

# Kernel Methods and Automata

Anton Xue

## 1 Introduction

In this sketch we study the relation between graph kernels and finite state machines. We study previous work in algebraic formulation for automata theory and kernel methods, with an interest in how these techniques may be extended to formal methods.

## 2 Preliminaries

### 2.1 Algebraic Automata Theory

Previous work in formulation in algebraic foundations for automata theory exist in literature [4], and much of our notation is taken from [2]. We overload notation for addition (+) and multiplication ( $\cdot$ ) in algebraic structures when possible to avoid clutter. Similarly, the additive identity (0) and multiplicative identity (1) are also overloaded when possible.

**Definition 1** (Monoid). *A monoid is an algebraic structure  $(\mathbb{K}, \cdot, 1)$  where:*

- $\mathbb{K}$  is closed under monoid multiplication  $\cdot : \mathbb{K} \times \mathbb{K} \rightarrow \mathbb{K}$ .
- 1 is the multiplicative identity.

*When possible, we elide the  $\cdot$  in monoid multiplication to write  $ab$  instead of  $a \cdot b$ . When multiplication  $\cdot$  is commutative, the system is known as a commutative monoid.*

**Example 1.** *Consider a finite alphabet  $\Sigma$ , then the system  $(\Sigma^*, \cdot, \varepsilon)$  forms a monoid, where monoid multiplication  $\cdot$  is string concatenation and  $\varepsilon$  is the empty string.*

**Definition 2** (Semiring). *A semiring  $(\mathbb{K}, +, \cdot, 0, 1)$  is a system where:*

- Semiring addition is a commutative monoid  $(\mathbb{K}, +, 0)$ .
- Semiring multiplication  $(\mathbb{K}, \cdot, 1)$  is a monoid.

- 0 annihilates semiring multiplication.

**Example 2.** For a finite alphabet  $\Sigma$ , the system  $(2^{\Sigma^*}, \cup, \cdot, \emptyset, \varepsilon)$  forms a semiring. Here semiring addition is set union, and semiring multiplication is defined as:

$$A \cdot B = \{a \cdot b : a \in A, b \in B\}$$

for all  $A, B \in 2^{\Sigma^*}$ . We take the additive identity to be the empty set  $\emptyset$ , and the multiplicative identity as the empty string  $\varepsilon$ .

A weighted finite state transducer (WFST) is a very general transition system defined using a semiring to specify transition weights.

**Definition 3** (Weighted Finite State Transducer). A weighted finite state transducer over a semiring  $\mathbb{K}$  is a system  $(\Sigma_I, \Sigma_O, Q, I, F, \Delta, \lambda, \rho)$  where:

- $\Sigma_I$  is a finite input alphabet.
- $\Sigma_O$  is a finite output alphabet.
- $Q$  is a finite set of states.
- $I \subseteq Q$  is the set of starting states.
- $F \subseteq Q$  is the set of final states.
- $\Delta \subseteq Q \times (\Sigma_I \cup \{\varepsilon\}) \times \mathbb{K} \times (\Sigma_O \cup \{\varepsilon\}) \times Q$  is the transition function weighted by  $\mathbb{K}$ .
- $\lambda : I \rightarrow \mathbb{K}$  is the initial state weight function.
- $\rho : F \rightarrow \mathbb{K}$  is the final state weight function.

Note that the transition function  $\Delta$  is defined to permit non-deterministic behavior by default.

It should be noted that  $\Sigma_I$  and  $\Sigma_O$  can be seen as the generators of the free monoids  $\Sigma_I^*$  and  $\Sigma_O^*$ , which represents the set of all strings over  $\Sigma_I$  and  $\Sigma_O$  respectively.

Although we are not yet interested in the full generality that a WFST offers, it is still nice to see what is available to us in terms of abstraction.

A special case of WFSTs that we are interested in are non-deterministic finite automata, whose specification in terms of WSFTs is by Cortes [2]. We introduce a simplified structure.

**Definition 4** (Non-deterministic Finite Automata). A non-deterministic finite automata is a system  $(\Sigma, Q, \Delta, I, F)$  where:

1.  $\Sigma$  is a finite alphabet.

2.  $Q$  is a finite set of states.
3.  $\Delta \subseteq Q \times \Sigma \times Q$  is the transition function.
4.  $I \subseteq Q$  is the set of initial states.
5.  $F \subseteq Q$  is the set of final states.

We do not permit  $\varepsilon$ -transitions in our definition, which can be eliminated anyways [5].

## 2.2 Reproducing Kernel Hilbert Spaces

Inner products spaces are nice because they allow us to measure projection. Additionally, an inner product induces a norm, from which a metric, and thus topology, can also be created.

**Definition 5** (Hilbert Space). *A Hilbert space is a inner product space that is complete with respect to the metric induced by its inner product.*

If  $\mathcal{H}$  is a Hilbert space, write  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  to denote its inner product. We drop the subscript when context is clear, and remark that another option is for complex-valued inner products.

A special type of Hilbert spaces are known as reproducing kernel Hilbert spaces (RKHS) [1]. In short, these are Hilbert spaces with a special function known as the reproducing kernel.

**Definition 6** (Reproducing Kernel). *Let  $\mathcal{H}$  be a Hilbert space. A reproducing kernel is a function  $k : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  that satisfies the property*

$$f(x) = \langle f, k(x, \cdot) \rangle_H$$

for all  $f \in \mathcal{H}$ .

A RKHS can be induced by the existence of a kernel function (not to be confused with the reproducing kernel). This type of kernel function can be defined over general sets, and provides a way of taking inner products by embedding them into a Hilbert space.

**Definition 7** (Kernel). *Let  $\mathcal{X}$  be a set. A function  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  if:*

1.  $k(x, y) = k(y, x)$  for all  $x, y \in \mathcal{X}$ . This is known as symmetric.
2. If for all  $x_1, x_2, \dots, x_n \in \mathcal{X}$  the Gram matrix  $K$  defined by:

$$K_{i,j} = k(x_i, x_j)$$

is positive definite.

Given a set  $\mathcal{X}$  and a kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , one can generate a Hilbert space. Define a “feature map”  $\phi : \mathcal{X} \rightarrow \mathbb{R}^{\mathcal{X}}$  as follows:

$$\phi(x) = k(x, \cdot)$$

Then  $\phi$  maps each element of  $\mathcal{X}$  into a Hilbert space  $\mathcal{H}$  consisting of the closure of the span of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$ :

$$\mathcal{H} = \overline{\text{span} \{f : \mathcal{X} \rightarrow \mathbb{R}\}} = \overline{\left\{ \sum_{i=1}^n a_i k(\cdot, x_i) : n \in \mathbb{N}, x_i \in \mathcal{X}, a_i \in \mathbb{R} \right\}}$$

Consider two functions  $f, g \in \mathcal{H}$  which would have form:

$$f(x) = \sum_{i \in I} a_i k(x, u_i) \quad g(x) = \sum_{j \in J} b_j k(x, v_j)$$

where each  $u_i, v_j \in \mathcal{X}$ . Note that the summation may be countably infinite over the index sets  $I$  and  $J$ , since Hilbert spaces are complete with respect to the norm induced by the inner product. The inner product between  $f$  and  $g$  is then defined as:

$$\langle f, g \rangle = \left\langle \sum_{i \in I} a_i k(\cdot, u_i), \sum_{j \in J} b_j k(\cdot, v_j) \right\rangle = \sum_{i \in I} \sum_{j \in J} a_i b_j k(u_i, v_j)$$

This is further described and proven in previous work [2, 6].

