error if the mass will get far away or not



# A kernel-based Machine

## Perceptron training

---

$\vec{w}_0 \leftarrow \vec{0}; b_0 \leftarrow 0; k \leftarrow 0; R \leftarrow \max_{1 \leq i \leq l} \|\vec{x}_i\|$

do

for  $i = 1$  to  $\ell$

if  $y_i(\vec{w}_k \cdot \vec{x}_i + b_k) \leq 0$  then

$$\vec{w}_{k+1} = \vec{w}_k + \eta y_i \vec{x}_i$$

$$b_{k+1} = b_k + \eta y_i R^2$$

$k = k + 1$

endif

endfor

while an error is found

return  $k, (\vec{w}_k, b_k)$



# Kernel Function Definition

---

**Def. 2.26** *A kernel is a function  $k$ , such that  $\forall \vec{x}, \vec{z} \in X$*

$$k(\vec{x}, \vec{z}) = \phi(\vec{x}) \cdot \phi(\vec{z})$$

*where  $\phi$  is a mapping from  $X$  to an (inner product) feature space.*

- Kernels are the product of mapping functions such as

$$\vec{X} \in \mathfrak{R}^n, \quad \vec{\phi}(\vec{X}) = (\phi_1(\vec{X}), \phi_2(\vec{X}), \dots, \phi_m(\vec{X})) \in \mathfrak{R}^m$$



# The Kernel Gram Matrix

---

- With KM-based learning, the sole information used from the training data set is the Kernel Gram Matrix

$$K_{\text{training}} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_m) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_m) \\ \dots & \dots & \dots & \dots \\ k(\mathbf{x}_m, \mathbf{x}_1) & k(\mathbf{x}_m, \mathbf{x}_2) & \dots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

- If the kernel is valid,  $K$  is symmetric definite-positive .





# Valid Kernels

---

## **Def. B.11** *Eigen Values*

Given a matrix  $\mathbf{A} \in \mathbb{R}^m \times \mathbb{R}^n$ , an eigenvalue  $\lambda$  and an eigenvector  $\vec{x} \in \mathbb{R}^n - \{\vec{0}\}$  are such that

$$\mathbf{A}\vec{x} = \lambda\vec{x}$$

## **Def. B.12** *Symmetric Matrix*

A square matrix  $\mathbf{A} \in \mathbb{R}^n \times \mathbb{R}^n$  is symmetric iff  $\mathbf{A}_{ij} = \mathbf{A}_{ji}$  for  $i \neq j$   $i = 1, \dots, m$  and  $j = 1, \dots, n$ , i.e. iff  $\mathbf{A} = \mathbf{A}'$ .

## **Def. B.13** *Positive (Semi-) definite Matrix*

A square matrix  $\mathbf{A} \in \mathbb{R}^n \times \mathbb{R}^n$  is said to be positive (semi-) definite if its eigenvalues are all positive (non-negative).



# Mercer's condition

---

**Proposition 2.27** (*Mercer's conditions*)

Let  $X$  be a finite input space with  $K(\vec{x}, \vec{z})$  a symmetric function on  $X$ . Then  $K(\vec{x}, \vec{z})$  is a kernel function if and only if the matrix

$$k(\vec{x}, \vec{z}) = \phi(\vec{x}) \cdot \phi(\vec{z})$$

is positive semi-definite (has non-negative eigenvalues).

- If the Gram matrix:  $G = k(\vec{X}_i, \vec{X}_j)$  is positive semi-definite there is a mapping  $\phi$  that produces the target kernel function



# Mercer's Theorem (finite space)

---

- Let us consider  $\mathbf{K} = \left( K(\vec{x}_i, \vec{x}_j) \right)_{i,j=1}^n$
- $\mathbf{K}$  symmetric  $\Rightarrow \exists \mathbf{V}: \mathbf{K} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}'$  for Takagi factorization of a complex-symmetric matrix, where:
  - $\mathbf{\Lambda}$  is the diagonal matrix of the eigenvalues  $\lambda_t$  of  $\mathbf{K}$
  - $\vec{\mathbf{V}}_t = \left( v_{ti} \right)_{i=1}^n$  are the eigenvectors, i.e. the columns of  $\mathbf{V}$
- Let us assume lambda values non-negative

$$\phi : \vec{x}_i \rightarrow \left( \sqrt{\lambda_t} v_{ti} \right)_{t=1}^n \in \mathfrak{R}^n, i = 1, \dots, n$$



# Mercer's Theorem (sufficient conditions)

---

- Therefore

$$\Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j) = \sum_{t=1}^n \lambda_t v_{ti} v_{tj} = (\mathbf{V} \mathbf{\Lambda} \mathbf{V}')_{ij} = \mathbf{K}_{ij} = K(\vec{x}_i, \vec{x}_j)$$

- which implies that  $K$  is a kernel function



# Mercer's Theorem (necessary conditions)

---

- Suppose we have negative eigenvalues  $\lambda_s$  and eigenvectors  $\vec{v}_s$  the following point

$$\vec{z} = \sum_{i=1}^n v_{si} \Phi(\vec{x}_i) = \sum_{i=1}^n v_{si} \left( \sqrt{\lambda_t} v_{ti} \right)_t = \sqrt{\Lambda} \mathbf{V}' \vec{v}_s$$

- has the following norm:

$$\|\vec{z}\|^2 = \vec{z} \cdot \vec{z} = \sqrt{\Lambda} \mathbf{V}' \vec{v}_s \sqrt{\Lambda} \mathbf{V}' \vec{v}_s = \vec{v}_s' \mathbf{V} \sqrt{\Lambda} \sqrt{\Lambda} \mathbf{V}' \vec{v}_s =$$

$$\vec{v}_s' \mathbf{K} \vec{v}_s = \vec{v}_s' \lambda_s \vec{v}_s = \lambda_s \|\vec{v}_s\|^2 < 0$$

this contradicts the geometry of the space.



