# When is the Normalized Edit Distance over Non–Uniform Weights a Metric?

Dana Fisman and <u>**Ilay Tzarfati**</u>

# The Edit (Levenshtein) Distance

Types of operations:                                    Uniform weight

- Delete the letter $\sigma$                    $\begin{bmatrix} \sigma \\ \varepsilon \end{bmatrix}$     1

- Insert the letter $\sigma$                    $\begin{bmatrix} \varepsilon \\ \sigma \end{bmatrix}$     1

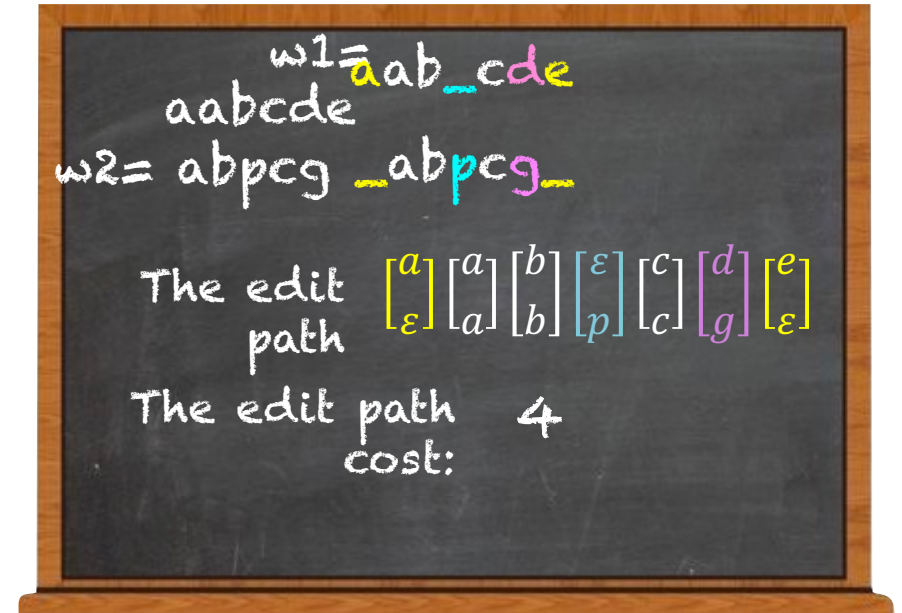- Substitute the letter $\sigma$ with $\sigma'$    $\begin{bmatrix} \sigma \\ \sigma' \end{bmatrix}$     1

