

Tries

Construction, manipulation, and visualization of prefix trees

Anton Antonov

Mathematica for Prediction blog

Mathematica for Prediction project at GitHub

October 2013

November 2013

December 2013

Introduction

This document is a guide for using the functions for creation, manipulation, and visualization of prefix trees (tries) of the package `TriesWithFrequencies.m` provided by the project `MathematicaForPrediction` at GitHub; see [1].

The first section gives a brief discussion of what are prefix trees. See also the Wikipedia entry [2] for an alternative introduction and references. The second section provides basic examples of construction, search, and visualization with prefix trees. The third section gives a concrete example with a trie construction with the text of the U.S. constitution. The fourth section discusses data mining and machine learning applications. (To be written...) The last section provides details of the trie implementation in the package [1]. (To be written...)

Prefix trees are also known as “tries”. We are going to use “trie” and “prefix tree” as synonyms.

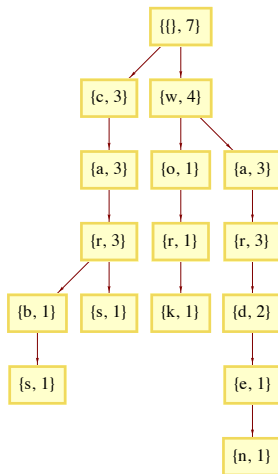
What are prefix trees (tries)?

In computer science a prefix tree (or a trie) is a data structure that allows quick retrieval of associations between a set of keys and a set of values. It is assumed that each key is a sequence of elements. (Most often the keys are considered to be strings i.e. sequences of characters.) The prefix tree nodes do not store the keys. Instead the position of a node in a prefix tree corresponds to a key.

Consider the following list of words:

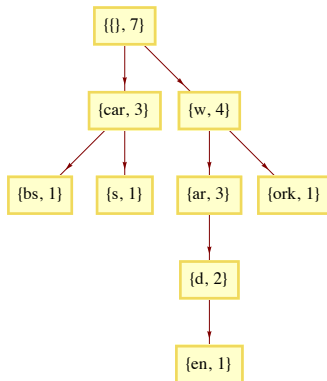
car
carbs
cars
war
ward
warden
work

Here is a prefix tree in which the keys are the list of words above and the values are the number of appearances of the corresponding prefix:

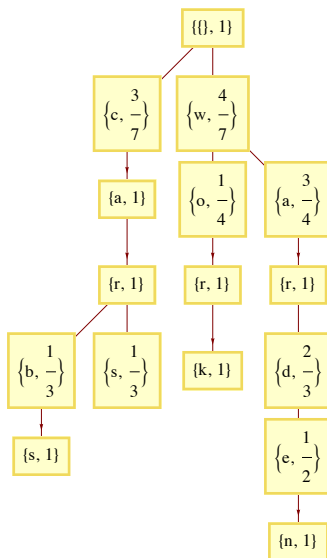


For example, the prefix “car” appears three times in the list of words; that is why we have path $\{(), 7\} \rightarrow \{“c”, 3\} \rightarrow \{“a”, 3\} \rightarrow \{“r”, 3\}$.

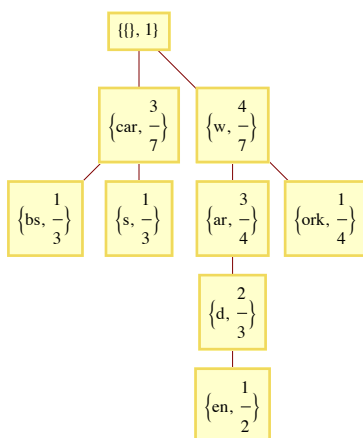
We can shrink the tree by merging the nodes that have only one child with their child. In this way we produce a tree in which each node would have an inseparable sub-sequences of elements. (In these examples these elements are characters.)



Another transformation of a trie is to replace the numbers of appearances with probabilities -- each value of a node is the probability to arrive to that node from the parent.



Of course we can do both transformations.



If we want to develop a trie package the most immediate trie operations to implement are: (i) trie creation from a list of words, (ii) finding the longest prefix of a word known to a trie, (iii) inserting a new word into a trie, (iv) finding the value corresponding to a key (i.e. word). The next section shows how to do with [1] these operations and the transformations described above. The package [1] implements trie operations and transformations for “words” that are strings or lists (of atom elements).