# Is it a valid kernel?

---

- It may not be a kernel so we can use  $M' \cdot M$

**Proposition B.14** *Let  $A$  be a symmetric matrix. Then  $A$  is positive (semi-) definite iff for any vector  $\vec{x} \neq 0$*

$$\vec{x}' A \vec{x} > 0 \quad (\geq 0).$$

From the previous proposition it follows that: If we find a decomposition  $A$  in  $M' M$ , then  $A$  is semi-definite positive matrix as

$$\vec{x}' A \vec{x} = \vec{x}' M' M \vec{x} = (M \vec{x})' (M \vec{x}) = M \vec{x} \cdot M \vec{x} = \|M \vec{x}\|^2 \geq 0.$$



# Valid Kernel operations

---

- $k(x,z) = k_1(x,z) + k_2(x,z)$
- $k(x,z) = k_1(x,z) * k_2(x,z)$
- $k(x,z) = \alpha k_1(x,z)$
- $k(x,z) = f(x)f(z)$
- $k(x,z) = k_1(\phi(x), \phi(z))$
- $k(x,z) = x'Bz$



# Basic Kernels for unstructured data

---

- Linear Kernel
- Polynomial Kernel
- Lexical kernel
- String Kernel





# Linear Kernel

---

- In Text Categorization documents are word vectors

$$\Phi(d_x) = \vec{x} = (0, \dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 1)$$

buy acquisition stocks sell market

$$\Phi(d_z) = \vec{z} = (0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0)$$

buy company stocks sell

- The dot product  $\vec{x} \cdot \vec{z}$  counts the number of features in common
- This provides a sort of *similarity*



# Feature Conjunction (polynomial Kernel)

---

- The initial vectors are mapped in a higher space

$$\Phi(\langle x_1, x_2 \rangle) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$$

- More expressive, as  $(x_1x_2)$  encodes

**Stock+Market vs. Downtown+Market** features

- We can smartly compute the scalar product as

$$\begin{aligned}\Phi(\vec{X}) \cdot \Phi(\vec{Z}) &= \\ &= (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1) \cdot (z_1^2, z_2^2, \sqrt{2}z_1z_2, \sqrt{2}z_1, \sqrt{2}z_2, 1) = \\ &= x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 + 2x_1z_1 + 2x_2z_2 + 1 = \\ &= (x_1z_1 + x_2z_2 + 1)^2 = (\vec{X} \cdot \vec{Z} + 1)^2 = K_{Poly}(\vec{X}, \vec{Z})\end{aligned}$$



# Document Similarity

---

**Doc 1**

**Doc 2**

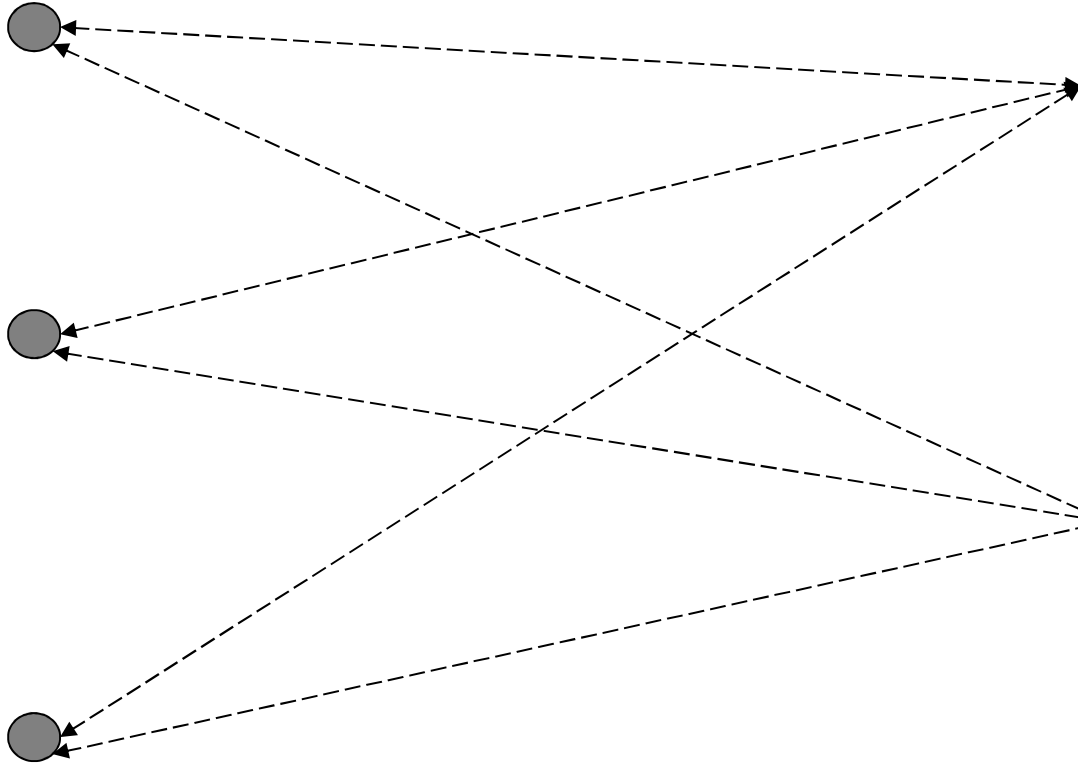
industry ●

company ●

telephone ●

product ●

market ●



# Lexical Semantic Kernel [CoNLL 2005]

---

- The document similarity is the SK function:

$$SK(d_1, d_2) = \sum_{w_1 \in d_1, w_2 \in d_2} s(w_1, w_2)$$

- where  $s$  is any similarity function between words, e.g. WordNet [Basili et al., 2005] similarity or LSA [Cristianini et al., 2002]
- Good results when training data is small



# Using character sequences

---

$$\phi(\text{"bank"}) = \vec{x} = (0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0)$$

bank      ank      bnk      bk      b

$$\phi(\text{"rank"}) = \vec{z} = (1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1)$$

rank      ank      rnk      rk      r

- $\vec{x} \cdot \vec{z}$  counts the number of common substrings

$$\vec{x} \cdot \vec{z} = \phi(\text{"bank"}) \cdot \phi(\text{"rank"}) = k(\text{"bank"}, \text{"rank"})$$



# String Kernel

---

- Given two strings, the number of matches between their substrings is evaluated
- E.g. Bank and Rank
  - B, a, n, k, Ba, Ban, Bank, Bk, an, ank, nk,...
  - R, a, n, k, Ra, Ran, Rank, Rk, an, ank, nk,...
- String kernel over sentences and texts
- Huge space but there are efficient algorithms



# Formal Definition

---

$$s = s_1, \dots, s_{|s|}$$

$$\vec{I} = (i_1, \dots, i_{|u|}) \quad u = s[\vec{I}]$$

$$\phi_u(s) = \sum_{\vec{I}:u=s[\vec{I}]} \lambda^{l(\vec{I})}, \text{ where } l(\vec{I}) = i_{|u|} - i_1 + 1$$

$$K(s, t) = \sum_{u \in \Sigma^*} \phi_u(s) \cdot \phi_u(t) = \sum_{u \in \Sigma^*} \sum_{\vec{I}:u=s[\vec{I}]} \lambda^{l(\vec{I})} \sum_{\vec{J}:u=t[\vec{J}]} \lambda^{l(\vec{J})} =$$

$$= \sum_{u \in \Sigma^*} \sum_{\vec{I}:u=s[\vec{I}]} \sum_{\vec{J}:u=t[\vec{J}]} \lambda^{l(\vec{I})+l(\vec{J})}, \text{ where } \Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n$$



# Kernel between Bank and Rank

---

B, a, n, k, Ba, Ban, Bank, an, ank, nk, Bn, Bnk, Bk and ak are the substrings of *Bank*.

R, a, n, k, Ra, Ran, Rank, an, ank, nk, Rn, Rnk, Rk and ak are the substrings of *Rank*.





# An example of string kernel computation

---

- $\phi_a(\text{Bank}) = \phi_a(\text{Rank}) = \lambda^{(i_1 - i_1 + 1)} = \lambda^{(2 - 2 + 1)} = \lambda,$
- $\phi_n(\text{Bank}) = \phi_n(\text{Rank}) = \lambda^{(i_1 - i_1 + 1)} = \lambda^{(3 - 3 + 1)} = \lambda,$
- $\phi_k(\text{Bank}) = \phi_k(\text{Rank}) = \lambda^{(i_1 - i_1 + 1)} = \lambda^{(4 - 4 + 1)} = \lambda,$
- $\phi_{an}(\text{Bank}) = \phi_{an}(\text{Rank}) = \lambda^{(i_2 - i_1 + 1)} = \lambda^{(3 - 2 + 1)} = \lambda^2,$
- $\phi_{ank}(\text{Bank}) = \phi_{ank}(\text{Rank}) = \lambda^{(i_3 - i_1 + 1)} = \lambda^{(4 - 2 + 1)} = \lambda^3,$
- $\phi_{nk}(\text{Bank}) = \phi_{nk}(\text{Rank}) = \lambda^{(i_2 - i_1 + 1)} = \lambda^{(4 - 3 + 1)} = \lambda^2$
- $\phi_{ak}(\text{Bank}) = \phi_{ak}(\text{Rank}) = \lambda^{(i_2 - i_1 + 1)} = \lambda^{(4 - 2 + 1)} = \lambda^3$

$$K(\text{Bank}, \text{Rank}) = (\lambda, \lambda, \lambda, \lambda^2, \lambda^3, \lambda^2, \lambda^3) \cdot (\lambda, \lambda, \lambda, \lambda^2, \lambda^3, \lambda^2, \lambda^3) \\ = 3\lambda^2 + 2\lambda^4 + 2\lambda^6$$

