Kernel Functions for Structured Data

Dr. Konrad Rieck

Technische Universität Berlin

May 23, 2011

Konrad Rieck Kernel Functions for Structured Data

イロト イロト イヨト イヨト

Outline

Brief Review: Kernels

Definition and Properties

Kernels for Strings

Generic String Kernel Bag-of-words, N-grams and Substrings Efficient Implementation

Kernels for Trees

Parse Tree Kernel Efficient Implementation

・ 同 ト ・ ヨ ト ・ ヨ ト

Structured Data

Structured data ubiquituous in applied sciences

- Bioinfomatics
 e.g. DNA and protein sequences
- Natural language processing e.g. text documents and parse trees
- Computer security e.g. network traffic and program behavior
- Chemoinformatics
 e.g. molecule structures and relations



イロト イポト イヨト イヨト

Structured data \neq vectors \Rightarrow No machine learning?

Brief review: Kernels

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・

æ

What is a Kernel?

Kernel function or short kernel:

- A positive semi-definite function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- Similarity measure for objects in a domain \mathcal{X}
- Basic building block of many learning algorithms

Definition

A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a *kernel* iff k is symmetric and positive semi-definite for any subset $\{x_1, \ldots, x_l\} \subset \mathcal{X}$, that is

$$\sum_{i,j=1}^m c_i c_j k(x_i,x_j) \geq 0 \text{ with } c_1,\ldots,c_m \in \mathbb{R}.$$

イロト イポト イヨト イヨト

э

Kernels and Feature Spaces

Theorem

A kernel k induces a feature map $\psi : \mathcal{X} \to \mathcal{F}$ to a Hilbert space, where k equals an inner product. That is, for all $x, y \in \mathcal{X}$

 $k(\mathbf{x},\mathbf{y}) = \langle \psi(\mathbf{x}), \psi(\mathbf{y}) \rangle \,.$

Interface to geometry in feature space

Access to inner products, vector norms and distances, e.g.,

$$||\psi(x)||_{2} = \sqrt{k(x,x)}$$
$$||\psi(x) - \psi(y)||_{2} = \sqrt{k(x,x) + k(y,y) - 2k(x,y)}$$

Classic Kernels

Let $\mathcal{X} \subseteq \mathbb{R}^d$. Then kernels $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ are given by

•
$$k(x,y) := \langle x, y \rangle = \sum_{i=1}^{d} x_i y_i$$
 (Linear kernel)

►
$$k(x,y) := (\langle x,y \rangle + \theta)^p$$
 (Polynomial kernel)

►
$$k(x,y) := \exp\left(\frac{||x-y||^2}{\gamma}\right)$$
 (Gaussian kernel)

►
$$k(x,y) := tanh(\langle x,y \rangle + \theta)$$
 (Sigmoidal kernel)

However: Domain \mathcal{X} not restricted to vectorial data!

イロト イポト イヨト イヨト

Kernels and Structured Data

Kernels for structured data

- Definition of kernel k over non-vectorial domain \mathcal{X}
- ► Any *k* valid, if symmetric and positive semi-definite
- Integration with kernel-based learning methods



Kernels for Strings

・ロト ・ 同ト ・ ヨト ・ ヨト

Strings

Alphabet

An alphabet \mathcal{A} is a finite set of discrete symbols

► DNA,
$$\mathcal{A} = \{A, C, G, T\}$$

► Natural language text, $A = \{a, b, c, ..., A, B, C, ...\}$

String or Sequence

A string x is concatenation of symbols from A

- \mathcal{A}^n = all strings of length *n*
- A^* = all strings of arbitary length
- |x| = length of a string

イロト イポト イヨト イヨト

Embedding Strings

Mapping of strings to a feature space

- Characterize strings using a language $L \subseteq A^*$.
- Feature space spanned by occurences of words $w \in L$

Feature map

A function $\phi:\mathcal{A}^*\to\mathbb{R}^{|L|}$ mapping strings to $\mathbb{R}^{|L|}$ given by

$$\phi: \mathbf{x} \longmapsto \left(\#_{w}(\mathbf{x}) \cdot \sqrt{N_{w}} \right)_{w \in \mathbb{N}}$$

where $\#_w(x)$ returns the occurences of w in string x and N_w is a weighting of individual words.

イロト イポト イヨト イヨト

String Kernels

Generic String Kernel

A generic string kernel $k : \mathcal{A}^* \times \mathcal{A}^* \to \mathbb{R}$ is given by

$$k(\mathbf{x},\mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \sum_{w \in L} \#_w(\mathbf{x}) \cdot \#_w(\mathbf{z}) \cdot N_w$$

Proof.

By definition k is an inner product in $\mathbb{R}^{|L|}$ and thus symmetric and positive semi-definite.

イロト イポト イヨト イヨト

э

Bag-of-Words

Characterization of strings using non-overlapping substrings

Bag-of-Words Kernel

String kernel using embedding language of words with delimiters D

$$k(x,y) = \sum_{w \in L} \#_w(x) \cdot \#_w(y) \cdot N_w$$
 with $L = (\mathcal{A} \setminus D)^*$

 Suitable for analysis of strings with known structure e.g., natural language text, tokenized data, log files

イロト イポト イヨト イヨト

N-grams

Characterization of strings using substrings of length n



N-gram Kernel

String kernel using embedding language of *n*-grams:

$$k(x,y) = \sum_{w \in L} \#_w(x) \cdot \#_w(y) \cdot N_w$$
 with $L = \mathcal{A}^n$

 Suitable for analysis of strings with unknown structure, e.g., DNA sequences, network attacks, binary data