Tries creation and functionality

Let us see how we can construct tries using [1].

Get ["~/MathFiles/MathematicaForPrediction/TriesWithFrequencies.m"]

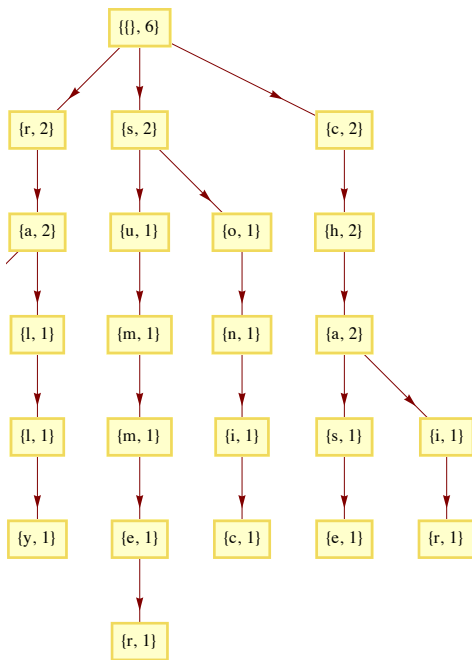
The functions of the package work with both strings and lists as keys. We are going to demonstrate functionalities using mostly strings, but also examples with lists would be shown.

First we are going to construct a simple trie.

```
mytrie = TrieCreate[{"race", "rally", "summer", "sonic", "chase", "chair"}]
{{{ {}, 6 }, {{ r, 2 }, {{ a, 2 }, {{ c, 1 }, {{ e, 1 } } }, {{ l, 1 }, {{ l, 1 }, {{ y, 1 } } } } } },
 {{ s, 2 }, {{ u, 1 }, {{ m, 1 }, {{ m, 1 }, {{ e, 1 }, {{ r, 1 } } } } } },
 {{ o, 1 }, {{ n, 1 }, {{ i, 1 }, {{ c, 1 } } } } } },
 {{ c, 2 }, {{ h, 2 }, {{ a, 2 }, {{ s, 1 }, {{ e, 1 } } }, {{ i, 1 }, {{ r, 1 } } } } } } }
```

Let us visualize the trie.

```
TrieForm[mytrie]
```



(Internally, the function `TrieForm` makes from its argument a list of rules that are suitable for visualization with `LayeredGraphPlot`.)

We can search for words in the trie with `TriePosition`.

```
pos = TriePosition[mytrie, "sun"]
{3, 2}
```

`TriePosition` returns the position of the subtree that has a root corresponding to the last character of the longest prefix of the word argument.

```
mytrie[[Sequence @@ pos]]
{{ u, 1 }, {{ m, 1 }, {{ m, 1 }, {{ e, 1 }, {{ r, 1 } } } } }
```

As in illustration of the prefix tree structure manipulation let us retrieve the prefix from the prefix tree (instead of just using `StringTake["sun", Length[pos]]`).

```
Map[mytrie[[Sequence@@Join[#, {1, 1}]]] &,
  FoldList[Append, {First[#]}, Rest[#]] &@TriePosition[mytrie, "sun"]]
{s, u}
```

Here is another example with a word known by the trie.

```
pos = TriePosition[mytrie, "summer"]
{3, 2, 2, 2, 2, 2}
```

Obviously, using the results `TriePosition` we can test which words are known by a trie. We can also use the function `TrieRetrieve` which returns the value corresponding to a key or an empty list if the key is not in the trie argument.

```
TrieRetrieve[mytrie, "summer"]
{r, 1}
```