---



# String Kernels for OCR

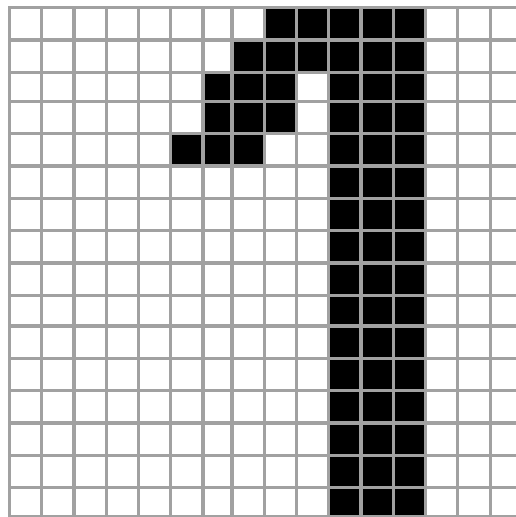
---

0123456789  
0123456789  
0123456789  
0123456789  
0123456789  
0123456789  
0123456789  
0123456789  
0123456789  
0123456789  
0123456789

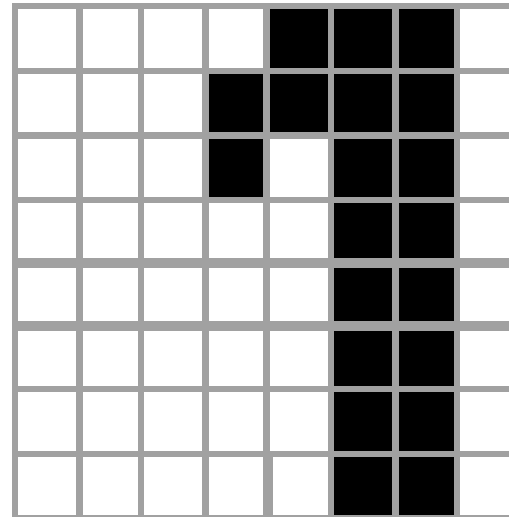


# Pixel Representation

---



(a)



(b)

Figure 6: Resampling of an image from  $16 \times 16$  to  $8 \times 8$  format



# Sequence of bits

---

L1	00011100
.	00111100
.	00101100
.	00001100
.	00001100
L8	00001100

$$SK(im_a, im_b) = \sum_{i=1..8} SK(L_a^i, L_b^i)$$



# Results

---

- Using columns+rows+diagonals

Digit	Precision	Recall	F1
0	97.78	97.78	97.78
1	95.45	93.33	94.38
2	93.62	97.78	95.65
3	93.33	93.33	93.33
4	97.83	100.00	98.90
5	97.67	93.33	95.45
6	100.00	97.78	98.88
7	91.84	100.00	95.74
8	93.18	91.11	92.13
9	93.02	88.89	90.91
<b>Multiclass accuracy</b>	<b>95.33</b>		