## 2.3 Tensor Products

Tensor algebra are treated in more detail in other texts [3], and multiple generalizations exist. In this sketch we are interested in a very specialized slice of tensor algebra. Note that we use the term linear space instead of vector space, since much of linear algebra, including tensor algebra, does not necessarily have to be defined over fields.

**Definition 8** (Tensor Product of Linear Spaces). *Let  $V$  be a linear space with countable basis  $\{v_1, v_2, \dots\}$ , and  $W$  be a linear space with countable basis  $\{w_1, w_2, \dots\}$ . The tensor product  $V \otimes W$  is the linear space spanned by the countable basis  $\{v_i \otimes w_j\}$ .*

Note that we treat  $v_i \otimes w_j$  as symbols without special meaning attributed to them. For instance, if  $V$  has basis  $\{v_1, v_2\}$  and  $W$  has basis  $\{w_1, w_2, w_3\}$ , then  $V \otimes W$  has basis:

$$\{v_1 \otimes w_1, v_1 \otimes w_2, v_1 \otimes w_3, v_2 \otimes w_1, v_2 \otimes w_2, v_2 \otimes w_3\}$$

Note that the above basis written in this manner can be considered an ordered basis of  $V \otimes W$ .

Additionally, if  $S : V \rightarrow W$  and  $T : X \rightarrow Y$  are linear transformations, then the tensor product would have type:

$$S \otimes T : V \otimes X \rightarrow W \otimes Y$$

Finite-dimensional linear transformations are often envisioned as matrices that act on column vectors by left multiplication. For two matrices  $A_{m \times n}$  and  $B_{p \times q}$  that represent linear transforms  $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $B : \mathbb{R}^q \rightarrow \mathbb{R}^p$  respectively, the Kronecker product is a way to lift the linear transform represented by  $A$  and  $B$  into the tensor product space. We write this in terms of block matrices as follows:

$$A \otimes B = \begin{bmatrix} A_{1,1}B & \dots & A_{1,n}B \\ \vdots & \ddots & \vdots \\ A_{m,1}B & \dots & A_{m,n}B \end{bmatrix}$$

This is then a linear transformation with respect to the ordered basis  $\mathbb{R}^n \otimes \mathbb{R}^q$  and  $\mathbb{R}^m \otimes \mathbb{R}^p$ .

The Kronecker product on matrices is also used as a tensor product over adjacency matrices of graphs to define the notion of a product graph.

## 3 Embedding Automata

### 3.1 Simple Graphs

We first consider kernel methods over graphs, which is the main focus of Vishwanathan [6]. The main idea is that the space of simple graphs can be treated using kernel methods described above. Here, simple graphs are taken to mean (un)-directed graphs that are representable by 0/1 adjacency matrices.

Let  $G = (V_G, E_G)$  and  $H = (V_H, E_H)$  be two finite simple graphs. Here  $V_G$  is the vertex set of  $G$ , and  $E_H \subseteq V_G \times V_G$  is the edge set of  $G$ , with similar notation for  $H$ . Use the same letters to represent their adjacency matrices in an abuse of notation. Then a kernel on this space of graph adjacency matrices can be defined as follows:

$$k(G, H) = \sum_{k=0}^{\infty} \mu(k) q^\top (G \otimes H)^k p$$

where  $\mu : \mathbb{N} \rightarrow \mathbb{R}$  is a function over  $\mathbb{N}$  to help ensure that the summation converges,  $q$  and  $p$  are fixed and of the appropriate dimensions intended to mean something along the lines of final and initial weights with respect to a basis on  $V_G \otimes H_G$ .

Note that adjacency matrices may either act on column vectors by left multiplication or row vectors by right multiplication, whichever is convenient.

