

Content and Context: Two-pronged Bootstrapped Learning for Regex-formatted Entity Extraction



Stanley Simoes[†], Deepak P[#], Munu Sairamesh[†], Deepak Khemani[†], Sameep Mehta[‡]

[†]Indian Institute of Technology Madras

[#]Queen's University Belfast

[‡]IBM Research - India

IBM Research

1. REGEXES FOR ENTITY EXTRACTION

- Regular Expressions (regexes)* can be used to elegantly characterize entities having an underlying syntactical pattern
- Easy to develop, interpret, maintain, and incorporate domain knowledge
- High quality regexes have:
 - ▷ **high precision:** all matches are instances
 - ▷ **high recall:** all instances are matches
- Difficult to manually design a high quality regex for entity extraction
 - ▷ A human would generalize some instances to a high precision regex, but would miss out on some variants
 - ▷ eg: `CS\d{3}` does not match CS courses *CS6880* and *CSE578* (high precision, low recall)
- Need a system that automatically learns and suggests:
 - overlooked instances
 - high quality regexes
- Use cases:
 - ▷ design a regex for a new entity extraction system
 - ▷ enhance regex in an existing entity extraction system

ENTITY	REGEX	MATCHED INSTANCE
CS course	<code>CS\d{3}</code>	<i>CS376</i>
Intel CPU	<code>i\d-\d{3}</code>	<i>i5-750</i>

2. PROBLEM STATEMENT

- Given:**
 - ▷ a high precision seed regex r_{seed}
 - ▷ a document corpus \mathcal{D} i.e., a set of seed matches \mathcal{M}_{seed}
- Goal:** To identify
 - \mathcal{M}_{exp} : instances in \mathcal{D} not matched by r_{seed}
 - \mathcal{R} : set of diverse, high quality regexes
- Proposed Approach:** (Two stages)
 - Stage I: Match Set Expansion**

CS376 *CS290* \mathcal{M}_{seed} → MATCH SET EXPANDER → \mathcal{M}_{exp} *CS6880* *CSE578*
 - Stage II: Regex Recommendation**