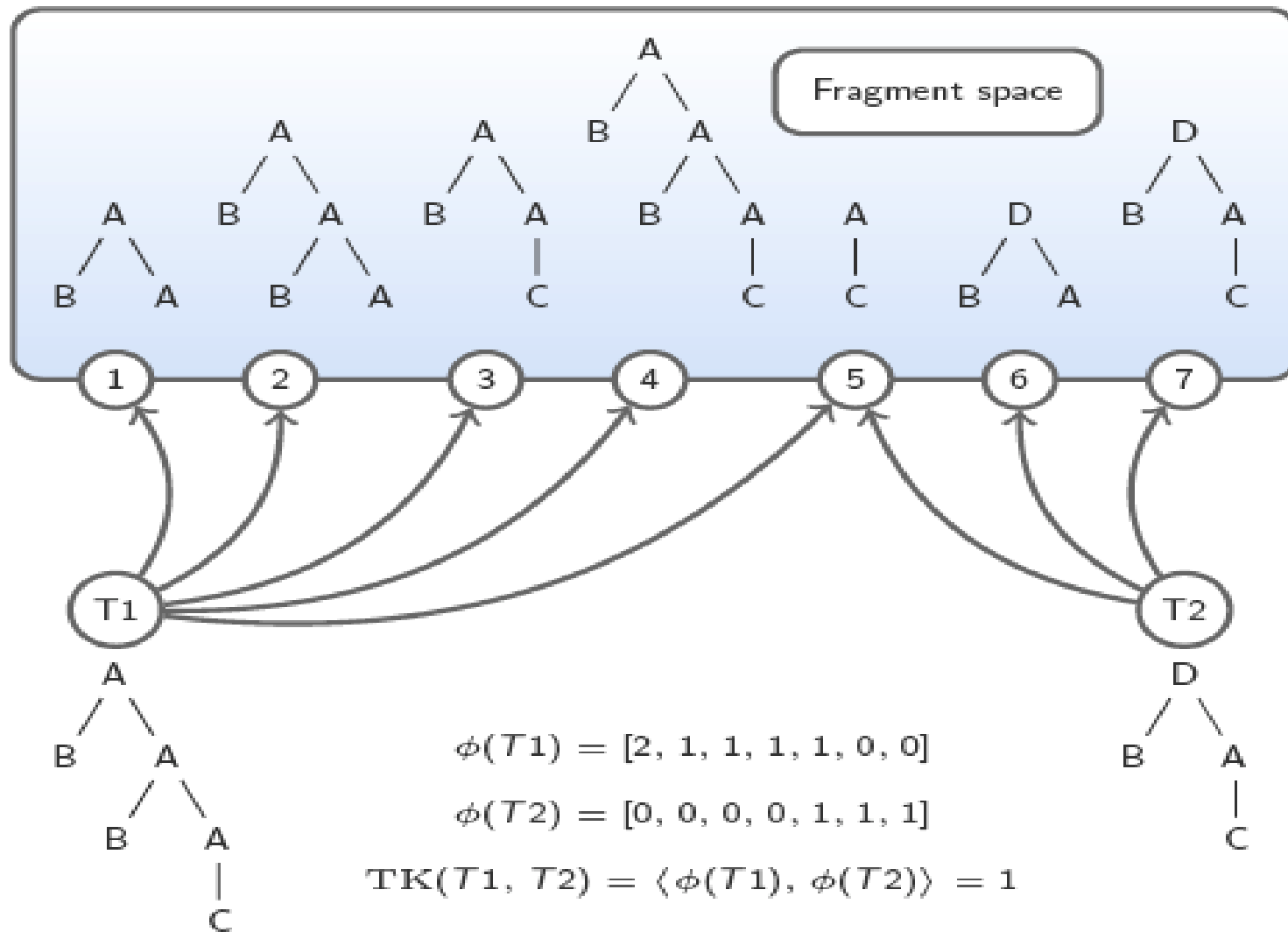
# Tree kernels

---

- Subtree, Subset Tree, Partial Tree kernels
- Efficient computation



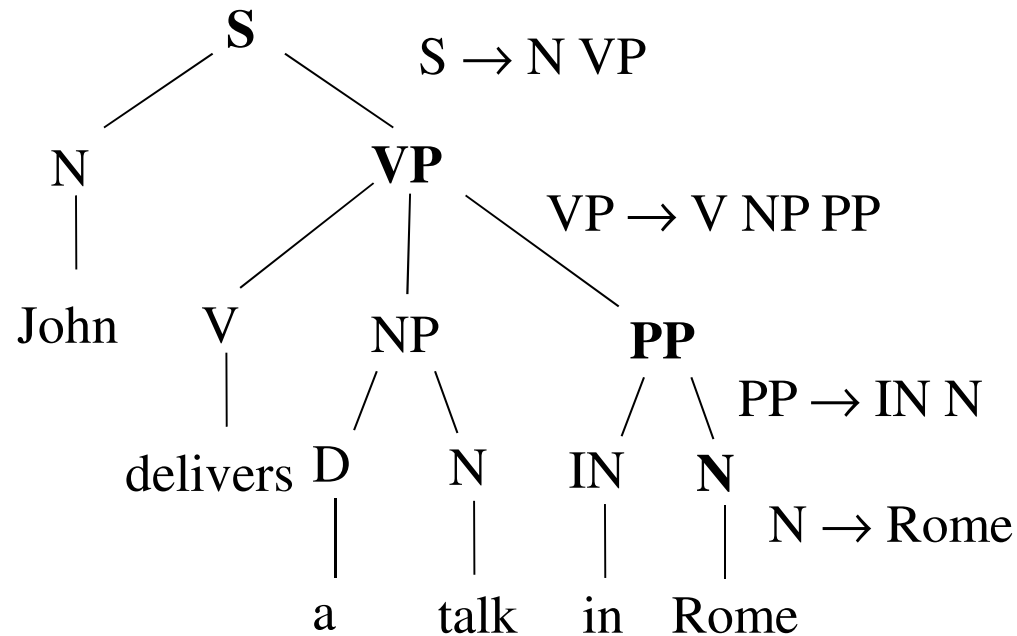
# Main Idea of Tree Kernels



# Example of a syntactic parse tree

---

- “John delivers a talk in Rome”