Suppose that  $G$  and  $H$  are known, then a normalized version of the kernel can be defined as:

$$k(G, H) = \frac{1}{|V_G|^2 |V_H|^2} \sum_{k=0}^{\infty} \mu(k) q^\top (G \otimes H)^k p$$

Since  $G$  and  $H$  represent adjacency matrices respectively, if  $q$  and  $p$  take on values in  $[0, 1]$ , then each iteration of the summation does not exceed the dimension of  $G \otimes H$ , which is  $|V_G|^2 |V_H|^2$ . In

addition, if  $\mu$  is a probability distribution on  $\mathbb{N}$ , then:

$$0 \leq k(G, H) \leq 1$$

for all simple graphs  $G$  and  $H$ .

### 3.2 Non-deterministic Finite Automata

A similar trick can be used for NFAs. However a few considerations have to be made when representing NFAs as matrices. In particular, because semiring multiplication is not commutative, we must be extra careful in our construction.

Consider an NFA  $(\Sigma, Q, \Delta, S, F)$ . We want to represent this as a matrix  $A_{|Q| \times |Q|}$  where  $A_{i,j}$  denotes the transition from state  $q_i$  to state  $q_j$ . We arbitrarily choose to read strings left to right, so this implies that  $A$  should act by right multiplication on row vectors.

Additionally, we must account for the fact that edges in an NFA may be weighted by multiple elements (or strings, in a more general setting) from the alphabet. This suggests that each entry of  $A$  should instead be a set:

$$A_{i,j} = \{\sigma : (q_i, \sigma, q_j) \in \Delta\}$$

In fact, it is possible to view each entry of the  $A$  matrix as a member of the  $(2^{\Sigma^*}, \cup, \cdot, \emptyset, \varepsilon)$  semiring, which was discussed earlier. Note that technically this would technically permit each transition to consume a string of letters rather than a single one, but we could just be careful in our construction to avoid these cases that would make analysis efforts more difficult. All together,  $A$  would then be typed as follows:

$$A : (\Sigma^*)^{|Q|} \rightarrow (\Sigma^*)^{|Q|}$$

where the semiring operations and identities are implicit. For each  $1 \leq i, j \leq |Q|$ , the entry  $A_{i,j}$  is still defined as before, and this  $A$  matrix acts on row vectors by right multiplication.

Additionally, take  $\alpha$  to be a  $|Q|$ -sized vector representing the initial states, where:

$$\alpha_i = \begin{cases} \{\varepsilon\} & q_i \in I \\ \emptyset & \text{otherwise} \end{cases}$$

Define  $\beta$  to be the vector representing the final states in a similar manner.

We will sometimes overload notation to have  $A$  mean both the transition matrix and the NFA that it represents. Inspired by the normalized kernel definition for simple graphs above, this suggests that we can extend it to NFA as well.

**Definition 9** (Normalized NFA Kernel). *Given two non-deterministic finite automata:*

$$A = (\Sigma_A, Q_A, \Delta_A, I_A, F_A) \qquad B = (\Sigma_B, Q_B, \Delta_B, I_B, F_B)$$

The normalized NFA kernel is defined as:

$$k(A, B) = \frac{1}{|Q_A|^2 |Q_B|^2} \sum_{k=0}^{\infty} \frac{\mu(k)}{|\Sigma_A|^k |\Sigma_B|^k} \left| (\alpha_A \otimes \alpha_B)^\top (A \otimes B)^k (\beta_A \otimes \beta_B) \right|$$

where  $A \otimes B$  is the Kronecker product of the transition matrices of  $A$  and  $B$  over the  $2^{(\Sigma_A \cup \Sigma_B)^*}$  semiring,  $\alpha_A \otimes \alpha_B$  is the Kronecker product of the initial state vector of  $A$  and  $B$ ,  $\beta_A \otimes \beta_B$  is the Kronecker product of the final state vector of  $A$  and  $B$ , and  $\mu$  is a probability distribution on  $\mathbb{N}$ .

**Theorem 1.** *The normalized NFA kernel induces a RKHS.*

*Proof.* It suffices to prove that  $k$  is symmetric and positive definite.