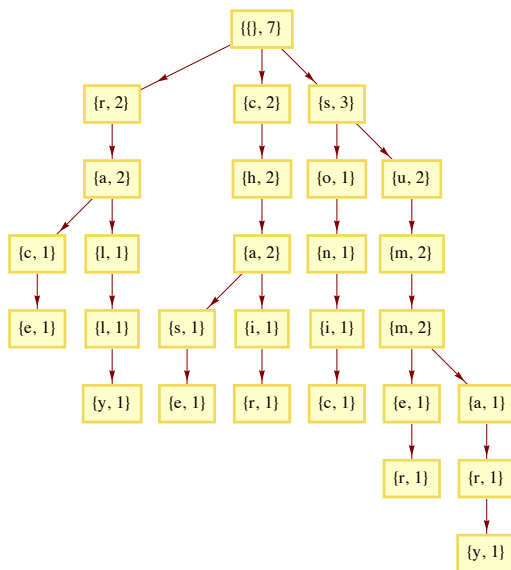
```
TrieRetrieve[mytrie, "serenity"]
{}
```

We can insert a new word in the trie with the command `TrieInsert`:

```
mytrie = TrieInsert[mytrie, "summary"]
{{{ {}, 7}, {{r, 2}, {{a, 2}, {{c, 1}, {{e, 1}}}, {{l, 1}, {{l, 1}, {{y, 1}}}}}},
  {{c, 2}, {{h, 2}, {{a, 2}, {{s, 1}, {{e, 1}}}, {{i, 1}, {{r, 1}}}}}},
  {{s, 3}, {{o, 1}, {{n, 1}, {{i, 1}, {{c, 1}}}}}, {{u, 2},
    {{m, 2}, {{m, 2}, {{e, 1}, {{r, 1}}}, {{a, 1}, {{r, 1}, {{y, 1}}}}}}}}
```

Here is the graph form of the updated trie:

```
TrieForm[mytrie]
```

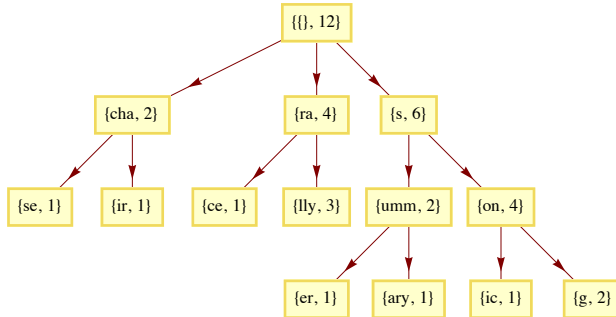


For the examples that follow let us insert several more words into the trie `mytrie`.

```
mytrie = Fold[TrieInsert, mytrie, {"son", "song", "song", "rally", "rally"}];
```

Let us shrink the trie and visualize the result

```
TrieForm[TrieShrink[mytrie]]
```

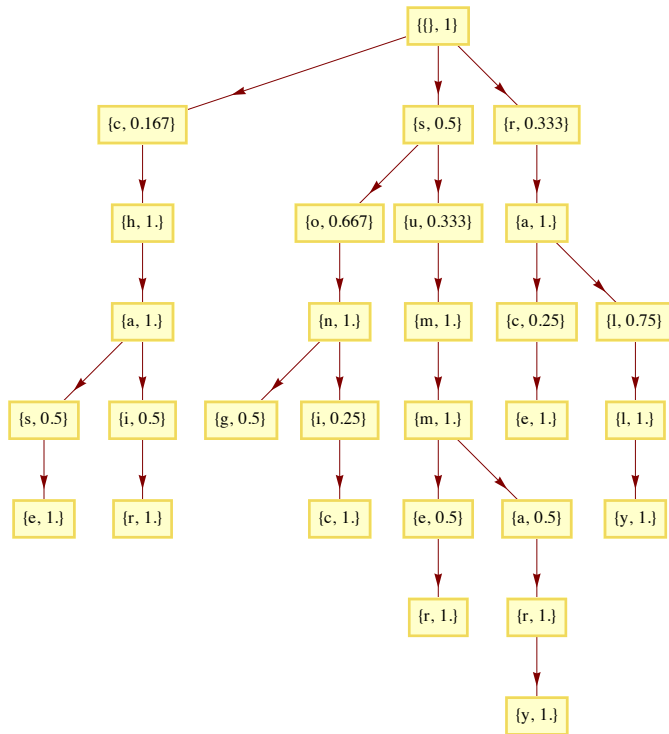


The function `TrieShrink` groups internal nodes and leaves into “prefixes”.

