

Dynamic Programming

Edit distance and its variants

Tyler Moore

CS 2123, The University of Tulsa

Some slides created by or adapted from Dr. Kevin Wayne. For more information see

<http://www.cs.princeton.edu/~wayne/kleinberg-tardos>. Some code reused from [Python Algorithms](#) by Magnus Lie

Hetland.

Edit distance

- Misspellings make approximate pattern matching an important problem
- If we are to deal with inexact string matching, we must first define a cost function telling us how far apart two strings are, i.e., a distance measure between pairs of strings.
- The edit distance is the minimum number of changes required to convert one string into another

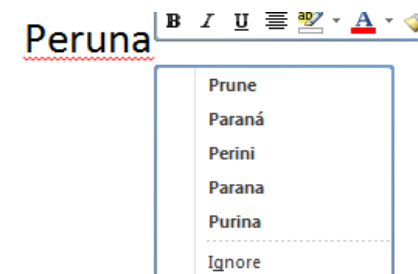
2 / 18

String edit operations

- We consider three types of changes to compute edit distance:
 - 1 Substitution: Change a single character from pattern s to a different character in text t , such as changing "shot" to "spot"
 - 2 Insertion: Insert a single character into pattern s to help it match text t , such as changing "ago" to "agog".
 - 3 Deletion: Delete a single character from pattern s to help it match text t , such as changing "hour" to "our"
- This definition of edit distance is also called Levenshtein distance
- Can you think of any other natural changes that might capture a single misspelling?

3 / 18

Edit distance application #1



- Spell checkers identify words in a dictionary with close edit distance to the misspelled word
- But how do they order the list of suggestions?

4 / 18

Towards a dynamic programming alternative

- We note that there are only $|s|$ possible values for i and $|t|$ possible values for j when invoking `edit_dist(i, j)` recursively
- This means there are at most $|s| \cdot |t|$ recursive function calls to cache in an iterative version
- The table is a two-dimensional matrix C where each of the $|s| \cdot |t|$ cells contains the cost of the optimal solution of this subproblem
- We just need a clever way to calculate the cost for each entry based on only a small subset of already-computed values.

9/18

Evaluation order

- To determine the value of cell (i, j) we need three values to already be computed: the cells $(i-1, j-1)$, $(i, j-1)$, and $(i-1, j)$.
- Any evaluation order with this property will do, including the row-major order used in the upcoming code

10/18

Edit distance: dynamic programming code

```
def iter_string_compare_lists(s, t):
    C, s, t = [], " " + s, " " + t #prepend empty character for edge case
    C.append(range(len(t)+1)) #initialize cost data structure
    for i in range(len(s)):
        C.append([i+1])
    for i in range(1, len(s)): #go through all characters of s
        for j in range(1, len(t)):
            #case 1: check for match at i and j
            if s[i]==t[j]: c_match = C[i-1][j-1]
            else: c_match = C[i-1][j-1]+1
            #case 2: there is an extra character to insert
            c_ins = C[i][j-1]+1
            #case 3: there is an extra character to remove
            c_del = C[i-1][j]+1
            c_min=min(c_match, c_ins, c_del)
            C[i].append(c_min)
    return C[i][j]
```

11/18

Edit distance: DP with cost table as dictionary

```
def iter_string_compare(s, t):
    C, s, t = {}, " " + s, " " + t #prepend empty character for edge case
    for j in range(len(t)): #initialize cost data structure
        C[0, j]=j
    for i in range(1, len(s)):
        C[i, 0]=i
    for i in range(1, len(s)): #go through all chars of s
        for j in range(1, len(t)):
            #case 1: check for match at i and j
            if s[i]==t[j]: c_match = C[i-1, j-1]
            else: c_match = C[i-1, j-1]+1
            #case 2: there is an extra character to insert
            c_ins = C[i, j-1]+1
            #case 3: there is an extra character to remove
            c_del = C[i-1, j]+1
            c_min=min(c_match, c_ins, c_del)
            C[i, j]=c_min
    return C[i, j]
```

12/18

Building edit distance cache

s: run
t: drain

C	-	d	r	a	i	n
-	0	← 1	← 2	← 3	← 4	← 5
r	↑ 1	↖ 1	↖ 1	← 2	← 3	← 4
u	↑ 2	↖ 2	↖ 2	↖ 2	← 3	← 4
n	↑ 3	↖ 3	↖ 3	↖ 3	↖ 3	↖ 3

Steps to turn “run” into “drain”

- 1 Insert d
- 2 Keep r
- 3 Substitute a for u
- 4 Insert i
- 5 Keep n

13 / 18

Edit distance exercises

- Build cost table by hand following DP algorithm
 - 1 s: bear, t: pea
 - 2 s: farm, t: for
- Performance cost of DP edit distance
 - Operations: $\Theta(|s| \cdot |t|)$
 - Storage: $\Theta(|s| \cdot |t|)$

14 / 18

Variation of edit distance: approximate substring matching

- Suppose we want to find the best close match to a smaller word in a larger string (e.g., find the closest match to “Tulsa” in “SMU Tulda Rice”)
- We need to modify our existing code in two ways
 - 1 Cost table initialization: all starting costs $C[0,j]$ should be set to 0
 - 2 Return the finishing cell $C[i,k]$ that minimizes the overall cost

15 / 18

Substring matching code

```
def iter_substring_match(s,t):
    C,s,t = {}," "+s," "+t #prepend empty character for edge cases
    for j in range(len(t)): #initialize cost data structure
        C[0,j]=0 #changed: ignore cost of preceding unmatched t
    for i in range(1,len(s)):
        C[i,0]=i
    for i in range(1,len(s)): #go through all chars of s
        for j in range(1,len(t)):
            #case 1: check for match at i and j
            if s[i]==t[j]: c_match = C[i-1,j-1]
            else: c_match = C[i-1,j-1]+1
            #case 2: there is an extra character to insert
            c_ins = C[i,j-1]+1
            #case 3: there is an extra character to remove
            c_del = C[i-1,j]+1
            c_min=min(c_match,c_ins,c_del)
            C[i,j]=c_min
    finj = min([(C[i,k],k) for k in range(1,len(t)-1)])
    return "with edit dist %i, %s morphs into %s finishing at %p"
```

16 / 18

Excercise: substring matching cache

s: Tulsa
t: SMU Tulda Rice

c	_	S	M	U	_	T	u	l	d	a	_	R	i	c	e
-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
T	↑1	↑1	↑1	↑1	↑1	↑0	↑1	↑1	↑1	↑1	↑1	↑1	↑1	↑1	↑1
u	↑2	↑2	↑2	↑2	↑2	↖1	↖0	←1	↑2	↑2	↑2	↑2	↑2	↑2	↑2
l	↑3	↑3	↑3	↑3	↑3	↑2	↑1	↖0	←1	←2	←3	↑3	↑3	↑3	↑3
s	↑4	↑4	↑4	↑4	↑4	↑3	↑2	↑1	↖1	←2	←3	↑4	↑4	↑4	↑4
a	↑5	↑5	↑5	↑5	↑5	↑4	↑3	↑2	↖2	↖1	←2	←3	←4	↑5	↑5

Substring ending at position 9 ("Tulda") is the closest substring to "Tulsa"

Variation of edit distance: longest common subsequence

- We might want to find the longest scattered sequence of characters within both strings
- For example, the longest common subsequence of "republican" and "democrat" is "eca"
- To get the longest subsequence, we can still allow insertions and deletions, but substitutions are forbidden
- We can change the edit distance code to behave as before on matches where the last characters are the same, but never select a substitution