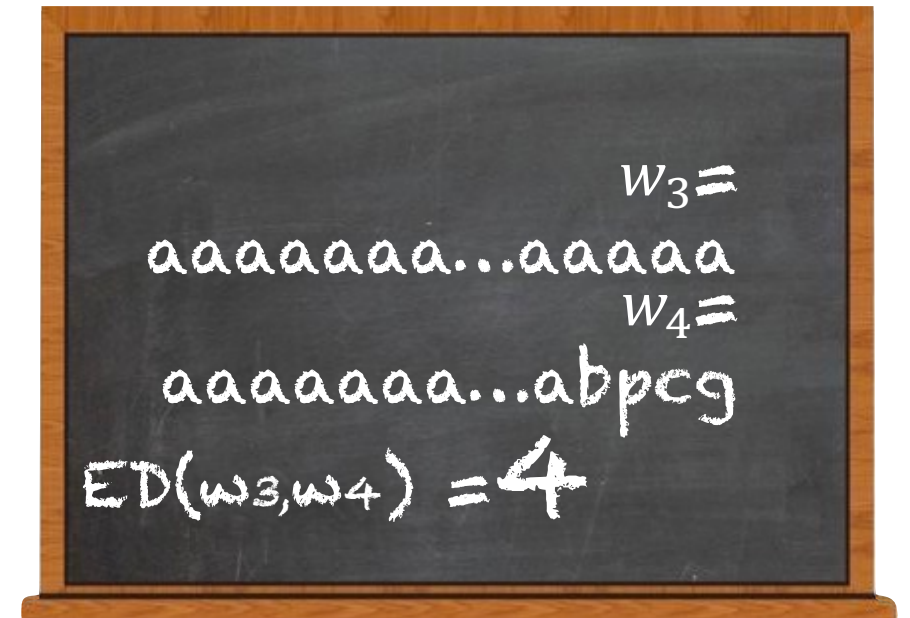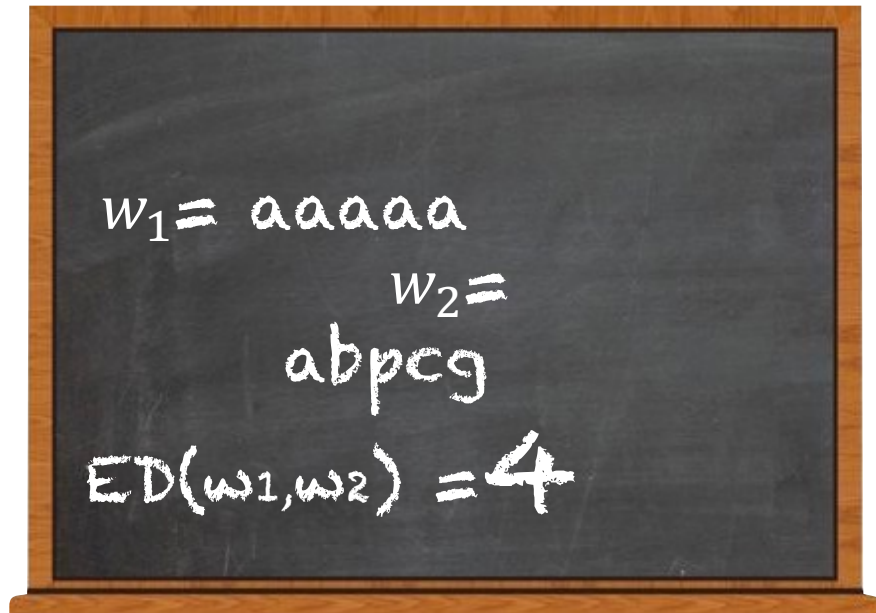- No change                              $\begin{bmatrix} \sigma \\ \sigma \end{bmatrix}$     0

$w1 = \overline{a}ab\_cde$
$aabcde$
$w2 = abpcg \_abpcg\_$

The edit path    $\begin{bmatrix} a \\ \varepsilon \end{bmatrix} \begin{bmatrix} a \\ a \end{bmatrix} \begin{bmatrix} b \\ b \end{bmatrix} \begin{bmatrix} \varepsilon \\ p \end{bmatrix} \begin{bmatrix} c \\ c \end{bmatrix} \begin{bmatrix} d \\ g \end{bmatrix} \begin{bmatrix} e \\ \varepsilon \end{bmatrix}$

The edit path cost:    4

$ED(w_1, w_2) = \min \{ \text{weight of ops needed to transforms } w_1 \text{ to } w_2 \}$

