# Chomsky normal form

Back (for review), by popular demand

# What is Chomsky Normal Form?

It's a way to express the rules of a CFG

Every CFG may be written in normal form

# What is the structure of Chomsky Normal Form?

CFG is in Chomsky normal form if every rule takes form:

$$A \rightarrow BC$$

$$A \rightarrow a$$

- B and C may not be the start variable

- Only the start variable can transition to $\varepsilon$

- Each variable goes to two other variables or two one terminal

- The start variable may not point back to the start variable

# Consider a typical CFG

- Variables and terminals mix          $A \rightarrow xBy$

- Some variables point to other single variables

    $A \rightarrow C$

- Start variable can point to itself      $S \rightarrow SS \mid y$

- Any variable can transition to $\varepsilon$      $B \rightarrow \varepsilon$

# Converting to Chomsky Normal Form

- $S_0 \rightarrow S$ where $S$ was original start variable

- Remove $A \rightarrow \varepsilon$

<span style="color:red">We won't use this rule in this class, will use all others</span>

- For each multiple-occurrence of A, add new rules with A deleted

    $R \rightarrow uAvAw$    change to  $R \rightarrow uvAw \mid uAvw \mid uvw$

- Shortcut all unit rules

    Given $A \rightarrow B$ and $B \rightarrow u$ , add $A \rightarrow u$

- Replace rules $A \rightarrow u_1 u_2 u_3 \dots u_k$ with:

    $A \rightarrow u_1 A_1, A_1 \rightarrow u_2 A_2, A_2 \rightarrow u_3 A_3, \dots, A_{k-2} \rightarrow u_{k-1} u_k$

# Let's Chomsky-ize a non-Chomsky form grammar

G1:

S -> AB

A -> cAn | c

B -> BB | n

**Replace terminals-variable mixes with variables only**

S -> AB

A -> $U_c$A$U_n$ | c

B -> BB | n

**$U_c$ -> c**

**$U_n$ -> n**

**Convert 3-variable rules to 2-variable rules**

S -> AB

A -> **$U_c$D** | c

B -> BB | n

**$U_c$ -> c**

**$U_n$ -> n**

**D -> A$U_n$**

More on G1:

G1:

S -> AB

A -> $U_c$D | c

B -> BB | n

$U_c$ -> c

$U_n$ -> n

D -> A$U_n$

This already fits normal form!

Typical to add new start state

$S_0$ -> S

S -> AB

A -> $U_c$D | c

B -> BB | n

$U_c$ -> c

$U_n$ -> n

D -> A$U_n$

Replace S in $S_0$ -> S rule

Final answer!

$S_0$ -> AB

S -> AB

A -> $U_c$D | c

B -> BB | n

$U_c$ -> c

$U_n$ -> n

D -> A$U_n$

Reminder:

CFG is in Chomsky normal form if every rule takes form:

$$A \rightarrow BC$$

$$A \rightarrow a$$

- B and C may not be the start variable
- Only the start variable can transition to $\varepsilon$

- Each variable goes to two other variables or two one terminal
- The start variable may not point back to the start variable

# Let's get more complicated with grammar G2

**Let's add new start state first this time:**

G2:

S -> AB

A -> cAn | c | $\varepsilon$

B -> BB | n

$S_0$ **-> S**

S -> AB

A -> cAn | c | $\varepsilon$

B -> BB | n

$S_0$ **-> S**

S -> AB | **$\varepsilon$B**

A -> cAn | c | $\varepsilon$ | c**$\varepsilon$**n

B -> BB | n

**To remove $\varepsilon$, first plug it in wherever it applies**

More on G2

**Finish removing ε**

$S_0 \to S$

$S \to AB \mid$ **B**

$A \to cAn \mid c \mid$ **cn**

$B \to BB \mid n$

$S_0 \to S$

$S \to AB \mid \boldsymbol{\varepsilon B}$

$A \to cAn \mid c \mid \varepsilon \mid \boldsymbol{c\varepsilon n}$

$B \to BB \mid n$

**Replace terminals-variable mixes with variables only**

$S_0 \to S$

$S \to AB \mid$ **B**

$A \to \mathbf{U_c}\, A\, \mathbf{U_n} \mid c \mid \mathbf{U_c}\,\mathbf{U_n}$

$B \to BB \mid n$

$\mathbf{U_c \to c}$

$\mathbf{U_n \to n}$

# More on G2

$S_0 \rightarrow S$

$S \rightarrow AB \mid B$

$A \rightarrow U_c A U_n \mid c \mid U_c U_n$

$B \rightarrow BB \mid n$

$U_c \rightarrow c$

$U_n \rightarrow n$


$S_0 \rightarrow S$

$S \rightarrow AB \mid B$

$A \rightarrow U_c D \mid c \mid U_c U_n$

$B \rightarrow BB \mid n$

$U_c \rightarrow c$

$U_n \rightarrow n$

$D \rightarrow A U_n$


$S_0 \rightarrow AB \mid BB \mid n$

$S \rightarrow AB \mid BB \mid n$

$A \rightarrow U_c D \mid c \mid U_c U_n$

$B \rightarrow BB \mid n$

$U_c \rightarrow c$

$U_n \rightarrow n$

$D \rightarrow A U_n$

# G2 final answer

$S_0 \to AB \mid BB \mid n$

$S \to AB \mid BB \mid n$

$A \to U_c D \mid c \mid U_c U_n$

$B \to BB \mid n$

$U_c \to c$

$U_n \to n$

$D \to AU_n$

**Produces same CFL as:**

$S \to AB$

$A \to cAn \mid c \mid \varepsilon$

$B \to BB \mid n$

Reminder:

CFG is in Chomsky normal form if every rule takes form:

$$A \rightarrow BC$$

$$A \rightarrow a$$

- B and C may not be the start variable
- Only the start variable can transition to $\varepsilon$

- Each variable goes to two other variables or two one terminal
- The start variable may not point back to the start variable

# Let's get even <u>more</u> complicated

G3:

S -> AB | BS

A -> cAn | c | $\varepsilon$

B -> BB | n

**Add new start state:**

**$S_0$ -> S**

S -> AB | BS

A -> cAn | c | $\varepsilon$

B -> BB | n

**Remove $\varepsilon$**

**$S_0$ -> S**

S -> AB | BS | **B**

A -> cAn | c | **cn**

B -> BB | n

# More on G3

**Replace terminals-variable mixes with variables only**

**$S_0$ -> S**

S -> AB | BS | **B**

A -> cAn | c | **cn**

B -> BB | n

$\longrightarrow$

**$S_0$ -> S**

S -> AB | BS | **B**

A -> **$U_c$A$U_n$** | c | **$U_c$$U_n$**

B -> BB | n

**$U_c$ -> c**

**$U_n$ -> n**

# More on G3

**Replace single variables on right side (S, B)**

$S_0 \rightarrow S$

$S \rightarrow AB \mid BS \mid B$

$A \rightarrow U_c A U_n \mid c \mid U_c U_n$

$B \rightarrow BB \mid n$

$U_c \rightarrow c$

$U_n \rightarrow n$

$S_0 \rightarrow AB \mid BS \mid BB \mid n$

$S \rightarrow AB \mid BS \mid BB \mid n$

$A \rightarrow U_c A U_n \mid c \mid U_c U_n$

$B \rightarrow BB \mid n$

$U_c \rightarrow c$

$U_n \rightarrow n$

# More on G3

**Convert 3-variable rules to 2-variable rules**

$S_0$ -> **AB | BS | BB | n**

S -> AB | BS | **BB | n**

A -> $U_c$A$U_n$ | c | $U_c U_n$

B -> BB | n

$U_c$ -> c

$U_n$ -> n

---

$S_0$ -> **AB | BS | BB | n**

S -> AB | BS | **BB | n**

A -> $U_c$**D** | c | $U_c U_n$

B -> BB | n

$U_c$ -> **c**

$U_n$ -> **n**

**D -> A$U_n$**

# Reminder:

CFG is in Chomsky normal form if every rule takes form:

$$A \rightarrow BC$$

$$A \rightarrow a$$

- B and C may not be the start variable
- Only the start variable can transition to $\varepsilon$

- Each variable goes to two other variables or two one terminal
- The start variable may not point back to the start variable