We first show symmetry. Observe that exchanging the order of the Kronecker product on matrices only serves to change the order of the bases that the new linear transformation acts on. Consider two automata given as:

$$A = (\Sigma_A, Q_A, \Delta_A, I_A, F_A) \quad B = (\Sigma_B, Q_B, \Delta_B, I_B, F_B)$$

where we use  $A$  and  $B$  to also denote the transition matrix with each element taken from the  $2^{(\Sigma_A \cup \Sigma_B)^*}$  semiring, with  $\alpha_A$  and  $\alpha_B$  for the initial states and  $\beta_A$  and  $\beta_B$  for the final states as defined earlier. Since the semiring operation here is union, this implies that:

$$(\alpha_A \otimes \alpha_B)^\top (A \otimes B)^k (\beta_A \otimes \beta_B) = (\alpha_B \otimes \alpha_A)^\top (B \otimes A)^k (\beta_B \otimes \beta_A)$$

for all  $k$ . The immediate consequence is that  $k$  is symmetric.

Now take  $A$  as before and consider the following:

$$(\alpha_A \otimes \alpha_A)^\top (A \otimes A)^k (\beta_A \otimes \beta_A)$$

This will only be the empty set if  $A$  lacks strings of length  $k$ . If it is the case that  $k(A, A) = 0$ , then the implication is that  $A$  accepts no strings. In this case,  $\phi(A)$  is taken to mean the zero element of  $\mathcal{H}$ , and so we conclude that  $k$  is positive definite.

□

**Theorem 2.** *The normalized NFA kernel is normalized in the sense that:*

$$0 \leq k(A, B) \leq 1$$

for all NFA  $A$  and  $B$ .

*Proof.* Given two automata as:

$$A = (\Sigma_A, Q_A, \Delta_A, I_A, F_A) \quad B = (\Sigma_B, Q_B, \Delta_B, I_B, F_B)$$

First we have that  $k$  is always non-negative by inspection.

Our construction forces the adjacency matrices  $A$  and  $B$  to have in each entry at most  $|\Sigma_A|$  and  $|\Sigma_B|$  elements respectively. Consequently, each entry of  $A^k$  and  $B^k$  will have at most  $|\Sigma_A|^k$  and  $|\Sigma_B|^k$  elements respectively, since that is the maximum number of strings of length  $k$  for each alphabet. Therefore,  $(A \otimes B)^k$  will have size  $(|Q_A| |Q_B|) \times (|Q_A| |Q_B|)$ , with at most  $|\Sigma_A|^k |\Sigma_B|^k$  elements in each entry, and so:

$$\frac{1}{|Q_A|^2 |Q_B|^2} \frac{1}{|\Sigma_A|^k |\Sigma_B|^k} \left| (\alpha_A \otimes \alpha_B)^\top (A \otimes B)^k (\beta_A \otimes \beta_B) \right| \leq 1$$

for all  $k$ . As  $\mu$  is a probability distribution on  $\mathbb{N}$ , we conclude that  $k$  is normalized.

□

## 4 Bounding Inner Products

**What do lower bounds on inner products mean??**

## References

- [1] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- [2] Corinna Cortes, Patrick Haffner, and Mehryar Mohri. “Rational kernels: Theory and algorithms”. In: *Journal of Machine Learning Research* 5.Aug (2004), pp. 1035–1062.
- [3] Mikhail Itskov. *Tensor algebra and tensor analysis for engineers*. Springer, 2007.
- [4] Werner Kuich and Arto Salomaa. *Semirings, automata, languages*. Vol. 5. Springer Science & Business Media, 2012.
- [5] John E Savage. *Models of computation*. Vol. 136. Addison-Wesley Reading, MA, 1998.
- [6] S Vichy N Vishwanathan et al. “Graph kernels”. In: *Journal of Machine Learning Research* 11.Apr (2010), pp. 1201–1242.