

Introduction

Provides an overview of the project and a short discussion on the pertinent literature

The project number 8 is composed of two parts. One of them is perturb appropriately a dataset of documents, the second one rebuilding original documents.

I focused on the first part.

The literature in merit [1] suggests the most common errors observed and a methodology to identify such errors centered around the Needleman-Wunsch algorithm for identifying the best alignment of two strings. The original NW implementation was in the study of genomics and DNA sequences. [2]

This algorithm revolves around the idea of inserting characters '-' in smart places in order to obtain the best alignment of two strings.

The literature focuses more on finding what kind of errors OCR produce rather than cleverly reproduce them, hence making this work innovative in the sense that it allows to reproduce errors made by an optical reader into a larger dataset. The approach relies solely on (true_text, read_text) pairs.

Research question and methodology

Provides a clear statement on the goals of the project, an overview of the proposed approach, and a formal definition of the problem

I used a supervised method to learn what are the mistakes an OCR makes and reproduce them on any text.

These mistakes depend on

1. The OCR implementation
2. The dataset

Hence, given that Perturbers are made to replicate the mistakes they see in a ground truth, the kind of mistakes depend on the OCR technology and the dataset the Perturber was trained on. More, they apply a probabilistic approach (as it has no access to the original images) to replicate these.

The ground truth required is in the form (true_text, read_text). The goal is to identify the errors made and replicate them effectively on a fresh dataset. The idea is that a well trained perturber will be able to predict the errors it will make on an entire dataset by just a part of it.

It of course requires data to work and training on one dataset can give different results.

The approach is based on the structure of nchar (mimicking ngrams) (unichars and bichars in particular). The idea is to study what an nchar is transformed to by the scanner.

This very approach can be used in translation

Fitting and Transforming

Given a list of true documents x and a matching list of read documents y, fitting goes as follows:

```
1  Function Fit(x, y)
2      ----
3      Inputs:
4          x, y : lists of documents
5      Outputs:
6          perturbations : dictionary of dictionaries
7          p : map from keys(perturbations) to interval [0,1]
```

```

8      s : map from keys(perturbations) to integers
9      ----
10
11     perturbations, p, s <- empty dictionary
12     FOR ALL unichar uc in true
13         perturbations[uc] <- empty dictionary
14     FOR ALL bichars bc in true
15         perturbations[bc] <- empty dictionary
16
17     FOR i = 0 to N
18         true, scan <- Align(x[i], y[i]) // true and scan are of equal length
19
20         FOR j = 0 to length(true)
21             perturbations[true[j]][scan[j]] += 1 // add a perturbation of unichar
22             true[j]
23
24             FOR j = 0 to length(true)-1
25                 perturbations[true[j, j+1]][scan[j, j+1]] += 1 // add a perturbation
26                 of bichar true[j, j+1]
27
28             FOR ALL nchar in keys(perturbations)
29                 p[nchar] = //Probability that nchar is not perturbed to itself
30                 s[nchar] = //Number of times nchar was observed in x
31
32             return perturbations, p, s
33     End Function
34
35     Function Align(t1, t2)
36         ----
37         Inputs:
38             t1, t2 : strings
39         Outputs:
40             a1, a2 : strings of equal length
41         ----
42
43         a1, a2 <- Needlemanwunsch(t1, t2)
44
45         return a1, a2
46     End Function
47
48     Function ClusterFit(x, y)
49         ----
50         Inputs:
51             x, y : lists of strings
52             Vectorize : maps a dictionary to an array
53             Dictionarize : inverse of Vectorize
54             nch2idx : injective map from nchars to integers
55             Cluster : clustering method
56         Outputs:
57             perturbations: dictionary
58         ----
59
60         perturbations, _, _ <- Fit(x, y)
61         D = zeros(length(keys(perturbations)), length(keys(perturbations)))
62
63         // Build a distance matrix of keys of perturbations
64         FOR (nc1, nc2) in product(keys(perturbations), keys(perturbations))
65             id1, id2 <- nch2idx(nc1), nch2idx(nc2)
66             v1, v2 <- Vectorize(nc1), Vectorize(nc2)
67             d[id1, id2] <- Continuous_overlapping_distance(v1, v2)
68
69         // Cluster the keys of perturbations using D

