

Handout 20: GDev and PCFGs

Grammar Development

1. The grammar developer

- a. It is an extended version of the HW 8 GDev class. Launching it:

```
1 $ python -m gdev g20
2 G>
```

- b. Start from a regression set of labeled sentences (`g20.sents`)

```
1 the dog barks
2 *the dogs barks
3 ...
```

- c. We also require a grammar: `g20.fcfg`

2. Commands

- a. Listing the sentences:

```
1 G> ss
2 => 0 OK the dog barks
3     1 * the dogs barks
4     ...
5     82 * you give the present to
```

- b. The current sentence

```
1 G> s
2 0 OK the dog barks
```

- c. Parsing it

```
1 Number of trees: 1
2
3 Tree 0
4 (Root[]
5   (S[]
6     (NP[f='sg']
7       (Det[] the)
8       (N2[f='sg'] (N1[f='sg'] (NC[f='sg'] (N[f='sg'] dog))))))
9     (VP[f='sg'] (V[f='sg', s=0, -t] barks))))
```

- d. Next sentence

```
1 G> n
2 1 * the dogs barks
```

- e. Going backwards

```
1 G> b
2 0 OK the dog barks
```

- f. Jumping to a particular sentence

```

1  G> 11
2    11 OK cats bark

```

3. Finding errors

- a. Sentences where parser prediction differs from sentence label

```

1  G> e
2    => 11 * (OK) cats bark
3      20 * (OK) Alice thinks that the white rabbit is in a hurry
4    ...

```

- b. First label is prediction, parenthesized is truth

4. Scoring

- a. Computing:

```

1  G> sc
2    Accuracy:    0.6904761904761905
3    Sensitivity: 0.46808510638297873
4    Specificity: 0.972972972972973

```

- b. Accuracy: proportion correct
- c. Errors of **sensitivity**: good sentence, but the grammar fails to parse it. Add rules.
- d. Errors of **specificity**: bad sentence, but the grammar accepts it. Add constraints (features).

5. Fixing errors

- a. If it parses, but shouldn't, examine the parse tree
- b. If it doesn't parse, look at the chart

```

1  G> c
2    cats
3      N[f='pl'] -> 'cats'
4      NC[f='pl'] -> N[f='pl']
5      N1[f='pl'] -> NC[f='pl']
6      N2[f='pl'] -> N1[f='pl']
7    bark
8      V[f='pl', s=0, -t] -> 'bark'
9      V[f='base', s=0, -t] -> 'bark'
10     VP[f='base'] -> V[f='base', s=0, -t]
11     VP[f='pl'] -> V[f='pl', s=0, -t]

```

- c. After editing grammar or sents, reload:

```

1  G> r

```

6. Examining vocabulary

a. Finding words that need to be added to the grammar:

```
1 G> unk
2 I
3 did
4 do
5 ...
```

b. Finding words that are not tested in any of the sentences:

```
1 G> unt
2 Frodo
3 Gertrude
4 ...
```

c. Listing parts of speech:

```
1 G> pos
2 Adj
3 Aux
4 Det
5 ...
```

7. Parsing detects errors of omission; generation, errors of commission.

a. Generating a sentence

```
1 >>> g.gen()
2 G> gen
3 0 Gertrude thinks about these aardvarks
4 1 will the first one brighter wine know
5 ...
6 9 Alice knew
```

b. If it is not clear how the string was generated, look at the tree:

```
1 G> t9
2 Number of trees: 10. Tree 9:
3 (Root[]
4   (S[]
5     (NP[f='sg'] (PropN[] Alice))
6     (VP[f='sg'] (V[f='sg', s=0, -t] knew))))
```

c. Interesting sentences may be added to the regression set. Keep a good sentence:

```
1 G> k0
```

d. Keep a sentence, but label it as bad:

```
1 G> *1
```

Probabilistic Parsing

8. Ambiguity

```
1 >>> p = load_parser('g1.cfg')
2 >>> trees = list(p.parse('Mary walked the cat in the park'.split()))
3 >>> print(trees[0])
4 (S
5   (NP (Name Mary))
6   (VP
7     (V walked)
8     (NP (Det the) (N cat))
9     (PP (P in) (NP (Det the) (N park)))))
10 >>> print(trees[1])
11 (S
12   (NP (Name Mary))
13   (VP
14     (V walked)
15     (NP (Det the) (N cat) (PP (P in) (NP (Det the) (N park))))))
```

9. Weighted PCFG

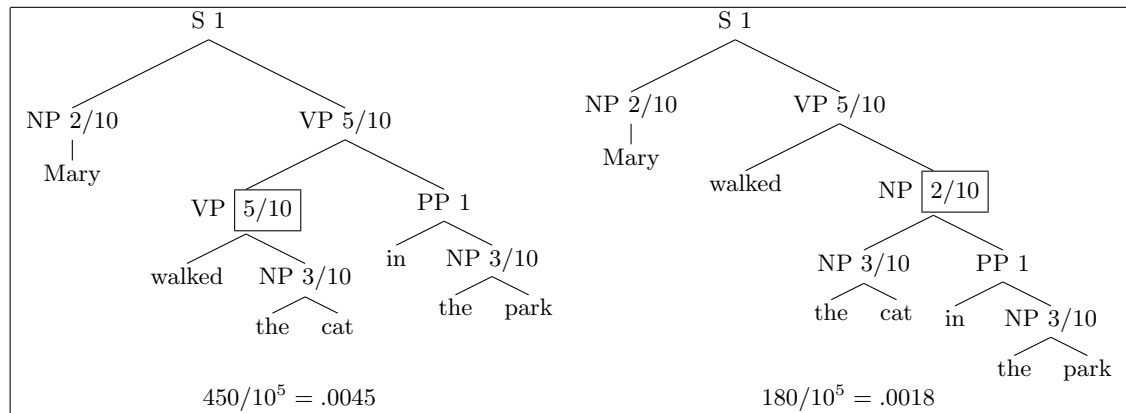
a. Example—put it in the file g2.pcfg:

```
1 % start S
2 S -> NP VP [1]
3 NP -> 'Mary' [.2]
4 NP -> 'the' 'cat' [.3]
5 NP -> 'the' 'park' [.3]
6 NP -> NP PP [.2]
7 PP -> 'in' NP [1]
8 VP -> 'walked' NP [.5]
9 VP -> VP PP [.5]
```

b. Parse

```
1 >>> p = load_parser('g2.pcfg')
2 >>> trees = list(p.parse('Mary walked the cat in the park'.split()))
3 >>> print(trees[0])
4 (S
5   (NP Mary)
6   (VP (VP walked (NP the cat)) (PP in (NP the park)))) (p=0.0045)
7 >>> print(trees[1])
8 (S
9   (NP Mary)
10  (VP walked (NP (NP the cat) (PP in (NP the park))))) (p=0.0018)
```

10. Compute the probabilities of the trees



11. Suppose we want the other parse on top

a. Change the VP weights (g2-alt.pcfg):

```
VP -> 'walked' NP [.9]
VP -> VP PP      [.1]
```

b. Now:

```
1  >>> p = load_parser('g2-alt.pcfg')
2  >>> trees = list(p.parse('Mary walked the cat in the park'.split()))
3  >>> print(trees[0])
4  (S
5   (NP Mary)
6   (VP walked (NP (NP the cat) (PP in (NP the park))))) (p=0.00324)
7  >>> print(trees[1])
8  (S
9   (NP Mary)
10  (VP (VP walked (NP the cat)) (PP in (NP the park)))) (p=0.00162)
```

12. Estimate a grammar from the following treebank:

