

# Optimizations and Extensions for Weighted CFG Parsers

## Practical Course in Natural Language Processing

**Boris Petrov**

Chair of Algorithmics  
Institute of Theoretical Computer Science  
TU Dresden

**2026-06-16**

# Outline

- 1 Pruning
  - Pruning for CKY parsing
  - Pruning for deductive parsing

- 2 A\*-parsing

# Pruning

- During CKY or deductive parsing many items are explored which are not part of the best derivation

# Pruning

- During CKY or deductive parsing many items are explored which are not part of the best derivation
- Idea: avoid items that are not part of the best derivation to speed up parsing

# Pruning

- During CKY or deductive parsing many items are explored which are not part of the best derivation
- Idea: avoid items that are not part of the best derivation to speed up parsing
- Problem: How can we know these items in advance?

# Pruning

- During CKY or deductive parsing many items are explored which are not part of the best derivation
- Idea: avoid items that are not part of the best derivation to speed up parsing
- Problem: How can we know these items in advance?
- Practical solution: Use simple methods but take the risk of finding suboptimal derivation.

# Pruning for CKY parsing

Consider this slightly modified (red) version of the CKY algorithm:

**Require:** weighted binary cfg  $(N, \Sigma, P, S, \mu)$ , word  $t_1 \dots t_n$  where  $t_1, \dots, t_n \in \Sigma$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid 0 \leq i < j \leq n, A \in N)$  such that, for all  $i, j, A$ ,

$$c_{i,j,A} = \max\{\mu(d) \mid d \in D_G^A(t_{i+1} \dots t_j)\} \cup \{0\}$$

```
1: function CKY( $P, \mu, t_1 \dots t_n$ )
2:   ( $c_{i,j,A} := 0 \mid 0 \leq i < j \leq n, A \in N$ )
3:   for  $1 \leq i \leq n$  do
4:     for  $A \rightarrow t_i \in P$  do
5:        $c_{i-1,i,A} := \max\{c_{i-1,i,A}, \mu(A \rightarrow t_i)\}$ 
6:     for  $2 \leq r \leq n$  do
7:       for  $0 \leq i \leq n - r$  do
8:          $j := i + r$ 
9:         for  $m \in \{i + 1, i + 2, \dots, j - 1\}$  do
10:          for  $B, C \in N$  do
11:            for  $A \in N$  such that  $A \rightarrow BC \in R$  do
12:               $c_{i,j,A} := \max\{c_{i,j,A}, \mu(A \rightarrow BC) \cdot c_{i,m,B} \cdot c_{m,j,C}\}$ 
13:   return  $c$ 
```

# Pruning for CKY parsing

An abstract “pruning operation“ on the table can be included (lines 6, 15).  
We speed up the algorithm by skipping the application of cfg rules to table cells with 0 probability (line 12).

**Require:** weighted binary cfg  $(N, \Sigma, P, S, \mu)$ , word  $t_1 \dots t_n$  where  $t_1, \dots, t_n \in \Sigma$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid 0 \leq i < j \leq n, A \in N)$  such that, for all  $i, j, A$ ,  $c_{i,j,A} \leq \max\{\mu(d) \mid d \in D_G^A(t_{i+1} \dots t_j)\} \cup \{0\}$

```
1: function CKY( $P, \mu, t_1 \dots t_n$ )
2:   ( $c_{i,j,A} := 0 \mid 0 \leq i < j \leq n, A \in N$ )
3:   for  $1 \leq i \leq n$  do
4:     for  $A \rightarrow t_i \in P$  do
5:        $c_{i-1,i,A} := \max\{c_{i-1,i,A}, \mu(A \rightarrow t_i)\}$ 
6:       ( $c_{i-1,i,A} \mid A \in N$ ) := prune( $(c_{i-1,i,A} \mid A \in N)$ )
7:       for  $2 \leq r \leq n$  do
8:         for  $0 \leq i \leq n - r$  do
9:            $j := i + r$ 
10:          for  $m \in \{i + 1, i + 2, \dots, j - 1\}$  do
11:            for  $B, C \in N$  do
12:              if  $c_{i,m,B} = 0$  or  $c_{m,j,C} = 0$  then continue
13:              for  $A \in N$  such that  $A \rightarrow BC \in R$  do
14:                 $c_{i,j,A} := \max\{c_{i,j,A}, \mu(A \rightarrow BC) \cdot c_{i,m,B} \cdot c_{m,j,C}\}$ 
15:                ( $c_{i,j,A} \mid A \in N$ ) := prune( $(c_{i,j,A} \mid A \in N)$ )
16:          return  $c$ 
```

# Pruning for CKY parsing

How can the pruning operation be implemented?