© Ilay Tzarfati

# Issue with ED

$w_1 = aaaaa$

$w_2 = abpcg$

$ED(w_1, w_2) = 4$

$w_3 = aaaaaaa...aaaaa$

$w_4 = aaaaaaa...abpcg$

$ED(w_3, w_4) = 4$

But these are much more similar

# How can we normalize ED?

- The **sum** edit-distance: $ED_{sum}(w_1, w_2)$

$$= \frac{ED(w_1, w_2)}{|w_1| + |w_2|}$$

- The **max** edit-distance: $ED_{max}(w_1, w_2)$

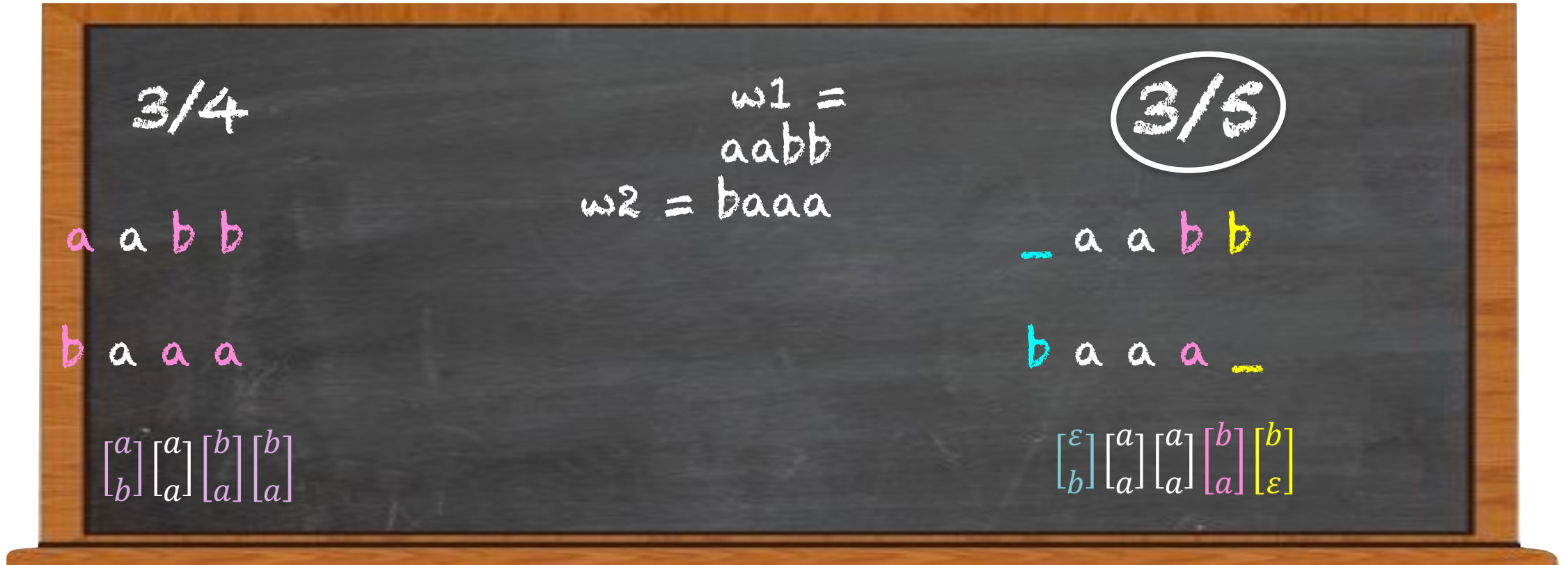$$= \frac{ED(w_1, w_2)}{max(|w_1|, |w_2|)}$$

All the above <u>do not</u> satisfy the triangle inequality! (hence is not a metric)

How can this be solved?
Marzal and Vidal suggested
dividing by the length of the edit path

# Normalizing using path length

3/4

$$w1 = aabb$$
$$w2 = baaa$$

$$\boxed{3/5}$$

a a b b

b a a a

_ a a b b

b a a a _

$$\begin{bmatrix} a \\ b \end{bmatrix} \begin{bmatrix} a \\ a \end{bmatrix} \begin{bmatrix} b \\ a \end{bmatrix} \begin{bmatrix} b \\ a \end{bmatrix}$$

$$\begin{bmatrix} \varepsilon \\ b \end{bmatrix} \begin{bmatrix} a \\ a \end{bmatrix} \begin{bmatrix} a \\ a \end{bmatrix} \begin{bmatrix} b \\ a \end{bmatrix} \begin{bmatrix} b \\ \varepsilon \end{bmatrix}$$

$$_2) = \min \{ \text{cost}(P) : P \text{ is an edit path transforming } w_1 \text{ to } w_2 \}$$

$$\text{Cost}(P) = \frac{weight(P)}{length(P)}$$

8

# NED- Normalized Edit Distance

$(w_1, w_2) = \min \{ cost(P): P \text{ is an edit path transforming } w_1 \text{ to } w_2 \}$

$$Cost(P) = \frac{weight(P)}{length(P)}$$

$w1 =$
aaaaa
$w2 =$
abpcg
$NED(w_1, w_2) = 4/5$

$w3 =$
aaaaaa..aaaaa
$w4 =$ aaaaaaa...
abpcg
$NED(w_3, w_4) = 4/10$
0

Matches our intuition ☺

But is it a metric ???