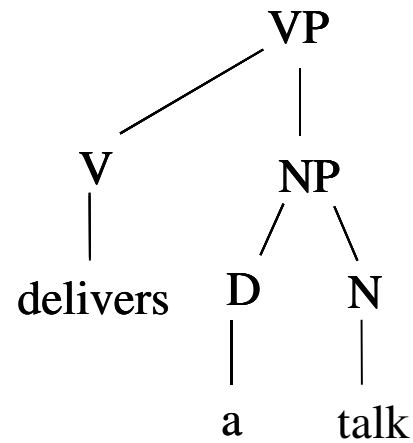




# The Syntactic Tree Kernel (STK)

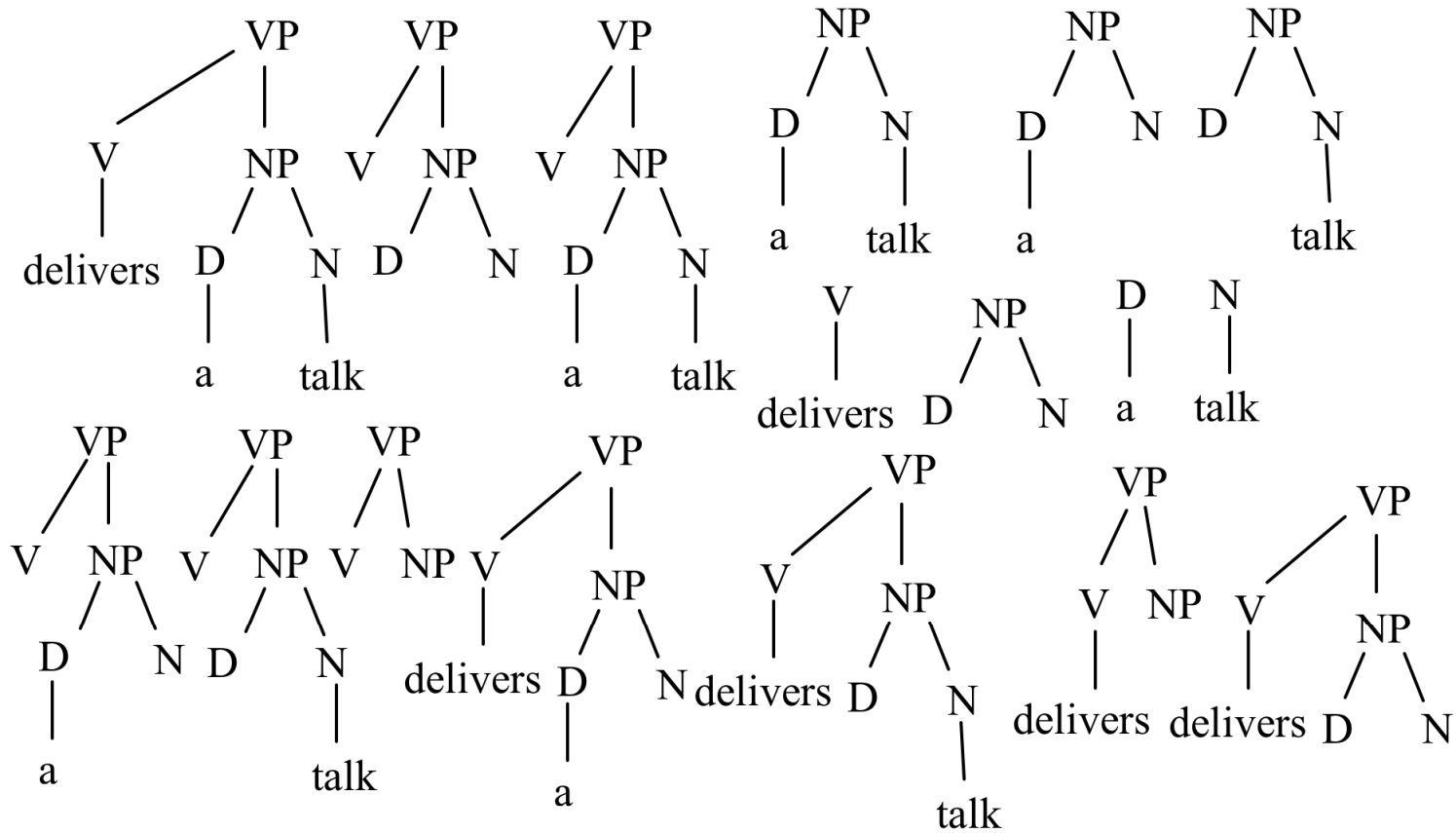
[Collins and Duffy, 2002]