All Substrings

Characterization of strings using all possible substrings



All-Substring Kernel

String kernel using embedding language of all strings:

$$k(x,y) = \sum_{w \in L} \#_w(x) \cdot \#_w(y) \cdot N_w$$
 with $L = \mathcal{A}^*$

- Suitable for analysis of generic string data
- Encoding of prior knowledge in weighting N_w

Implementing String Kernels

Efficient computation of string kernel k(x, z)

- Feature space high-dimensional but sparsely populated
- Sufficient to consider only w with $\#_w(x) \neq 0$ and $\#_w(z) \neq 0$
- Application of special data structures for strings

Implementation strategies

- **1.** Explicit but sparse representation of feature vectors, \longrightarrow hash tables, tries and sorted arrays
- 2. Implicit representation of feature vectors,
 - \longrightarrow suffix trees and arrays

イロト イポト イヨト イヨト

Sorted Arrays

Example: strings *x* and *y* with embedding language $L = A^3$

Extracted words w stored with $\#_w(x)$ in sorted array

$$\phi(\mathbf{x}) \leftarrow aa|1 - ab|1 - ba|1 - bb|1 \rightarrow \phi(\mathbf{z}) \leftarrow aa|3 - ab|1 - ba|1 - ba|1 \rightarrow ba|1 - b$$

- ▶ Explicit kernel computation → parallel loop over arrays
- ▶ Run-time O(|x| + |z|) for words with no or bounded overlap

イロト イポト イヨト イヨト

Suffix Trees

Strings jointly stored in generalized suffix tree



- ► Implicit kernel computation → depth first traversal
- ▶ Run-time O(|x| + |z|) for arbitrary embedding languages

Kernels for Trees

イロト イワト イヨト イヨト

Trees and Parse Trees

Tree

A tree $x = (V, E, v^*)$ is an acyclic graph (V, E) rooted at $v^* \in V$.

Parse tree

A tree x deriving from agrammar, such that each node $v \in V$ is associated with a production rule p(v).

Further notation

- $v_i = i$ -th child of node $v \in V$
- |v| = number of children of $v \in V$
- T = set of all possible parse trees

イロト 不得 トイヨト イヨト

Parse Trees

Tree representation of "sentences" derived from a grammar



Parse tree for "mary ate lamb" with production rules

•
$$p_1 : A \longrightarrow B$$

▶
$$p_2: B \longrightarrow$$
 "mary" "ate" C

イロト イポト イヨト イヨト

з.

▶
$$p_3 : C \longrightarrow$$
 "lamb"

Common data structure in several application domains, e.g., natural language processing, compiler design, ...

Embedding Trees

Characterization of parse trees using contained subtrees



Feature map

A function $\phi : T \to \mathbb{R}^{|T|}$ mapping trees to $\mathbb{R}^{|T|}$ given by

$$\phi: \mathbf{x} \longmapsto (\#_t(\mathbf{x}))_{t \in T}$$

where $\#_t(x)$ returns the occurences of subtree *t* in *x*.

Parse Tree Kernel

Parse Tree Kernel

A tree kernel $k : T \times T \to \mathbb{R}$ is given by

$$k(\mathbf{x},\mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \sum_{t \in T} \#_t(\mathbf{x}) \cdot \#_t(\mathbf{z})$$

Proof.

By definition k is an inner product in the space of all trees T and thus symmetric and positive semi-definite.

Counting shared subtrees

Parse tree kernel and counting

- Parse tree kernel counts the number of shared subtrees
- For each pair (v, w) determine shared subtrees at v and w.

$$k(\mathbf{x}, \mathbf{z}) = \sum_{t \in T} \#_t(\mathbf{x}) \cdot \#_t(\mathbf{z}) = \sum_{\mathbf{v} \in \mathsf{V}_x} \sum_{\mathbf{w} \in \mathsf{V}_z} c(\mathbf{v}, \mathbf{w})$$

Counting function

- c(v, w) = 0 if $p(v) \neq p(w)$ (different production)
- ► c(v, w) = 1 if |v| = |w| = 0 (leaf nodes)

otherwise

$$c(v,w) = \prod_{i=1}^{|v|} (1 + c(v_i,w_i))$$

Counting in Detail

- ► First base case: c(v, w) = 0 if $p(v) \neq p(w)$ ⇒ trivial, no match = no shared subtrees
- ► Second base case: c(v, w) = 1 if |v| = |w| = 0⇒ trivial, one leave = one subtree
- Recursion: $c(v, w) = \prod_{i=1}^{|v|} (1 + c(v_i, w_i))$



Count all combinations of shared subtrees below node A

$$c(v_A, w_A) = (n+1) \cdot (m+1)$$

Implementation of Tree Kernels

Efficient implementation using dynamic programming

- Explicit feature vector representations intractable
- Implicit kernel computation by counting shared subtrees



Matrix of counts c(v, w) for all shared subtrees sorted by height

- Count small subtrees first
- Gradually aggregate counts

イロト イポト イヨト イヨト

э

Run-time $\mathcal{O}(|V_x| \cdot |V_z|)$.

Conclusions

Kernels for strings and trees

- Effective means for learning with structured data
- Several efficient kernels and implementations

More interesting kernels for structured data

- Kernel for graphs, images, sounds, ...
- Convolution kernels, approximate kernels, ...

Interesting applications (upcoming lectures)

- "Catching hackers": Network intrusion detection
- "Discovering genes": Analysis of DNA sequences

References



Shawe-Taylor, J. and Cristianini, N. (2004). *Kernel methods for pattern analysis*. Cambridge University Press.

イロト イポト イヨト イヨト

з.