MATH 2301

* Non-regular languages

Q: How do we know that non-regular languages exist?

A1: Counting!

Say $\Sigma$ = alphabet (finite)

$\Sigma^*$ = all possible strings.

$\Sigma^*$ has $\infty^{ly}$ many, but countably many elements.
This means that there is a procedure to sequentially number $(1,2,3, \cdots)$ the elements of $\Sigma^*$, namely, the lexicographic order.
[There is a bijection $\Sigma^* \xrightarrow{1:1} \mathbb{N}$]

A language $L$ is just any subset of $\Sigma^*$.

How many languages? Infinitely many, as many as there are subsets of $\Sigma^*$, or subsets of $\mathbb{N}$.
Theorem: There are uncountably many subsets of $\mathbb{N}$.

$\Rightarrow$ there is no way to number the subsets of $\mathbb{N}$ by $1, 2, 3, \cdots$ [you'll always miss something if you try.]

(Proof goes via an argument called Cantor's diagonalisation argument.)

$\Rightarrow$ There are uncountably many languages!

Q: How many are regular?

Observation: At most as many as there are regular expressions.

A regex is just a special string in the alphabet $\Sigma$, together with extra symbols: $*, |, (,)$

Number of regexes $\leq$ Number of strings in $\Sigma \cup \{*, |, (,)\}$

Countable.                          Countable!
                                    Lexicographically orderable.

But  Number of regular languages $\leq$ Number of regexes.

Countable.

Since there are uncountably many languages, at least one (in fact, uncountably many) are non-regular.

## A2: The pumping lemma

We exploit a non-obvious feature that every regular language shares _(↙ if and only if)_

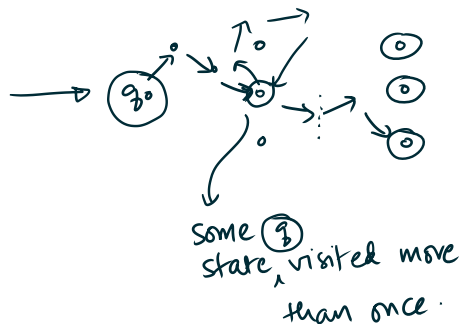**Recall:** A language is regular iff there is a DFA that recognises it.

Suppose $M$ is a DFA. Schematic below.

Suppose it has $n$ states.



If $w$ is any word of length $k$, then the computation of $w$ through $M$ passes through $(k+1)$ states, and gives a path of length $k$ through the DFA.

(Recall from graphs): If $|w| > n$, then it must repeat a state.



If $w$ is long enough, there is a "loop" within the calculation path

some $\text{⑨}$ state, visited more than once.

$\Rightarrow w = xyz$ :
$x$ = portion before $\text{⑨}$
$y$ = loop portion ⟵ nonzero length.
$z$ = portion after $\text{⑨}$

---

Suppose $w$ has length $\geq n$ and $w$ is accepted by $M$.
$w = xyz$ as before.

$\Rightarrow xz$ is also accepted by $M$.
& $xyyz$, $xyyyyyz$, etc are all accepted by $M$.

$\Rightarrow$ any pattern of the form $x\,y^*z$ is accepted by $M$.

$\Rightarrow$ If $M$ accepts a string $w$ of length $> n$ then it accepts all the strings fitting into the pattern $x\,y^*z$ : there is a non-empty "$y$" portion that can be pumped.

Use this to detect non-regular languages:
If there are long enough strings that can't be pumped, then the language is not regular.

** __Theorem__ (Pumping lemma): Let $L$ be a regular language. Then there exists some $n_L \in \mathbb{N}$, such that: if $w \in L$ and $|w| \geq n_L$, then:  ↱ # states in a DFA recognising $L$

$w = xyz$ with

1) $|y| \geq 1$

2) $|xy| < n_L$

3) $xy^k z \in L$ for every $k \in \mathbb{N}$ (including $k = 0$)!

__Example__: $\{0^k 1^k \mid k \in \mathbb{N}\} = L$

If $L$ were regular, there would be a DFA $M$, such that $L = L(M)$.

Let $n = $ # states in this hypothetical DFA.

Consider the string $0^{n+1} 1^{n+1} = w$.
   ↱ already longer than $n$.

Running $w$ through $M$, we will encounter a repeated state even before we get to the $1$s.

$\Rightarrow w = xyz$
      ↱ involves only zeroes.

$\Rightarrow x = 0^a$
   $y = 0^b$ $(b > 0)$
   $z = 0^c 1^{n+1}$

$(a + b + c = n+1)$

$\left.\begin{array}{l}\end{array}\right\}$ previous argument
$\Rightarrow xyyz, xyyyz$ etc are all in $L$.

$xyyz = 0^a 0^b 0^b 0^c 1^{n+1}$

$\quad = 0^{n+1+b} 1^{n+1}$

↗ DFA calculation $\Rightarrow$ in $L$

↘ doesn't satisfy def!

__Contradiction__: