

# Measures on Languages

## 1 Introduction

### 1.1 Notation

Let  $\Sigma$  denote a countable alphabet. Unless otherwise specified, assume  $|\Sigma| < \infty$ .

Let  $\Sigma^*$  be the set of all finite words from  $\Sigma$ . Write words as  $w$ .

Let  $L$  denote a language, implicitly over  $\Sigma$ . In other words,  $L \subseteq \Sigma^*$ .

Let  $\mathcal{L}$  be a family of languages.

## 2 Measures

Given a family of languages  $\mathcal{L}$ , let  $\sigma(\mathcal{L})$  be the  $\sigma$ -algebra generated on  $\mathcal{L}$  satisfying the following:

(1)

$$\emptyset, \Sigma^* \in \sigma(\mathcal{L})$$

(2)

$$L \in \sigma(\mathcal{L}) \implies L^c = \Sigma^* \setminus L \in \sigma(\mathcal{L})$$

(3)

$$L_1, L_2, \dots \in \sigma(\mathcal{L}) \implies \bigcup_{k=1}^{\infty} L_k \in \sigma(\mathcal{L})$$

Then  $(\mathcal{L}, \sigma(\mathcal{L}))$  is a measurable space.

*Remark 1.* If  $\mathcal{L}$  happened to be a family of regular languages, there is no guarantee that  $\sigma(\mathcal{L})$  will still be a family of regular languages. A counter example is the following:

$$L_1 = \{ab\} \quad L_2 = \{aabb\} \quad L_3 = \{aaabbb\} \quad \dots \quad L_k = \{a^k b^k\} \quad \dots$$

But taking the countable union yields:

$$\bigcup_{k=1}^{\infty} L_k = \{a^n b^n : n \in \mathbb{Z}_{\geq 0}\}$$

Which is not regular. □

## 2.1 Defining Measures

### 2.1.1 From Non-negative Integers

We first consider the non-negative integers  $\mathbb{Z}_{\geq 0}$ . Let  $\eta$  be a  $\sigma$ -finite measure on  $\mathbb{Z}_{\geq 0}$ . The  $\sigma$ -finite conditions ensures that no strange singularities occur for any integers under consideration. We may later restrict  $\eta$  to be finite if necessary, if we want nicer conditions.

Observe that, by abuse of notation:

$$\Sigma^* = \bigcup_{k=1}^{\infty} \Sigma^k$$

In English:  $\Sigma^*$  is the union of the set (language) of finite strings of length  $k$ , denoted  $\Sigma^k$ .

Because we assumed  $|\Sigma| < \infty$ , this also means that:  $|\Sigma^k| = |\Sigma|^k$ .

Consider now some language  $L \in \sigma(\mathcal{L})$ . Also decompose  $L$  into disjoint sub-languages by length as follows:

$$L = \bigcup_{k=1}^{\infty} L_k$$

Of course,  $L_k \subseteq \Sigma^k$ .

Because we are able to precisely calculate  $|\Sigma^k|$ , one “natural” way of defining a measure  $\lambda$  on the measurable space  $(\mathcal{L}, \sigma(\mathcal{L}))$  is as follows:

$$\lambda(L) = \sum_{k=1}^{\infty} \lambda(L_k) = \sum_{k=1}^{\infty} \frac{|L_k|}{|\Sigma^k|} \eta(k)$$