- **Threshold beam** (set all cell probabilities to 0 if worse than  $\theta \cdot$  best cell probability)

**Require:** family  $c = (c_{i,j,A} \in \mathbb{R} \mid A \in N)$ , threshold  $\theta \in [0, 1]$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid A \in N)$

```
1: function PRUNE( $c$ )
2:    $m = \max_{A \in N} \{c_{i,j,A} \mid A \in N\}$ 
3:   for  $A \in N$  do
4:     if  $c_{i,j,A} < m \cdot \theta$  then
5:        $c_{i,j,A} := 0$ 
6:   return  $c$ 
```

# Pruning for CKY parsing

How can the pruning operation be implemented?

- Threshold beam (set all cell probabilities to 0 if worse than  $\theta \cdot$  best cell probability)

**Require:** family  $c = (c_{i,j,A} \in \mathbb{R} \mid A \in N)$ , threshold  $\theta \in [0, 1]$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid A \in N)$

```
1: function PRUNE( $c$ )
2:    $m = \max_{A \in N} \{c_{i,j,A} \mid A \in N\}$ 
3:   for  $A \in N$  do
4:     if  $c_{i,j,A} < m \cdot \theta$  then
5:        $c_{i,j,A} := 0$ 
6:   return  $c$ 
```

- Fixed-sized beam (set all but  $n$  best cell probabilities to 0)

**Require:** family  $c = (c_{i,j,A} \in \mathbb{R} \mid A \in N)$ , size  $1 \leq n \leq |N|$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid A \in N)$

```
1: function PRUNE( $c$ )
2:    $[s_1, \dots, s_n] = n\text{-best}\{c_{i,j,A} \mid A \in N\}$ 
3:   for  $A \in N$  do
4:     if  $c_{i,j,A} < s_n$  then
5:        $c_{i,j,A} := 0$ 
6:   return  $c$ 
```

# Pruning for CKY parsing– implementation considerations

- No changes to data structures required.
- More speed-ups might be obtained by not adding items to the table which for sure would later be pruned:
  - Threshold beam: store weight  $m$  of currently best item. If new item has weight  $m'$  below  $\theta \cdot m$ , it is safe to prune immediately.
  - Fixed-size beam: store weights of the  $n$  best items. If the weight of new is below of worst item, prune immediately.

# Pruning for deductive parsing

## Original algorithm for deductive parsing.

**Require:** weighted binary cfg  $(N, \Sigma, P, S, \mu)$ , word  $t_1 \dots t_n$  where  $t_1, \dots, t_n \in \Sigma$

**Ensure:** family  $(c_{i,j,A} : \mathbb{R} \mid 0 \leq i < j \leq n, A \in N)$  such that  $c_{i,j,A} = \max\{\mu(d) \mid d \in D_G^A(t_{i+1} \dots t_j)\} \cup \{0\}$

```
1: function DEDUCE( $P, \mu, t_1 \dots t_n$ )
2:    $queue := \{(i-1, A, i, \mu(A \rightarrow t_i)) \mid 1 \leq i \leq n, A \rightarrow t_i \in P\}$ 
3:    $(c_{i,j,A} := 0 \mid 0 \leq i < j \leq n, A \in N)$ 
4:   while  $queue \neq \emptyset$  do
5:      $(i, A, j, w) := \operatorname{argmax}_{(i,A,j,w) \in queue} w$ 
6:      $queue \setminus = \{(i, A, j, w)\}$ 
7:     if  $c_{i,j,A} = 0$  then
8:        $c_{i,j,A} := w$ 
9:        $queue \cup = \{(i, A', j', \mu(A' \rightarrow AC) \cdot w \cdot c_{j,j',C}) \mid A' \rightarrow AC \in P\}$ 
10:       $queue \cup = \{(i', A', j, \mu(A' \rightarrow BA) \cdot c_{i',i,B} \cdot w) \mid A' \rightarrow BA \in P\}$ 
11:       $queue \cup = \{(i, A', j, \mu(A' \rightarrow A) \cdot w) \mid A' \rightarrow A \in P\}$ 
12:   return  $c$ 
```

# Pruning for deductive parsing

Pruning operation on queue is added.

**Require:** weighted binary cfg  $(N, \Sigma, P, S, \mu)$ , word  $t_1 \dots t_n$  where  $t_1, \dots, t_n \in \Sigma$

**Ensure:** family  $(c_{i,j,A} : \mathbb{R} \mid 0 \leq i < j \leq n, A \in N)$  such that  $c_{i,j,A} \leq \max\{\mu(d) \mid d \in D_G^A(t_{i+1} \dots t_j)\} \cup \{0\}$

```
1: function DEDUCE( $P, \mu, t_1 \dots t_n$ )
2:    $queue := \{(i-1, A, i, \mu(A \rightarrow t_i)) \mid 1 \leq i \leq n, A \rightarrow t_i \in P\}$ 
3:    $(c_{i,j,A} := 0 \mid 0 \leq i < j \leq n, A \in N)$ 
4:   while  $queue \neq \emptyset$  do
5:      $(i, A, j, w) := \operatorname{argmax}_{(i,A,j,w) \in queue} w$ 
6:      $queue \setminus= \{(i, A, j, w)\}$ 
7:     if  $c_{i,j,A} = 0$  then
8:        $c_{i,j,A} := w$ 
9:        $queue \cup= \{(i, A', j', \mu(A' \rightarrow AC) \cdot w \cdot c_{j,j',C}) \mid A' \rightarrow AC \in P\}$ 
10:       $queue \cup= \{(i', A', j, \mu(A' \rightarrow BA) \cdot c_{i',i,B} \cdot w) \mid A' \rightarrow BA \in P\}$ 
11:       $queue \cup= \{(i, A', j, \mu(A' \rightarrow A) \cdot w) \mid A' \rightarrow A \in P\}$ 
12:      prune( $queue$ )
13:   return  $c$ 
```

# Pruning for deductive parsing

Again two options for pruning:

- **Threshold beam** (remove each queue item if its probability is worse than  $\theta \cdot \text{prob. of best queue item}$  )

**Require:** set  $queue \subseteq \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{R}$ , threshold  $\theta \in [0, 1]$

**Ensure:** set  $queue' \subseteq \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{R}$

- 1: **function** PRUNE( $queue$ )
- 2:  $m = \max_{(i,A,j,w) \in queue} w$
- 3: **return**  $\{(i, A, j, w) \in queue \mid w > \theta \cdot m\}$

# Pruning for deductive parsing

Again two options for pruning:

- **Threshold beam** (remove each queue item if its probability is worse than  $\theta \cdot \text{prob. of best queue item}$  )

**Require:** set  $queue \subseteq \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{R}$ , threshold  $\theta \in [0, 1]$

**Ensure:** set  $queue' \subseteq \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{R}$

- 1: **function** PRUNE( $queue$ )
- 2:  $m = \max_{(i,A,j,w) \in queue} w$
- 3: **return**  $\{(i, A, j, w) \in queue \mid w > \theta \cdot m\}$

- **Fixed-sized beam** (only keep the  $n$  most probable queue items)

**Require:** set  $queue \subseteq \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{R}$ , size  $n \in \mathbb{N}$

**Ensure:** set  $queue' \subseteq \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{R}$

- 1: **function** PRUNE( $queue$ )
- 2:  $[i_1, \dots, i_n] = n\text{-best}(queue)$  w.r.t. 4th tuple component
- 3: **return**  $\{i_1, \dots, i_n\}$

## Pruning for deductive parsing – implementation considerations

- Best implement queue as Min-max heap [Atk+86], because access to best and worst elements is required.

## Pruning for deductive parsing – implementation considerations

- Best implement queue as Min-max heap [Atk+86], because access to best and worst elements is required.
- Again, don't add items to queue if they would be pruned immediately.

## Pruning for deductive parsing – implementation considerations

- Best implement queue as Min-max heap [Atk+86], because access to best and worst elements is required.
- Again, don't add items to queue if they would be pruned immediately.
- Alternatively, one can shrink the queue only occasionally and not in each iteration of the main loop.

## Pruning for deductive parsing – implementation considerations

- Best implement queue as Min-max heap [Atk+86], because access to best and worst elements is required.
- Again, don't add items to queue if they would be pruned immediately.
- Alternatively, one can shrink the queue only occasionally and not in each iteration of the main loop.
  
- Beware: Items for small spans are often more probable than items for large spans.  
Risk of pruning “good” large items in favour of “bad” small items.

## Pruning for deductive parsing – implementation considerations

- Best implement queue as Min-max heap [Atk+86], because access to best and worst elements is required.
- Again, don't add items to queue if they would be pruned immediately.
- Alternatively, one can shrink the queue only occasionally and not in each iteration of the main loop.
  