CSE578 *CS6880* \mathcal{M}_{exp} → REGEX RECOMMENDER → \mathcal{R}

\mathcal{R} includes: `CS\d{3}` (seed), `CS[A-Z]{0,1}\d{3}`, `CS\d{3,4}`

3. HIGH QUALITY REGEXES

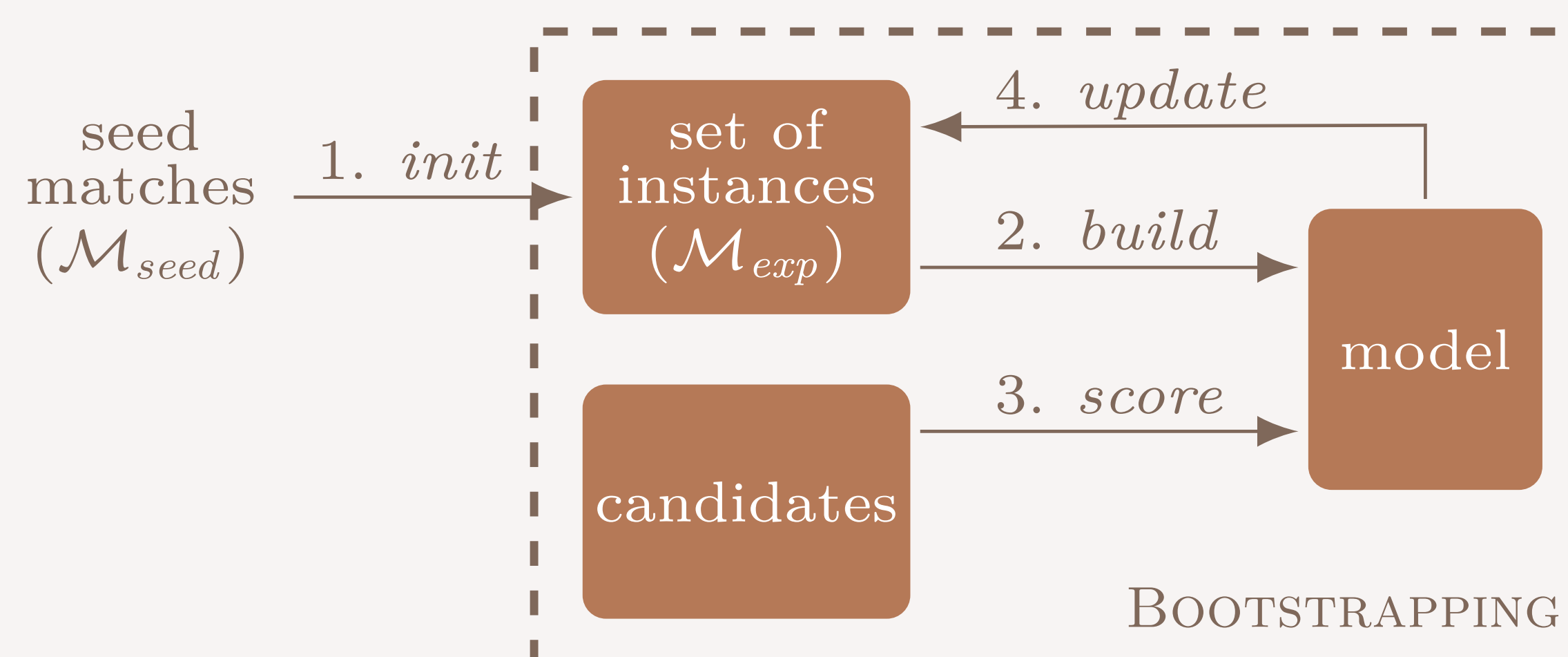
- Search the *generalization space* of r_{seed}
- Regex units: **character class** & **quantifier**
 - ▷ eg: `[A-Z]{2}\d{3,4}`
- Generalizing a regex:
 - allow more characters in character class
 - decrease quantifier's lower bound
 - increase quantifier's upper bound
 - ▷ eg:

$[A-Z]{2,4}$
 a) `[a-zA-Z]{2,4}`

$[A-Z]{1,4}$
 b) `[A-Z]{1,4}`

$[A-Z]{2,5}$
 c) `[A-Z]{2,5}`

4. MATCH SET EXPANSION

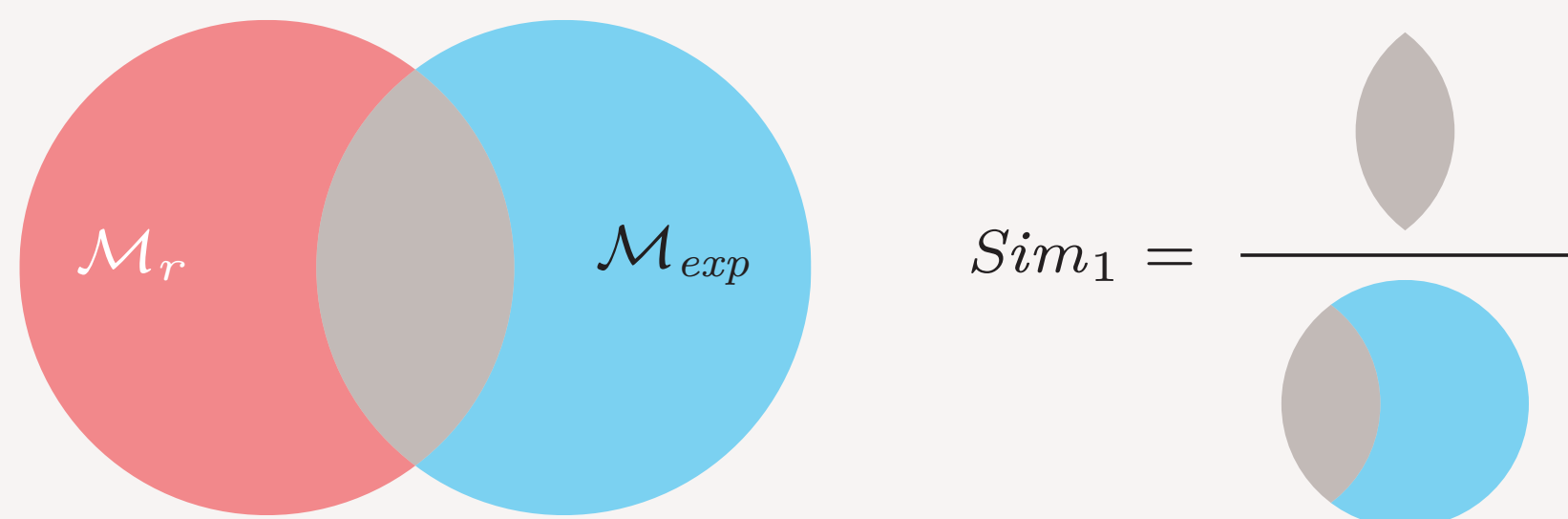


- Expand using **Bootstrapping**
 - ▷ repeat {build → score → update}
- candidates:**
 - ▷ matches of r_{seed} 's generalizations
- model: Logistic Regression**
 - ▷ \mathcal{M}_{exp} as 1
 - ▷ candidates as 0
- Matches represented using:
 - ▷ content
 - ▷ context