In order to replace the numbers of appearances with probabilities we use the function `TrieNodeProbabilities`.

```
mytriep =
  TrieNodeProbabilities[mytrie, "ProbabilityModifier" → (Round[#, 0.001] &)]
{{{ {}, 1}, {{c, 0.167},
  {{h, 1.}, {{a, 1.}, {{s, 0.5}, {{e, 1.}}}, {{i, 0.5}, {{r, 1.}}}}}},
 {{s, 0.5}, {{u, 0.333}, {{m, 1.},
  {{m, 1.}, {{e, 0.5}, {{r, 1.}}}, {{a, 0.5}, {{r, 1.}, {{y, 1.}}}}}},
 {{o, 0.667}, {{n, 1.}, {{i, 0.25}, {{c, 1.}}}, {{g, 0.5}}}}, {{r, 0.333},
 {{a, 1.}, {{c, 0.25}, {{e, 1.}}}, {{l, 0.75}, {{l, 1.}, {{y, 1.}}}}}}}
```

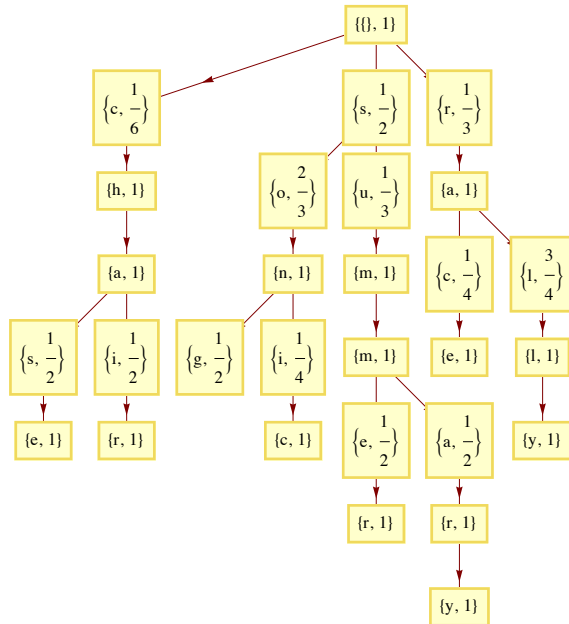
TrieForm[mytrie]



The function `TrieNodeProbabilities` takes the option "ProbabilityModifier" which can be used to change the appearance of the calculated probabilities. For example, we might prefer to see the probabilities as rational numbers

TrieNodeProbabilities[mytrie, "ProbabilityModifier" → Rationalize]

```
{{{}, 1}, {{c, 1/6}, {{h, 1}, {{a, 1}, {{s, 1/2}, {{e, 1}}}, {{i, 1/2}, {{r, 1}}}}}},
{{s, 1/2}, {{u, 1/3},
{{m, 1}, {{m, 1}, {{e, 1/2}, {{r, 1}}}, {{a, 1/2}, {{r, 1}, {{y, 1}}}}}},
{{o, 2/3}, {{n, 1}, {{i, 1/4}, {{c, 1}}}, {{g, 1/2}}}},
{{r, 1/3}, {{a, 1}, {{c, 1/4}, {{e, 1}}}, {{l, 3/4}, {{l, 1}, {{y, 1}}}}}}
```

TrieForm[%]

For data mining purposes we want to find the probabilities to arrive at the leaves from the root of a trie. The function `TrieLeafProbabilities` does that.

TrieLeafProbabilities[mytriep]

$$\left\{ \left\{ e, \frac{1}{12} \right\}, \left\{ r, \frac{1}{12} \right\}, \left\{ r, \frac{1}{12} \right\}, \left\{ y, \frac{1}{12} \right\}, \right. \\ \left. \left\{ c, \frac{1}{12} \right\}, \left\{ g, \frac{1}{6} \right\}, \left\{ n, \frac{1}{12} \right\}, \left\{ e, \frac{1}{12} \right\}, \left\{ y, \frac{1}{4} \right\} \right\}$$

