

Belief propagation on strings

Boris Yangel

Last modified: November 27, 2013

Abstract

This document describes a way of performing belief propagation in models with string variables, with factors corresponding to widely used operations on strings such as **Substring**, **Contains**, **Replace**, **Concat** etc.

Contents

1	Notation	1
2	Automata	2
2.1	Atomic automata	2
2.2	Operations on automata	2
2.3	Examples of compound automata	2
3	Transducers	2
3.1	Atomic transducers	3
3.2	Operations on transducers	3
3.3	Projection of an automaton onto a transducer	3
3.4	Some properties of operations on transducers	3
3.5	Examples of compound transducers	4
4	Message operators	4
4.1	Substring	4
4.2	Concat	5
4.3	ContainsAt	5

1 Notation

We denote automata by capital letters A, B etc or, in case of a specific automaton, by a string starting from a capital letter, such as Unif. However, when an automaton represents a BP message, we denote it by μ with a subscript describing the source and the destination of the message, like $\mu_{f \rightarrow s}$. Notation for transducers is the same, except that we use a different font: \mathcal{T} , \mathcal{U} , *Copy*. Lowercase letters and sequences of lowercase letters denote variables: s , t , str .

The notation $s \cdot t$ stands for the concatenation of strings s and t .

2 Automata

We represent marginals and messages from and to string variables by weighted finite state automata. Weighted finite state automaton, which we'll refer to simply as automaton from now on, is a function mapping strings to real values. Not every such function is a finite state automaton though. For an explanation of what a finite state automaton really is and what are the limitations of this class of functions we refer the reader to [?], while this paper takes a different approach. We instead define a set of atomic automata, a number of operations w.r.t. which automata are closed, and then simply use all the automata that can be obtained from the atomic ones using the proposed operations.

2.1 Atomic automata

Name	Notation	Definition
Uniform of characters	Unif_C	Maps every string consisting entirely of characters from the set C to 1 and all other strings to 0.
Single character	Char_C	Maps every string consisting of a single character from the set C to 1 and all other strings to 0.

If the set C consists of all possible characters, we denote Unif_C simply by Unif , and Char_C by Char .

2.2 Operations on automata

Name	Notation	Definition
Scale	αA	$(\alpha A)(s) := \alpha A(s)$
Sum	$A + B$	$(A + B)(s) := A(s) + B(s)$
Product	AB	$(AB)(s) := A(s)B(s)$
Concatenation	$A \cdot B$	$(A \cdot B)(s) := \sum_{s_1 \cdot s_2 = s} A(s_1)B(s_2)$.
Normalizer	$\mathcal{Z}(A)$	$\mathcal{Z}(A) := \sum_s A(s)$ (defined only when the corresponding sum converges).

We denote $A + (-1)B$ simply by $A - B$.

2.3 Examples of compound automata

Notation	Definition	Value
OfLen_N	$\underbrace{\text{Char} \cdot \text{Char} \cdot \dots \cdot \text{Char}}_{N \text{ times}}$	1 on any string of length N , zero on everything else.

3 Transducers

In order to compute messages from a factor to a variable, we also need to employ weighted finite state transducers, which are, essentially, functions mapping pairs of sequences to real values. Following the section 2, we define the set of atomic transducers and a number of operations that allows us to produce all the transducers of interest. For a more general introduction, we refer the reader to [?].

3.1 Atomic transducers

Name	Notation	Definition
Copy	$Copy$	$Copy(s, t) := \begin{cases} 1, & s = t \\ 0 & \text{otherwise} \end{cases};$
Copy char	$CopyChar$	$CopyChar(s, t) := \begin{cases} 1, & s = t, s = t = 1 \\ 0 & \text{otherwise} \end{cases};$
Consume an automaton	$Consume_A$	$Consume_A(s, t) := \begin{cases} A(s), & t = "" \\ 0 & \text{otherwise} \end{cases};$
Produce an automaton	$Produce_A$	$Produce_A(s, t) := \begin{cases} A(t), & s = "" \\ 0 & \text{otherwise} \end{cases}.$

3.2 Operations on transducers

Name	Notation	Definition
Concatenation	$\mathcal{A} \cdot \mathcal{B}$	$(\mathcal{A} \cdot \mathcal{B})(s, t) := \sum_{\substack{s_1, t_1, s_2, t_2 \\ s_1 \cdot s_2 = s, t_1 \cdot t_2 = t}} \mathcal{A}(s_1, t_1) \mathcal{B}(s_2, t_2)$
Transpose	\mathcal{A}^T	$(\mathcal{A}^T)(s, t) := \mathcal{A}(t, s)$
Superposition	$\mathcal{B}(\mathcal{A})$	$(\mathcal{B}(\mathcal{A}))(s, t) := \sum_u A(s, u) B(u, t)$

3.3 Projection of an automaton onto a transducer

Additionally, an automaton can be projected onto a transducer. The result of the projection is another automaton, which is defined as

$$\text{proj}[A, \mathcal{T}](t) := \sum_s A(s) \mathcal{T}(s, t). \quad (1)$$

This operation can be very handy when computing BP messages.