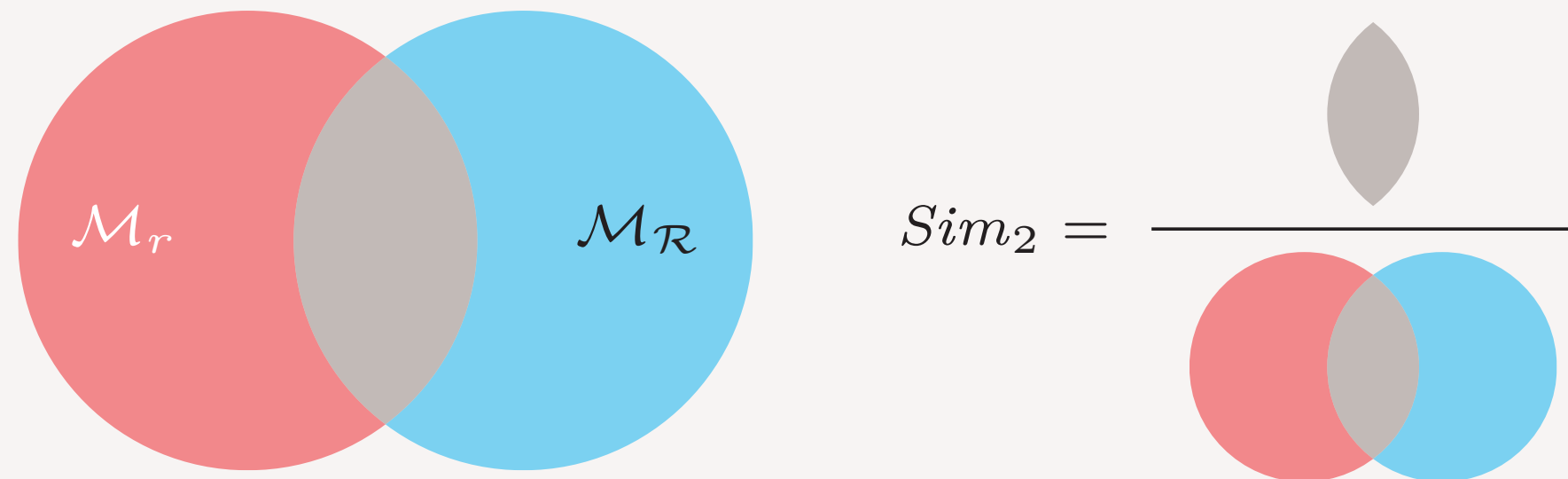
5. REGEX RECOMMENDATION

- Generate a set of regexes \mathcal{R} with
 - ▷ high recall: $\forall r \in \mathcal{R}, r$ is a generalization of r_{seed}
 - ▷ high precision: $\forall r \in \mathcal{R}, r$'s matches are in \mathcal{M}_{exp}
 - ▷ high diversity: $\forall r_1, r_2 \in \mathcal{R}, r_1$'s matches are not in r_2 's
- Using **Maximal Marginal Relevance**
 - ▷ add generalization r to \mathcal{R} that maximizes

$$\underbrace{\lambda \cdot \text{Sim}_1(r, \mathcal{M}_{exp})}_{\text{RELEVANCE TERM}} - \underbrace{(1 - \lambda) \cdot \text{Sim}_2(r, \mathcal{R})}_{\text{DIVERSITY TERM}}$$
 - ▷ Sim_1 : similarity to \mathcal{M}_{exp}



