

# Feature Engineering for Data-driven Dependency Parsing

Christian Rishøj Jensen  
August 2009

# Errata

	Presence			
	% of sentences with...			
Corpus	“	()	:	-
Law	0.0%	21.0%	2.0%	0.0%
Law2	<del>0.0%</del>	<del>21.0%</del>	<del>2.0%</del>	<del>0.0%</del>
Newspaper	14.0%	7.0%	11.0%	9.0%
Law+News	<del>14.0%</del>	<del>28.0%</del>	<del>13.0%</del>	<del>9.0%</del>
Both	7.0%	14.0%	6.5%	4.5%

“Presence” metrics for **Italian**  
(table 3.2)

# Errata

- **Bulgarian** (truncated treebank):
  - **Parenthesis** — 5% of sentences<sub>test</sub>  
*tiny* **adverse** effect [ $p \approx 0.49$ ]
  - **Colon** — 2% of sentences<sub>test</sub>  
LAS<sub>overall</sub> **-0.36** [ $p \approx 0.12$ ]



# Errata

- **Slovene (SDT):**
  - **Colon** — 2% of sentences  
LAS<sub>affected</sub> **+1.39** [ $p \approx 0.26$ ]
  - **DashApposition** — 2% of sentences  
LAS<sub>affected</sub> **+1.85** [ $p \approx 0.09$ ]

# Errata

- **Danish (DDT):**
  - **Colon** — 7% of sentences<sub>test</sub>  
LAS<sub>affected</sub> **+1.42** [ $p \approx 0.12$ ]
  - **DashApposition** — 0.3% (!)  
LAS<sub>overall</sub> **+0.34** [ $p \approx 0.10$ ]

# Errata

- **Spanish** (Cast3lb):
  - **Colon** — 5% of sentences<sub>test</sub>  
LAS<sub>affected</sub> **+1.11** [ $p \approx 0.23$ ]



# Unused Treebanks

- **BIO** (English)
- **Metu** (Turkish)
- **Alpino** (Dutch)
- **PDT** (Czech)
- **Tiger** (German)
- **Law2** (Italian, revised)

# Treebank Selection

- **Objective** selection criteria
  - Not explicitly stated
- Challenges:
  - **Presence** of linguistic phenomena
  - **Impact** on parsing accuracy



# Future Treebank Selection

- Gather data:
  - Treebank **characteristics**
  - Augmentation scheme **effects**
- **Automatically** determine applicable augmentations