3.4 Some properties of operations on transducers

Property 1.

$$\text{proj}[A, \mathcal{U}(\mathcal{T})] = \text{proj}[\text{proj}[A, \mathcal{T}], \mathcal{U}]. \quad (2)$$

Proof.

$$\begin{aligned} \text{proj}[A, \mathcal{U}(\mathcal{T})](s, t) &= \sum_s A(s) \sum_u \mathcal{T}(s, u) \mathcal{U}(u, t) = \sum_u \mathcal{U}(u, t) \sum_s A(s) \mathcal{T}(s, u) = \\ &= \sum_u \mathcal{U}(u, t) \text{proj}[A, \mathcal{T}](u) = \text{proj}[\text{proj}[A, \mathcal{T}], \mathcal{U}]. \end{aligned} \quad (3)$$

Property 2.

$$\mathcal{Z}(A) = \text{proj}[\text{Unif}, Consume_A \cdot Produce_{\text{Unif}}](t) \quad \forall t. \quad (4)$$

Proof.

$$\forall t : \quad \mathcal{Z}(A) = \sum_s A(s) = \sum_s \text{Unif}(s) \tilde{A}(s, t), \quad (5)$$

where

$$\tilde{A}(s, t) = A(s) = \text{Consume}_A \cdot \text{Produce}_{\text{Unif}}(s, t). \quad (6)$$

3.5 Examples of compound transducers

Here we present a number of compound transducers that are useful for message computation.

Notation	Definition	Value
$\text{Occurrences}(s, \text{str})$	$\text{Produce}_{\text{Unif}} \cdot \text{Copy} \cdot \text{Produce}_{\text{Unif}}$	The number of occurrences of s in str .
$\text{OccursAt}_p(s, \text{str})$	$\text{Produce}_{\text{OfLen}_p} \cdot \text{Copy} \cdot \text{Produce}_{\text{Unif}}$	1 if s occurs in str at the position p , 0 otherwise.
$\text{CopyOfLen}_N(s_1, s_2)$	$\underbrace{\text{CopyChar} \cdot \text{CopyChar} \cdot \dots \cdot \text{CopyChar}}_{N \text{ times}}$	1 if $s_1 = s_2$ and both strings are of length N .
$\text{IsSubstring}_{p,l}(s, \text{str})$	$\text{Produce}_{\text{OfLen}_p} \cdot \text{CopyOfLen}_l \cdot \text{Produce}_{\text{Unif}}$	1 if s occurs in str from the position p to $p + l$, 0 otherwise.

4 Message operators

4.1 Substring

4.1.1 Factor

$$f(\text{str}, s, p, l) = \mathbb{1}[\text{str contains } s \text{ from } p \text{ to } p + l]. \quad (7)$$

We assume that the position p and the length l are observed.

4.1.2 Message to str

$$\begin{aligned} \mu_{f \rightarrow \text{str}}(\text{str}) &= \sum_s \mu_{s \rightarrow f}(s) \mathbb{1}[\text{str contains } s \text{ from } p \text{ to } p + l] = \\ &= \text{proj}[\mu_{s \rightarrow f}, \text{IsSubstring}_{p,l}]. \end{aligned} \quad (8)$$

4.1.3 Message to s

$$\begin{aligned} \mu_{f \rightarrow s}(s) &= \sum_{\text{str}} \mu_{\text{str} \rightarrow f}(s) \mathbb{1}[\text{str contains } s \text{ from } p \text{ to } p + l] = \\ &= \text{proj}[\mu_{\text{str} \rightarrow f}, \text{IsSubstring}_{p,l}^T]. \end{aligned} \quad (9)$$