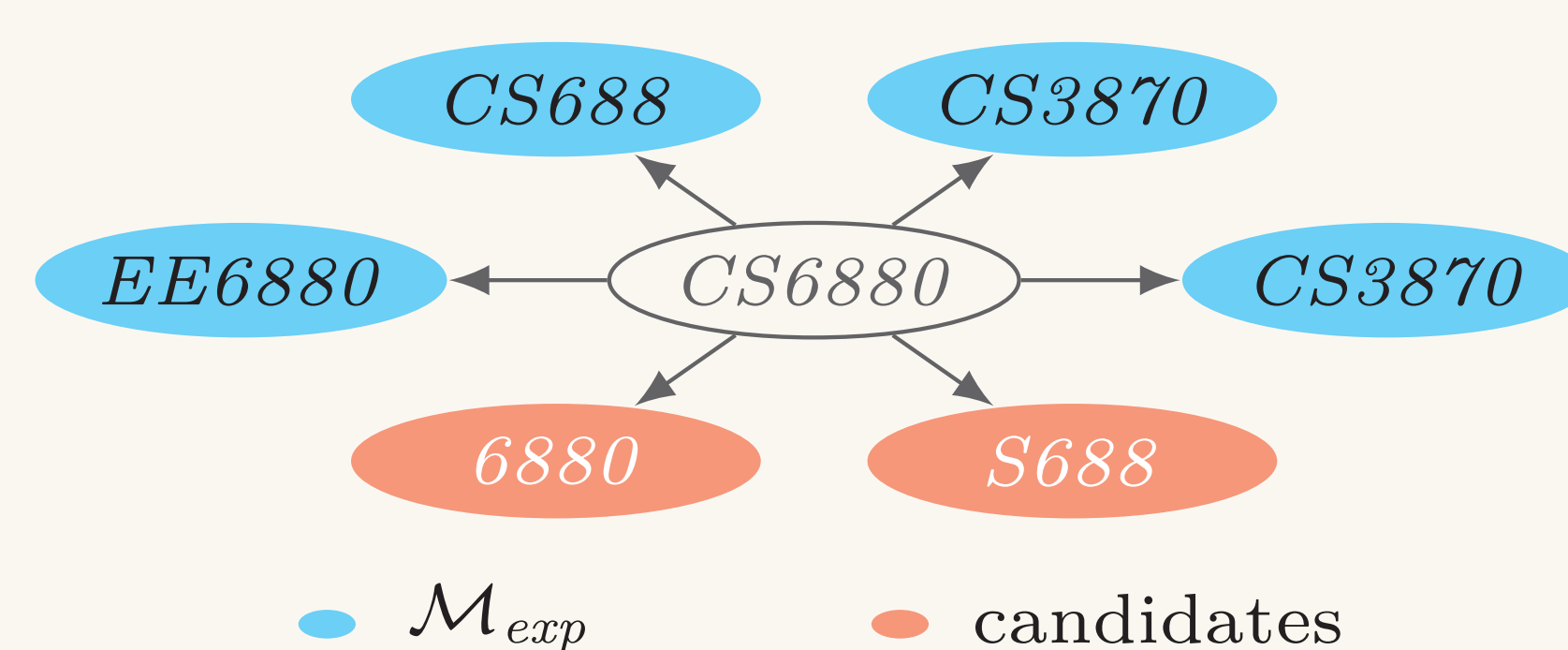
▷ Sim_2 : similarity to \mathcal{R}



* \mathcal{M}_r = r 's matches in \mathcal{D}
 $\mathcal{M}_\mathcal{R}$ = matches in \mathcal{D} for all $r \in \mathcal{R}$

CONTENT FEATURES

- Matches in \mathcal{M}_{exp} within k edit distance



$$\text{Content-Score}_k(\text{CS6880}) = \frac{4}{6} = 0.67$$

- F1:** Content-Score_2 (2 edit distance)
- F2:** Content-Score_3 (3 edit distance)

CONTEXT FEATURES

- Similarity to contexts of matches in \mathcal{M}_{exp}



- Language model L_P : \mathcal{M}_{exp} 's left contexts

$\dots w_1 w_2 w_3 \text{ CS6680 } \dots$

$$\text{Context-Score}_L(\text{CS6680}) = \sum_i \ln L_P(w_i)$$

- F3:** Context-Score_L (left contexts)
- F4:** Context-Score_R (right contexts)

6. EXPERIMENTS

- Two extraction tasks:
 - ▷ DATE_{WEBKB} - `\d{2}/\d{2}/\d{2}`
 - ▷ COURSE_{WEBKB} - `CS\d{3}`
- Parameters:
 - ▷ ≤ 4 regex units generalized
 - ▷ 1% expansion per iteration
 - ▷ 150 iterations
 - ▷ $\lambda = 0.7$

COMPARISON OF OUR MATCH SET EXPANDER'S ACCURACY WITH BASELINES

TASK	FREQ	PRECISION			\mathcal{M}_{seed}	RECALL			\mathcal{M}_{seed}	F-SCORE			
		GM ₁₀	GM _{1%}	OURS		GM ₁₀	GM _{1%}	OURS		FREQ	GM ₁₀	GM _{1%}	OURS
DATE _{WEBKB}	0.032	0.186	0.479	1.000	0.251	0.436	0.330	0.857	0.402	0.063	0.261	0.391	0.923
COURSE _{WEBKB}	0.070	0.672	0.633	0.994	0.342	0.348	0.349	0.855	0.509	0.131	0.459	0.450	0.919

▷ PRECISION of $\mathcal{M}_{seed} = 1$

▷ RECALL of FREQ = 1

*best values **boldfaced**

RECOMMENDED REGEXES

TASK	DATE _{WEBKB}	COURSE _{WEBKB}
r_{seed}	<code>\d{2}/\d{2}/\d{2}</code> eg: 12/13/01	<code>CS\d{3}</code> eg: CS376
\mathcal{R}	<code>\d{1,2}/\d{1,2}/\d{2}</code> eg: 2/3/01 <code>\d{2}/\d{1,2}/\d{2}</code> eg: 12/3/01 <code>\d{1,2}/\d{2}/\d{2}</code> eg: 2/13/01	<code>C?[a-zA-Z]{1,2}\d{3}</code> eg: CSE578 <code>[a-zA-Z]{1,2}S?\d{3}</code> eg: cs615 <code>CS\u{w}{1,3}\d\u{w}</code> eg: CS382M

*generalized units in **red and underlined**



✉ stanleysimoes@gmail.com