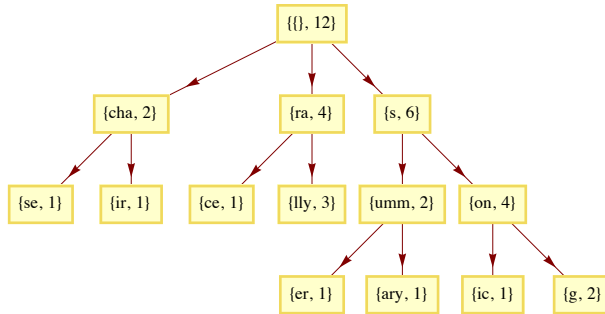
Let us verify that the sum of the probabilities is 1

Total[%[[All, 2]]]

1

Note that the internal node $\{n, \frac{1}{12}\}$ is an “internal leaf”: the sum of the probabilities of its children is smaller than 1. This can be seen in also in the shrunk frequencies trie -- we have the word “son” inserted in the trie.

TrieForm[TrieShrink[mytrie]]



Using lists

All the operations on tries for collections of strings also work for lists. This is demonstrated in this subsection with commands that mirror the ones used for tries of strings.

First let us convert a list of words into a list of integer lists.

```

lwords = Flatten /@Map[ToCharacterCode,
  Characters /@{"race", "rally", "summer", "sonic", "chase", "chair"}, {2}]
{{114, 97, 99, 101}, {114, 97, 108, 108, 121}, {115, 117, 109, 109, 101, 114},
 {115, 111, 110, 105, 99}, {99, 104, 97, 115, 101}, {99, 104, 97, 105, 114}}

```

Next is trie creation:

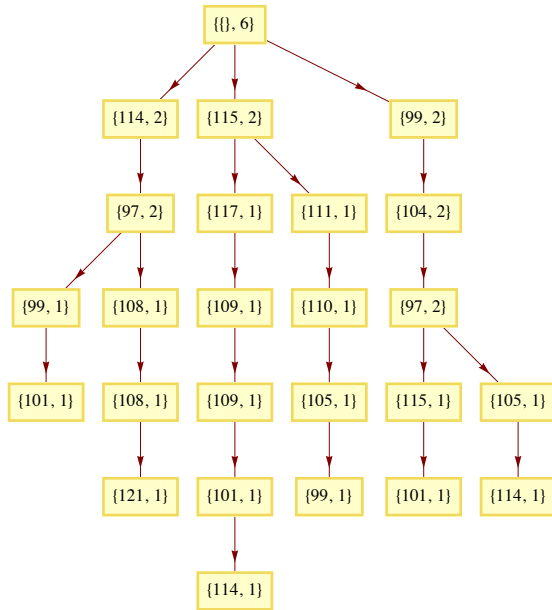
```

ltrie = TrieCreate[lwords]
{{{ }, 6}, {{114, 2},
  {{97, 2}, {{99, 1}, {{101, 1}}}, {{108, 1}, {{108, 1}, {{121, 1}}}}},
  {{115, 2}, {{117, 1}, {{109, 1}, {{109, 1}, {{101, 1}, {{114, 1}}}}}},
  {{111, 1}, {{110, 1}, {{105, 1}, {{99, 1}}}}}, {{99, 2},
  {{104, 2}, {{97, 2}, {{115, 1}, {{101, 1}}}, {{105, 1}, {{114, 1}}}}}}}}

```

Visualize the created trie:

TrieForm[ltrie]



For a list of integers find the position of the longest prefix known by the trie:

TriePosition[ltrie, {114, 97, 200}]