```

```

66     clusters <- Cluster(D)
67
68     // Modify perturbations
69     FOR C in clusters
70         centroid = Centroid(C)
71         FOR nchar in C
72             v <- Vectorize(C)
73             s <- Sum(v) // keep track of the sum of v
74             d <- Continuous_overlapping_distance(v, centroid)
75             v = d*v + (1-d)*centroid // move v towards the centroid
76             perturbations[nchar] = Dictionarize(v) // rebuild perturbations[nchar]
77 End Function
78
79 Function Continuous_overlapping_distance(v1, v2)
80     ----
81     Inputs:
82         v1, v2 : numerical arrays of equal length
83     Outputs:
84         d : distance
85     ----
86
87     Rescale v1, v2 to sum 1
88     mins = [min(v1[t], v2[t]) FOR t=0 to length(v1)]
89     return 1 - Mean(mins)
90 End Function
91
92 Function Vectorize(d)
93     ----
94     Inputs:
95         d : dictionary with string keys
96         nch2idx : injective map from nchars to integers
97     Outputs:
98         v : array
99     ----
100
101     FOR key, value in d
102         id <- nch2idx(key)
103         v[id] <- value
104
105     return v
106 End Function
107
108 Function Transform(text)
109     ----
110     Inputs:
111         text : string
112         perturbations, p, s : output of a call to Fit on some training data
113         beta : float between 0 and 1
114         num_rounds : integer >= 1
115     Outputs:
116         out : string
117     ----
118
119     perturbed_list <- empty list
120     perturbed_text <- empty string
121
122     FOR round in num_rounds
123         splits <- Split(text, p)
124         FOR nchar in splits
125             APPEND Perturb(nchar, beta, perturbations) TO perturbed_list
126     perturbed_text = CONCATENATE perturbed_list

```

```

127 End Function
128
129 Function Split(text, p)
130     ----
131     Inputs:
132         text : string
133         p : // as in Transform
134     Outputs:
135         splits : list of nchars
136     ----
137
138     // Splits `text` in nchars to maximize the probability of at least one
139     // of these nchars to be perturbed
140     // An nchar is perturbed if Perturb(nchar, 1, perturbations) != nchar
141
142 End Function
143
144 Function Perturb(nchar, beta, perturbations)
145     ----
146     Inputs:
147         nchar, beta, perturbations : // as in Transform
148     Outputs:
149         pert : perturbed `nchar`
150     ----
151
152     // with probability `beta`, map `nchar` to `pert` according to the
    distribution
153     // given by `perturbations[nchar]`
154     // with probability 1-`beta`, return `nchar`
155
156 End Function

```

Toy example of Fit and Transform

`beta` is set to 1 and so is `num_rounds`

```

1  x1 = 'unimi'
2  y1 = 'un Im i'
3
4  true, scan = Align(x1, y1)
5  // un-im-i
6  // un Im i
7
8  perturbations, p, s = Fit(true, scan)
9  // perturbations
10 /*{
11     'i': {'I': 1, 'i': 1, 'I': 1, 'i': 1},
12     'im': {'Im': 1},
13     'm': {'m': 1, 'm ': 1},
14     'n': {'n': 1, 'n ': 1},
15     'u': {'u': 1},
16     'un': {'un': 1}
17 }*/
18 // p
19 // {'i': 0.75, 'im': 1.0, 'm': 0.5, 'n': 0.5, 'u': 0.0, 'un': 0.0}
20 // s
21 // {'i': 4, 'im': 1, 'm': 2, 'n': 2, 'u': 1, 'un': 1}

```

```

1  fresh = 'umido'
2
3  ls = split(fresh)
4  // ['u', 'm', 'i', 'd', 'o']
5
6  l <- empty list
7  FOR _=0 to 4
8      APPEND Transform(fresh) TO l
9  // ['umIdo', 'um ido', 'umIdo', 'um ido', 'um Ido']
10
11 /*
12 probability of perturbing fresh when split as _ is:
13     ['u', 'm', 'i', 'd', 'o'] : 0.875
14     ['um', 'id', 'o'] : 0
15     ['um', 'i', 'd', 'o'] : 0.75
16 */

```

Transforming

Perturbing a text is not simple.

Perturbing character by character can work but actually no, it does not catch some properties of 2chars

But how to choose whether to use a 1char or a 2char? Computing all of the possibilities requires $F(n)$ computation, where n is length of string. Instead, I used a greedy algorithm maximizing at each step the probability of having a misinterpretation.

If wrks like this on trichars:

take the first 3 characters

compute all of the possible splittings

find the easiest-to-misinterpret split

add it to a list

then perturb everything

Transformation

Finally transform the list of splitted, tokenized (in uni-bichars) text.

Each element of the list is mapped randomly (distribution given by what the algorithm was fitted on) to another uni or bichar.

Then the results are recombined in a singel string

Metrics for evaluation

- A perturber is trained on a training set. It perturbs a test set. The transformations on the test set are compared with the ground truh ones (having two other perturbers learning the transformations from the test set)

Experimental results

Provides an overview of the dataset used for experiments, the metrics used for evaluating performances, and the experimental methodology. Presents experimental results as plots and/or tables

Concluding remarks

Provides a critical discussion on the experimental results and some ideas for future work

Ideas:

- Explore the markov chain process of perturbation
- explore the possibility of using a similarity matrix of nchars when perturbing
 - In order to perturb never seen before nchars

References

1. https://link.springer.com/chapter/10.1007/978-3-030-45442-5_13#Bib1
2. <https://gist.github.com/slowkow/06c6dba9180d013dfd82bec217d22eb5>