© Ilay Tzarfati

# Is NED a Metric ?

### Normalized Edit Distance (NED)



[Marzal & Vidal 93]

### Generalized Edit Distance (GED)



[Yujian & Bo '07]

### Contextual Edit Distance (CED)



[de la Higuera & Mico'08]

10

# The Problem

Weight function :
$$d: (\Sigma \cup \{\varepsilon\}) \times (\Sigma \cup \{\varepsilon\}) \to [0,1]$$

$NED_d$ that induced by d:
$$NED_d: \Sigma^* \times \Sigma^* \to [0,1]$$

When is $NED_d$ a metric ?

d:
( Pair of letters ) → ( [0,1] )

$NED_d$:
( Pair of words ) → ( [0,1] )

# If d is uniform, it is !

**The Normalized Edit Distance with Uniform Operation Costs is a Metric**

Dana Fisman ✉
Dept. of Computer Science, Ben-Gurion University, Beer-Sheva, Israel

Joshua Grogin ✉
Dept. of Computer Science, Ben-Gurion University, Beer-Sheva, Israel

Oded Margalit ✉
Dept. of Computer Science, Ben-Gurion University, Beer-Sheva, Israel

Gera Weiss ✉
Dept. of Computer Science, Ben-Gurion University, Beer-Sheva, Israel

—— Abstract ——

We prove that the normalized edit distance proposed in [Marzal and Vidal 1993] is a metric when the cost of all the edit operations are the same. This closes a long standing gap in the literature where several authors noted that this distance does not satisfy the triangle inequality in the general case, and that it was not known whether it is satisfied in the uniform case — where all the edit costs are equal. We compare this metric to two normalized metrics proposed as alternatives in the literature, when people thought that Marzal's and Vidal's distance is not a metric, and identify key properties that explain why the original distance, now known to also be a metric, is better for some applications. Our examination is from a point of view of formal verification, but the properties and their significance are stated in an application agnostic way.

OK.... But in the rest of the cases ???

# Our main result

A necessary and sufficient condition for $\text{NED}_d$ to be a Metric.

Surprisingly, d being a metric is neither sufficient nor necessary.

# Essential edit operations

An edit operation is **essential**, if there exist optimal edit path that use it.

$$a \xrightarrow[\substack{\begin{bmatrix}a\\b\end{bmatrix} \\ \begin{bmatrix}a\\\varepsilon\end{bmatrix} \begin{bmatrix}\varepsilon\\b\end{bmatrix}}]{} b$$

$\forall a \in \Sigma$, $(a,\varepsilon)$ and $(\varepsilon,a)$ are essential.

$d(a,b) \geq d(a,\varepsilon) + d(\varepsilon,b)$ iff $(a,b)$ is inessential.

$d(a,\varepsilon) = d(\varepsilon,a) = 0.4$ w1=a

w2=b

$d(b,\varepsilon) = d(\varepsilon,b) = 0.3$

$d(a,b) = d(b,a) = 0.8$ $\begin{bmatrix}a\\\varepsilon\end{bmatrix}\begin{bmatrix}\varepsilon\\b\end{bmatrix}$

$NED_d(\text{w1,w2}) = \frac{0.4+0.3}{2} = 0.35$

In the example both $(a,b)$, $(b,a)$ are inessential

# The Necessary and sufficient condition for $\text{NED}_d$ to be a Metric

**Weight function : Let** $d: (\Sigma \cup \{\varepsilon\}) \times (\Sigma \cup \{\varepsilon\})$
$\rightarrow [0,1]$ Let $a,c \in \Sigma \cup \{\varepsilon\}$ and $b$