---

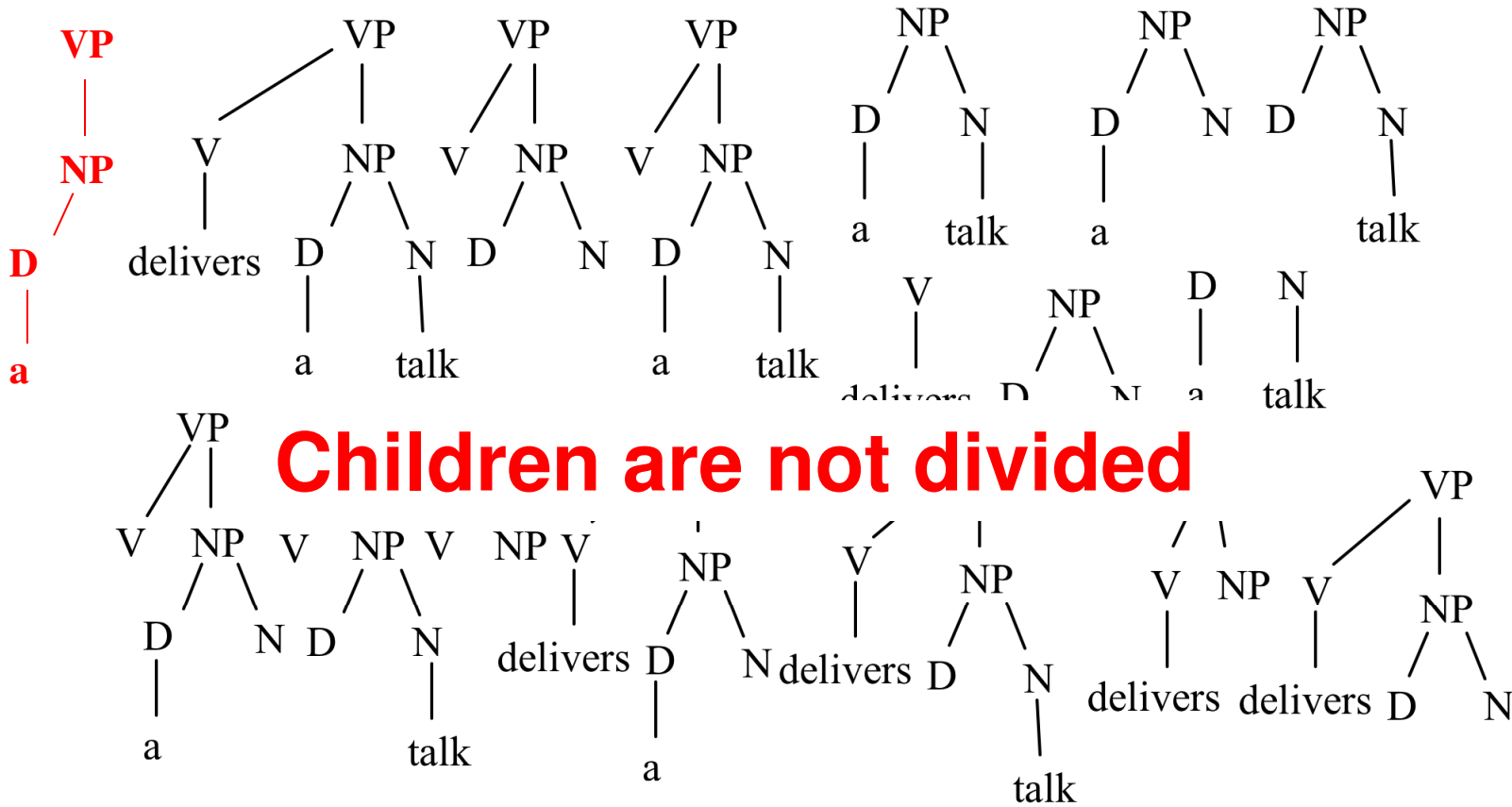


# The overall fragment set

---



# The overall fragment set

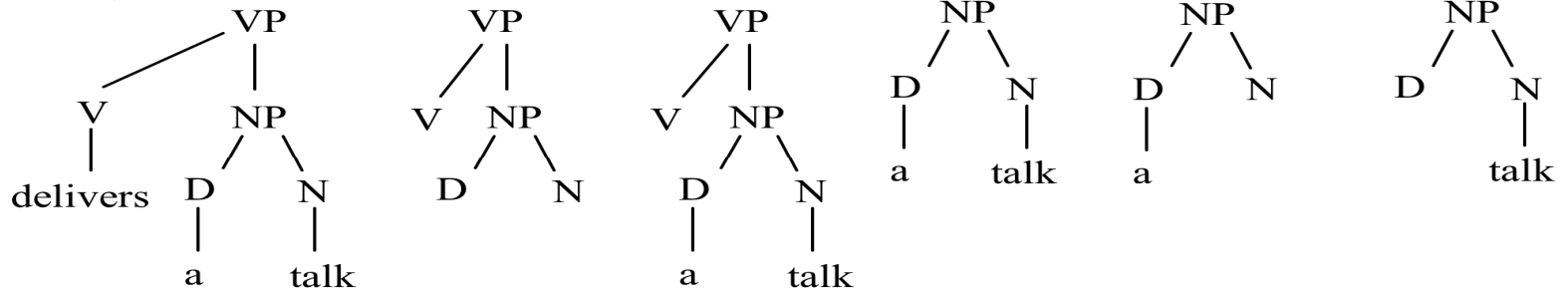


**Children are not divided**

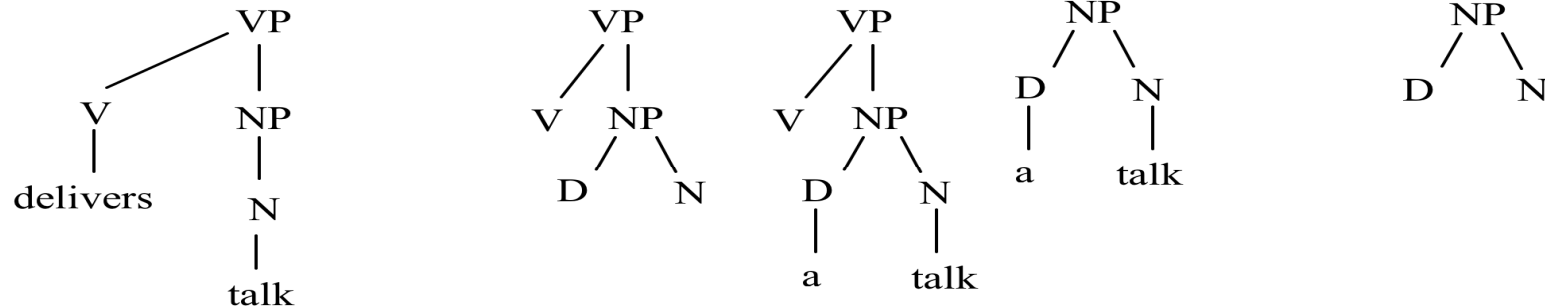


# Explicit kernel space

$$\phi(T_x) = \vec{x} = (0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0)$$



$$\phi(T_z) = \vec{z} = (1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0)$$



- $\vec{x} \cdot \vec{z}$  counts the number of common substructures



# Efficient evaluation of the scalar product

---

$$\begin{aligned}\vec{x} \cdot \vec{z} &= \phi(T_x) \cdot \phi(T_z) = K(T_x, T_z) = \\ &= \sum_{n_x \in T_x} \sum_{n_z \in T_z} \Delta(n_x, n_z)\end{aligned}$$



# Efficient evaluation of the scalar product

---

$$\begin{aligned}\vec{x} \cdot \vec{z} &= \phi(T_x) \cdot \phi(T_z) = K(T_x, T_z) = \\ &= \sum_{n_x \in T_x} \sum_{n_z \in T_z} \Delta(n_x, n_z)\end{aligned}$$

- [Collins and Duffy, ACL 2002] evaluate  $\Delta$  in  $O(n^2)$ :