- Beware: Items for small spans are often more probable than items for large spans.  
Risk of pruning “good” large items in favour of “bad” small items.  
(Solution: see A\*-star parsing below)

# Outline

- 1 Pruning
  - Pruning for CKY parsing
  - Pruning for deductive parsing

- 2 A\*-parsing

## A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .

## A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .
- How about *future costs*, i.e., the weight of  $S \Rightarrow_G^* t_1 \cdots t_i A t_{j+1} \cdots t_n$ ?

## A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .
- How about *future costs*, i.e., the weight of  $S \Rightarrow_G^* t_1 \cdots t_i \ A \ t_{j+1} \cdots t_n$ ?
- If future costs are taken into account, then maybe less items from the queue need to be processed.

## A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .
- How about *future costs*, i.e., the weight of  $S \Rightarrow_G^* t_1 \cdots t_i \ A \ t_{j+1} \cdots t_n$ ?
- If future costs are taken into account, then maybe less items from the queue need to be processed.
  - Why?: Usually items with small spans are more probable than items with large spans.

# A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .
- How about *future costs*, i.e., the weight of  $S \Rightarrow_G^* t_1 \cdots t_i \ A \ t_{j+1} \cdots t_n$ ?
- If future costs are taken into account, then maybe less items from the queue need to be processed.
  - Why?: Usually items with small spans are more probable than items with large spans.
  - This is counteracted by future costs which are higher for items with small spans.

## A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .
- How about *future costs*, i.e., the weight of  $S \Rightarrow_G^* t_1 \cdots t_i \ A \ t_{j+1} \cdots t_n$ ?
- If future costs are taken into account, then maybe less items from the queue need to be processed.
  - Why?: Usually items with small spans are more probable than items with large spans.
  - This is counteracted by future costs which are higher for items with small spans.
- Klein and Manning [KMO3] propose several admissible heuristics.

## A\*-parsing

- weighted deductive parsing computes for each item  $(i, j, A)$  in the table the weight of the most probable derivation from  $A$  to  $t_{i+1} \cdots t_j$ .
- How about *future costs*, i.e., the weight of  $S \Rightarrow_G^* t_1 \cdots t_i \ A \ t_{j+1} \cdots t_n$ ?
- If future costs are taken into account, then maybe less items from the queue need to be processed.
  - Why?: Usually items with small spans are more probable than items with large spans.
  - This is counteracted by future costs which are higher for items with small spans.
- Klein and Manning [KMO3] propose several admissible heuristics.
- A heuristic may also be useful when pruning items during CKY parsing.

# A\*-parsing – Viterbi outside score

We use the admissible heuristic out:

$$\text{out}(A) = \max_{d \in D_G, u, w \in \Sigma^* : S \xrightarrow{d}_G uAw} \text{weight}(d)$$

It can be computed by a variant of the inside/outside algorithm:

```
1: function INSIDE
2:   for  $A \in N$  do
3:      $\text{in}(A) := \max(\{\mu(A \rightarrow \alpha) \mid A \rightarrow \alpha \in R\} \cup \{0\})$ 
4:   while not converged do
5:     for  $A \in N$  do
6:        $\text{in}(A) = \max(\{\text{in}(A)\} \cup \{\mu(A \rightarrow BC) \cdot \text{in}(B) \cdot \text{in}(C) \mid A \rightarrow BC \text{ in } R\}$   

          $\cup \{\mu(A \rightarrow B) \cdot \text{in}(B) \mid A \rightarrow B \text{ in } R\})$ 
7: function OUTSIDE
8:   set  $\text{out}(B) := \begin{cases} 1 & B = S \\ 0 & \text{otherwise} \end{cases}$  for each  $B \in N$ 
9:   while not converged do
10:  for  $B \in N$  do
11:     $\text{out}(B) := \max(\{\text{out}(B)\} \cup \{\text{out}(A) \cdot \mu(A \rightarrow BC) \cdot \text{in}(C) \mid A \rightarrow BC \text{ in } R\}$   

       $\cup \{\text{out}(A) \cdot \mu(A \rightarrow CB) \cdot \text{in}(C) \mid A \rightarrow CB \text{ in } R\}$   

       $\cup \{\text{out}(A) \cdot \mu(A \rightarrow B) \mid A \rightarrow B \text{ in } R\})$ 
```