$\in \Sigma$

✓ 1. identity:
   d(a,c) = 0 iff a=c

✓ 2. symmetry:
   d(a,c) = d(c,a)

? 3. Relaxed triangle inequality :
   d(a,b) + d(b,c) ≥min { d(a,c), d(a,$\varepsilon$)+
   d($\varepsilon$,c)}

? 4. At least half:
   d($\varepsilon$,b) = d(b,$\varepsilon$) ≥ $\frac{1}{2}$

*A d satisfying those properties, is termed*

15

# Sketch of the idea:
# The (Relaxed) Triangle Inequality

$\min\{d(a,c), d(a,\varepsilon) + d(\varepsilon,c)\} \leq d(a,b) + d(b,c)$

Obviously to satisfy the triangle inequality we need:

$d(a,c) \leq d(a,b) + d(b,c)$

But since we know that $d(a,c)$ can be replace with $d(a,\varepsilon)$, $d(\varepsilon,c)$ we require that:

$$\min\{d(a,c), d(a,\varepsilon) + d(\varepsilon,c)\} \leq d(a,b) + d(b,c)$$

# Sketch of the idea:
# At least half

$$d(b,\varepsilon) = d(\varepsilon, b) \geq \frac{1}{2}$$

We want to make sure that inflating the edit path is not worthwhile .

$w_1 = a \quad w_2 = ac^{100} \quad w_3 = b$

d($\varepsilon$,c) = d(c,$\varepsilon$) = 0.2,
d(a,b) = d(b,a) = 0.55 ,
d(a,$\varepsilon$) = d($\varepsilon$,a) = d(b,$\varepsilon$) = d($\varepsilon$,b) = 0.6 ...

Optimal edit path $[w_1 \rightarrow w_3]$: $\begin{bmatrix} a \\ b \end{bmatrix}$
cost : 0.55.

edit path $[w_1 \rightarrow w_2]$:

$\begin{bmatrix} a \\ a \end{bmatrix} \begin{bmatrix} \varepsilon \\ c \end{bmatrix} \cdots \begin{bmatrix} \varepsilon \\ c \end{bmatrix} \begin{bmatrix} \varepsilon \\ c \end{bmatrix}$

cost $\frac{100*0.2}{101} = 0.198$

edit path $[w_2 \rightarrow w_3]$:

$\begin{bmatrix} a \\ b \end{bmatrix} \begin{bmatrix} c \\ \varepsilon \end{bmatrix} \cdots \begin{bmatrix} c \\ \varepsilon \end{bmatrix} \begin{bmatrix} c \\ \varepsilon \end{bmatrix}$

cost $\frac{100*0.2+0.5}{101} = 0.203$

$NED_d$ is not Metric
$0.198 + 0.203 < 0.55$

1. identity:
   d(a,c) = 0 iff a=c

2. symmetry:
   d(a,c) = d(c,a)

3. Relaxed triangle inequality :
   $d(a,b) + d(b,c) \geq \min \{ d(a,c),$
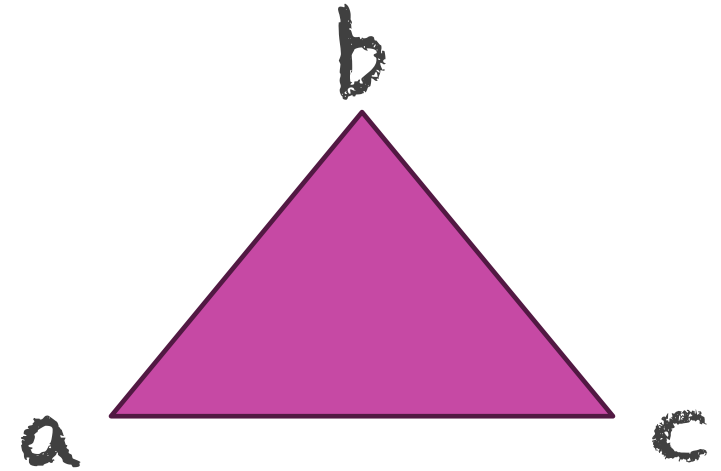   $d(a,\varepsilon)+ d(\varepsilon,c)\}$
4. At least half:
   $d(\varepsilon,b) = d(b,\varepsilon) \geq \frac{1}{2}$

Are necessary.