4.2 Concat

4.2.1 Factor

$$f(s, s_1, s_2) = \mathbb{1}[s = s_1 \cdot s_2]. \quad (10)$$

4.2.2 Message to s

$$\mu_{f \rightarrow s}(s) = \sum_{s_1} \mu_{s_1 \rightarrow f}(s_1) \sum_{s_2} \mu_{s_2 \rightarrow f}(s_2) \mathbb{1}[s = s_1 \cdot s_2] = \mu_{s_1 \rightarrow f} \cdot \mu_{s_2 \rightarrow f}. \quad (11)$$

4.2.3 Message to s_1

$$\begin{aligned} \mu_{f \rightarrow s_1}(s_1) &= \sum_s \mu_{s \rightarrow f}(s) \sum_{s_2} \mu_{s_2 \rightarrow f}(s_2) \mathbb{1}[s = s_1 \cdot s_2] = \\ &= \text{proj}[\mu_{s \rightarrow f}, \text{Copy} \cdot \text{Consume}_{\mu_{s_2 \rightarrow f}}]. \end{aligned} \quad (12)$$

4.2.4 Message to s_2

$$\begin{aligned} \mu_{f \rightarrow s_2}(s_2) &= \sum_s \mu_{s \rightarrow f}(s) \sum_{s_1} \mu_{s_1 \rightarrow f}(s_1) \mathbb{1}[s = s_1 \cdot s_2] = \\ &= \text{proj}[\mu_{s \rightarrow f}, \text{Consume}_{\mu_{s_1 \rightarrow f}} \cdot \text{Copy}]. \end{aligned} \quad (13)$$

4.3 ContainsAt

4.3.1 Factor

$$\begin{aligned} f(str, s, p, c) &= c \mathbb{1}[str \text{ contains } s \text{ at pos } p] + \\ &+ (1 - c)(1 - \mathbb{1}[str \text{ contains } s \text{ at pos } p]). \end{aligned} \quad (14)$$

We assume that the position p is always observed.

4.3.2 Message to c

$$\begin{aligned} \mu_{f \rightarrow c}(c = 1) &= \sum_{str} \mu_{str \rightarrow f}(str) \sum_s \mu_{s \rightarrow f}(s) \mathbb{1}[str \text{ contains } s \text{ at pos } p] = \\ &= \mathcal{Z}(\mu_{str \rightarrow f} \text{proj}[\mu_{s \rightarrow f}, \text{OccursAt}_p]), \end{aligned} \quad (15)$$

$$\begin{aligned} \mu_{f \rightarrow c}(c = 0) &= \sum_{str} \mu_{str \rightarrow f}(str) \sum_s \mu_{s \rightarrow f}(s) (1 - \mathbb{1}[str \text{ contains } s \text{ at pos } p]) = \\ &= \mathcal{Z}(\mu_{str \rightarrow f}) \mathcal{Z}(\mu_{s \rightarrow f}) - \mathcal{Z}(\mu_{str \rightarrow f} \text{proj}[\mu_{s \rightarrow f}, \text{OccursAt}_p]). \end{aligned} \quad (16)$$

4.3.3 Message to str

$$\begin{aligned}
\mu_{f \rightarrow str}(str) &= \mu_{c \rightarrow f}(c = 1) \sum_s \mu_{s \rightarrow f}(s) \mathbb{1}[str \text{ contains } s \text{ at pos } p] + \\
&\quad + \mu_{c \rightarrow f}(c = 0) \sum_s \mu_{s \rightarrow f}(s) (1 - \mathbb{1}[str \text{ contains } s \text{ at pos } p]) = \\
&= (\mu_{c \rightarrow f}(c = 1) - \mu_{c \rightarrow f}(c = 0)) \text{proj}[\mu_{s \rightarrow f}, \text{Occurs}\mathcal{A}t_p] + \mu_{c \rightarrow f}(c = 0) \text{Unif}. \quad (17)
\end{aligned}$$

4.3.4 Message to s

$$\begin{aligned}
\mu_{f \rightarrow s}(s) &= \mu_{c \rightarrow f}(c = 1) \sum_{str} \mu_{str \rightarrow f}(str) \mathbb{1}[str \text{ contains } s \text{ at pos } p] + \\
&\quad + \mu_{c \rightarrow f}(c = 0) \sum_{str} \mu_{str \rightarrow f}(str) (1 - \mathbb{1}[str \text{ contains } s \text{ at pos } p]) = \\
&= (\mu_{c \rightarrow f}(c = 1) - \mu_{c \rightarrow f}(c = 0)) \text{proj}[\mu_{str \rightarrow f}, \text{Occurs}\mathcal{A}t_p^T] + \mu_{c \rightarrow f}(c = 0) \text{Unif}. \quad (18)
\end{aligned}$$