$\Delta(n_x, n_z) = 0$ , if the productions are different else

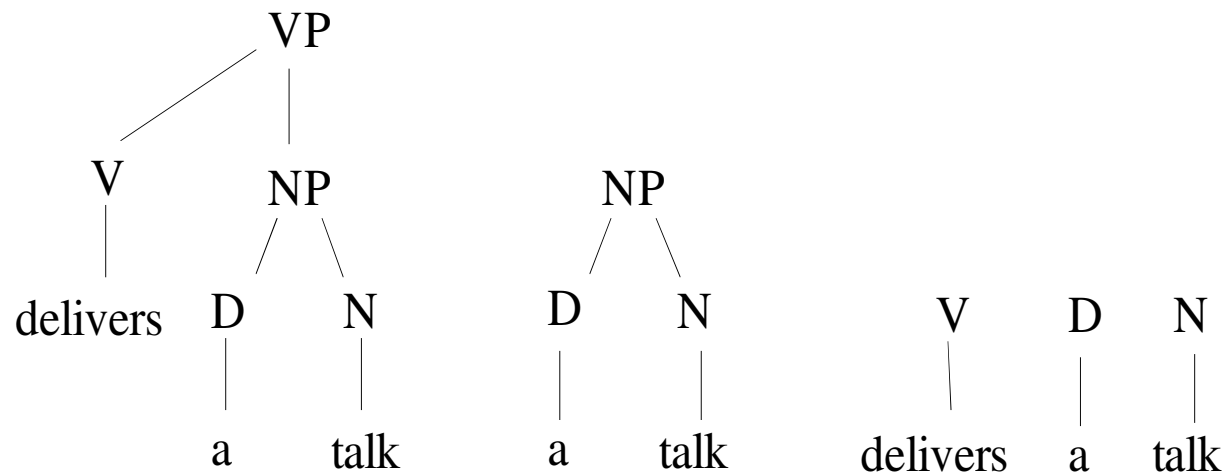
$\Delta(n_x, n_z) = 1$ , if pre-terminals else

$$\Delta(n_x, n_z) = \prod_{j=1}^{nc(n_x)} (1 + \Delta(ch(n_x, j), ch(n_z, j)))$$



# SubTree (ST) Kernel [Vishwanathan and Smola, 2002]

---



# Evaluation

---

- Given the equation for STK

$\Delta(n_x, n_z) = 0$ , if the productions are different else

$\Delta(n_x, n_z) = 1$ , if pre-terminals else

$$\Delta(n_x, n_z) = \prod_{j=1}^{nc(n_x)} (1 + \Delta(ch(n_x, j), ch(n_z, j)))$$





# SVM-light-TK Software

---

- Encodes ST, STK and combination kernels in SVM-light [Joachims, 1999]
- Available at <http://dit.unitn.it/~moschitt/>
- Tree forests, vector sets
- The new SVM-Light-TK toolkit will be released asap (email me to have the current version)



# Practical Example on Question Classification

---

- **Definition:** What does HTML stand for?
- **Description:** What's the final line in the Edgar Allan Poe poem "The Raven"?
- **Entity:** What foods can cause allergic reaction in people?
- **Human:** Who won the Nobel Peace Prize in 1992?
- **Location:** Where is the Statue of Liberty?
- **Manner:** How did Bob Marley die?
- **Numeric:** When was Martin Luther King Jr. born?
- **Organization:** What company makes Bentley cars?



# Conclusions

---

- Dealing with noisy and errors of NLP modules require robust approaches
  - SVMs are robust to noise and Kernel methods allows for:
    - Syntactic information via STK
    - Shallow Semantic Information via PTK
    - Word/POS sequences via String Kernels
  - When the IR task is complex, syntax and semantics are essential
- ⇒ Great improvement in Q/A classification
- SVM-Light-TK: an efficient tool to use them



# SVM-light-TK Software

---

- Encodes ST, SST and combination kernels in SVM-light [Joachims, 1999]
- Available at <http://dit.unitn.it/~moschitt/>
- Tree forests, vector sets
- New extensions: the PT kernel will be released asap



# References

---

- Alessandro Moschitti, Silvia Quarteroni, Roberto Basili and Suresh Manandhar, *Exploiting Syntactic and Shallow Semantic Kernels for Question/Answer Classification*, Proceedings of the 45th Conference of the Association for Computational Linguistics (ACL), Prague, June 2007.
- Alessandro Moschitti and Fabio Massimo Zanzotto, *Fast and Effective Kernels for Relational Learning from Texts*, Proceedings of The 24th Annual International Conference on Machine Learning (ICML 2007), Corvallis, OR, USA.
- Daniele Pighin, Alessandro Moschitti and Roberto Basili, *RTV: Tree Kernels for Thematic Role Classification*, Proceedings of the 4th International Workshop on Semantic Evaluation (SemEval-4), English Semantic Labeling, Prague, June 2007.
- Stephan Bloehdorn and Alessandro Moschitti, *Combined Syntactic and Semantic Kernels for Text Classification*, to appear in the 29th European Conference on Information Retrieval (ECIR), April 2007, Rome, Italy.
- Fabio Aioli, Giovanni Da San Martino, Alessandro Sperduti, and Alessandro Moschitti, *Efficient Kernel-based Learning for Trees*, to appear in the IEEE Symposium on Computational Intelligence and Data Mining (CIDM), Honolulu, Hawaii, 2007



# An introductory book on SVMs, Kernel methods and Text Categorization

---