# A\*-parsing – parsing algorithm with heuristic

The out-value is now simply included when selecting the best item from the queue:

**Require:** weighted binary cfg  $(N, \Sigma, P, S, \mu)$ , word  $t_1 \dots t_n$  where  $t_1, \dots, t_n \in \Sigma$

**Ensure:** family  $(c_{i,j,A} : \mathbb{R} \mid 0 \leq i < j \leq n, A \in N)$  such that  $c_{i,j,A} = \max\{\mu(d) \mid d \in D_G^A(t_{i+1} \dots t_j)\} \cup \{0\}$

```
1: function DEDUCE( $P, \mu, t_1 \dots t_n$ )
2:    $queue := \{(i-1, A, i, \mu(A \rightarrow t_i)) \mid 1 \leq i \leq n, A \rightarrow t_i \in P\}$ 
3:    $(c_{i,j,A} := 0 \mid 0 \leq i < j \leq n, A \in N)$ 
4:   while  $queue \neq \emptyset$  do
5:      $(i, A, j, w) := \operatorname{argmax}_{(i,A,j,w) \in queue} w \cdot \text{out}(A)$ 
6:      $queue \setminus= \{(i, A, j, w)\}$ 
7:     if  $c_{i,j,A} = 0$  then
8:        $c_{i,j,A} := w$ 
9:        $queue \cup= \{(i, A', j', \mu(A' \rightarrow AC) \cdot w \cdot c_{j,j',C}) \mid A' \rightarrow AC \in P\}$ 
10:       $queue \cup= \{(i', A', j, \mu(A' \rightarrow BA) \cdot c_{i',i,B} \cdot w) \mid A' \rightarrow BA \in P\}$ 
11:       $queue \cup= \{(i, A', j, \mu(A' \rightarrow A) \cdot w) \mid A' \rightarrow A \in P\}$ 
12:   return  $c$ 
```

# Pruning for deductive parsing – in the context of CKY

Can A\*-parsing be applied for CKY parsing?

Yes: in combination with pruning:

- Threshold beam

**Require:** family  $c = (c_{i,j,A} \in \mathbb{R} \mid A \in N)$ , threshold  $\theta \in [0, 1]$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid A \in N)$

```
1: function PRUNE(c)
2:    $m = \max_{A \in N} \{c_{i,j,A} \cdot \text{out}(A) \mid A \in N\}$ 
3:   for  $A \in N$  do
4:     if  $c_{i,j,A} \cdot \text{out}(A) < m \cdot \theta$  then
5:        $c_{i,j,A} := 0$ 
6:   return c
```

# Pruning for deductive parsing – in the context of CKY

Can A\*-parsing be applied for CKY parsing?

Yes: in combination with pruning:

- Threshold beam

**Require:** family  $c = (c_{i,j,A} \in \mathbb{R} \mid A \in N)$ , threshold  $\theta \in [0, 1]$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid A \in N)$

```
1: function PRUNE(c)
2:    $m = \max_{A \in N} \{c_{i,j,A} \cdot \text{out}(A) \mid A \in N\}$ 
3:   for  $A \in N$  do
4:     if  $c_{i,j,A} \cdot \text{out}(A) < m \cdot \theta$  then
5:        $c_{i,j,A} := 0$ 
6:   return c
```

- Fixed-size beam

**Require:** family  $c = (c_{i,j,A} \in \mathbb{R} \mid A \in N)$ , size  $1 \leq n \leq |N|$

**Ensure:** family  $(c_{i,j,A} \in \mathbb{R} \mid A \in N)$

```
1: function PRUNE(c)
2:    $[s_1, \dots, s_n] = n\text{-best}\{c_{i,j,A} \cdot \text{out}(A) \mid A \in N\}$ 
3:   for  $A \in N$  do
4:     if  $c_{i,j,A} \cdot \text{out}(A) < s_n$  then
5:        $c_{i,j,A} := 0$ 
6:   return c
```

- [Atk+86] M. D. Atkinson et al. “Min-max Heaps and Generalized Priority Queues”. In: *Commun. ACM* 29.10 (Oct. 1986), pp. 996–1000. ISSN: 0001-0782. DOI: 10.1145/6617.6621. URL: <http://doi.acm.org/10.1145/6617.6621>.
- [KM03] Dan Klein and Christopher D. Manning. “A\* Parsing: Fast Exact Viterbi Parse Selection”. In: *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*. 2003, pp. 119–126. URL: <https://www.aclweb.org/anthology/N03-1016>.