# Sufficient condition

$d$ is fine $\rightarrow NED_d$ is metric

Let $w_1, w_2, w_3 \in \Sigma^*$ :

1. $NED_d\ (w_1, w_2) = 0$ iff $w_1 = w_2$ ✓

2. $NED_d\ (w_1, w_2) = NED_d\ (w_2, w_1)$ ✓

3. $NED_d\ (w_1, w_2) + NED_d\ (w_2, w_3) \geq NED_d\ (w_1, w_3)$ ?

# $NED_d$ follows the triangle inequality

$$\begin{bmatrix} a \\ \cancel{a} \\ \cancel{b} \\ c \end{bmatrix}$$

**Compose:**

input : P(w1->w2) , P(w2->w3)

Output: P(w1->w3)

d fine -> cost(P(w1->w3)) ≤ cost(P(w1->w2) ) + cost(P(w2->w3))

# Proof idea:

w1 = aaa

w2 = bb

w3 = cccc

$P(1\rightarrow2) = \begin{bmatrix}a\\b\end{bmatrix}\begin{bmatrix}a\\\varepsilon\end{bmatrix}\begin{bmatrix}a\\b\end{bmatrix}$

$P(2\rightarrow3) = \begin{bmatrix}\varepsilon\\c\end{bmatrix}\begin{bmatrix}b\\c\end{bmatrix}\begin{bmatrix}b\\c\end{bmatrix}\begin{bmatrix}\varepsilon\\c\end{bmatrix}$

a a a
b _ b
_ b b _
c c c c

**Problem**

$P(1\rightarrow2) = \begin{bmatrix}\varepsilon\\\varepsilon\end{bmatrix}\begin{bmatrix}a\\b\end{bmatrix}\begin{bmatrix}a\\\varepsilon\end{bmatrix}\begin{bmatrix}a\\b\end{bmatrix}\begin{bmatrix}\varepsilon\\\varepsilon\end{bmatrix}$

$P(2\rightarrow3) =$

$P(1\rightarrow3) = \begin{bmatrix}\varepsilon\\c\end{bmatrix}\begin{bmatrix}a\\c\end{bmatrix}\begin{bmatrix}a\\\varepsilon\end{bmatrix}\begin{bmatrix}a\\c\end{bmatrix}\begin{bmatrix}\varepsilon\\c\end{bmatrix}$

a a a
b _ b
_ b _ b
c c _ c
c

*alignment*

# Examples for fine d's

Distances in [0,n] : Let $d: [0,n] \times [0,n] \rightarrow [0,1]$ be defined as follows :

$$d(n_1, n_2) = \frac{|n_1 - n_2|}{n + 1}$$

Distances in $\mathbb{N}$: Let $d: \mathbb{N} \times \mathbb{N} \rightarrow [0,1]$ be defined as follows :

$$d(n_1, n_2) = 1 - \frac{1}{|n_1 - n_2| + 1}$$

Distances in $\Sigma = 2^k$: Let $d: 2^k \times 2^k \rightarrow [0,1]$ be defined as follows :

$$d(v_1, v_2) = \frac{HD(v_1, v_2)}{k}$$

*NEDd is metric*

# Applications in Formal Verification

- FV requires $\omega$-words
- 
  Robustness requires $NED$ between omega-words

- [FGW23] suggested $\omega NED$

- $\omega NED$ can be generalized to $\omega NED_d$

- Same algorithms are applicable.

# Paper results

- We have shown necessary and sufficient conditions for $NED_d$ to be metric.

- We have shown several fine d's .

- We have shown that $NED_d$ can also be used for formal verification.

Thanks for listening ! Any questions ?