{2, 2}

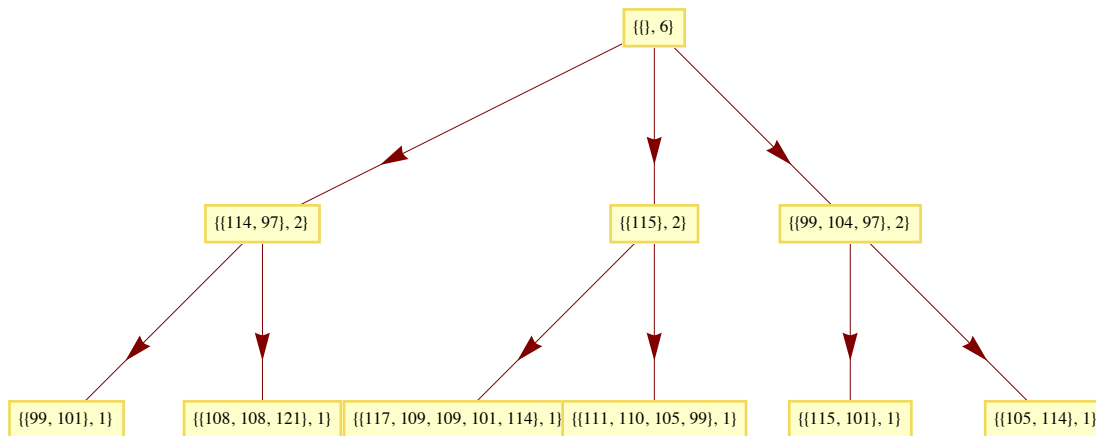
Shrink the trie:

sltrie = TrieShrink[ltrie]

```
{{{ {}, 6 }, {{{ {114, 97}, 2 }, {{{ {99, 101}, 1 }, {{{ {108, 108, 121}, 1 }}}},
  {{{ {115}, 2 }, {{{ {117, 109, 109, 101, 114}, 1 }, {{{ {111, 110, 105, 99}, 1 }}}},
  {{{ {99, 104, 97}, 2 }, {{{ {115, 101}, 1 }, {{{ {105, 114}, 1 }}}}}
```

Visualize the shrunk trie:

TrieForm[sltrie]



Find the probabilities to reach the leaves:

TrieLeafProbabilities[TrieNodeProbabilities[sltrie]] // Grid

{99, 101}	0.166667
{108, 108, 121}	0.166667
{117, 109, 109, 101, 114}	0.166667
{111, 110, 105, 99}	0.166667
{115, 101}	0.166667
{105, 114}	0.166667

Using tries as associative arrays

We can create tries that can be used as associative arrays (dictionaries). The function `TrieInsert` has also the signature `TrieInsert[trie_,key_,value_]`.

Let us make a trie in which the values corresponding to the inserted words are their reversals.

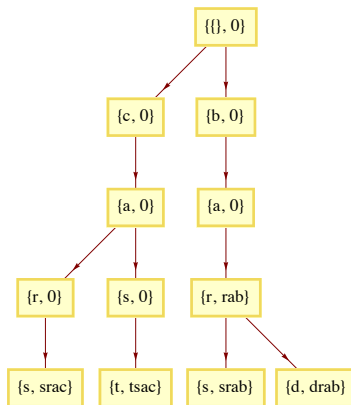
```

aaTrie = Fold[TrieInsert[#1, #2, StringReverse[#2]] &,
  TrieCreate[{}], {"cars", "cast", "bar", "bars", "bard"}]
{{{ {}, 0 }, {{c, 0}, {{a, 0}, {{r, 0}, {{s, srac}}}, {{s, 0}, {{t, tsac}}}}},
  {{b, 0}, {{a, 0}, {{r, rab}, {{s, srab}}, {{d, drab}}}}}

```

Let us visualize the trie -- note that the prefixes of the inserted words have values 0. For example, since we did not insert the word "car" its value in the trie is 0.

TrieForm[aaTrie]



Here are examples of retrieval of associated values:

TrieRetrieve[aaTrie, "car"]

{r, 0}

TrieRetrieve[aaTrie, "bars"]

{s, srab}

TrieRetrieve[aaTrie, "bard"]

{d, drab}

TrieRetrieve[aaTrie, "train"]

{}

Example with the U.S. constitution

Load some text:

```
text = Import["ExampleData/USConstitution.txt"];
```

Split the text into words:

```
words = ToLowerCase /@ StringSplit[text,  

  {" ", ".", ";", ":", "!", "?", " ", "-", "\n", "(", ")" }];
```

```
words //
```

```
Length
```

8719

Drop the words that are too short:

```
words = Select[words, StringLength[#] ≥ 1 &];
```

```
words // Length
```

7632

The number of unique words is

```
Length[Union[words]]
```

```
1190
```

Create a trie:

```
usctrie = TrieCreate[words];
```

Convert the trie nodes values into probabilities:

```
usctriep = TrieNodeProbabilities[usctrie];
```

Search for the prefix “universe”:

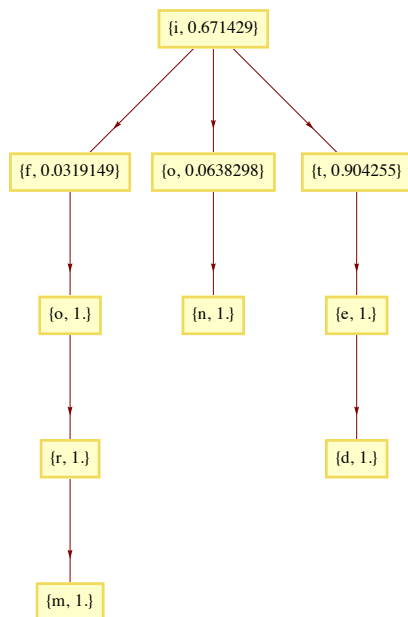
```
pos = TriePosition[usctrie, "universe"]
```

```
{29, 4, 7}
```

We can see that only the prefix “uni” is in the trie.

Here is the subtree of the trie with the last letter of “uni” as a root:

```
TrieForm[usctriep[[Sequence @@ pos]]]
```



Find the probabilities and positions of words that starts with “uni” and finish with “m”, “n”, and “d”:

```
probs = TrieLeafProbabilitiesWithPositions[usctriep[[Sequence @@ pos]]];
```

```
probs[[All, 2]] = probs[[All, 2]] / usctriep[[Sequence @@ pos, 1, 2]];
```

```
probs
```

```
{ {m, 0.0319149, {2, 2, 2, 2, 1}},
  {n, 0.0638298, {3, 2, 1}}, {d, 0.904255, {4, 2, 2, 1}} }
```

Here are the corresponding words:

```
Function[{ppos}, "uni" <> Apply[StringJoin,
  Map[usctrie[[Sequence @@ pos]] [Sequence @@ Join[#, {1, 1}]] &,
  FoldList[Append, {First[#]}, Rest[#]] & @ Most[ppos]]] /@ probs[[All, 3]]
{uniform, union, united}
```

Let us find the most frequently appearing words in the text. First we find the probabilities to reach each of the leafs of the trie:

```
probs = SortBy[TrieLeafProbabilitiesWithPositions[usctrie], -#[[2]] &];
probs // Length
1190
```

Here are the top 40 most frequent words in the text:

```
With[{ntop = 40}, t = Transpose[
  {Function[{ppos}, StringJoin @@ Map[usctrie[[Sequence @@ Join[#, {1, 1}]]] &,
    FoldList[Append, {First[#]}, Rest[#]] & @ Most[ppos]]] /@
    probs[[1 ;; ntop, 3]], probs[[1 ;; ntop, 2]]};
Magnify[Grid[Flatten /@ Transpose[Partition[t, Length[t] / 2]],
  Alignment -> Left, Spacings -> {{Automatic, Automatic, 2}, Automatic}], 0.8]]
```

the	0.0951258	not	0.0057652
of	0.0647275	may	0.0057652
shall	0.0400943	which	0.00563417
and	0.0345912	no	0.00550314
to	0.0264675	all	0.00537212
be	0.0234539	from	0.00537212
or	0.0209644	on	0.00511006
in	0.018999	law	0.00511006
states	0.0169025	office	0.00484801
president	0.0158543	vice	0.00471698
by	0.0131027	this	0.00458595
a	0.0127096	amendment	0.00458595
united	0.0111373	person	0.00445493
for	0.0110063	house	0.0043239
state	0.0103512	but	0.0043239
any	0.0103512	other	0.00406184
as	0.00838574	representatives	0.00379979
have	0.00825472	one	0.00379979
congress	0.00786164	he	0.00366876
such	0.00681342	article	0.00366876

Using tries in data mining

Details on trie visualization

References

- [1] Anton Antonov, Tries with frequencies *Mathematica* package, source code at GitHub, <https://github.com/antononcube/MathematicaForPrediction>, package *TriesWithFrequencies.m*, (2013).
- [2] Wikipedia, Trie, <http://en.wikipedia.org/wiki/Trie> .