

## README

IL-POST Data Version 1.0  
(c) Microsoft Corporation. All rights reserved.  
Microsoft Research India Pvt. Ltd.  
2010

### GOAL

To support the task of Part-of-Speech Tagging (POS) and other forms of data driven linguistic research on Indian Languages in general, MSR India has developed POS labeled data for Hindi, Bangla and Sanskrit as a part of the Indian Language – Part-of-Speech Tagset (IL-POST) project.

### CORPUS DESCRIPTION

This corpus is designed for those who need annotated text corpora for their work. The corpus is designed based on the IL-POST framework. IL-POST is a POS-tagset framework for Indian Languages, which has been designed to cover the morph-syntactic details of Indian Languages. It supports a three-level hierarchy of Categories, Types and Attributes. The corpus mainly consist therefore of two different level of information for each lexical token (a) lexical Category and Types, and (b) set morphological attributes and their associated values in the context.

#### Example:

विभिन्न\JQ.0.pl.dir.n.nnm राज्यअ\NC.fem.pl.obl.0 में\PP.0.sg.loc अनेक\JQ.0.pl.dir.n.nnm  
महानगरीय\JJ.0.sg.dir परियोजनाएं\NC.0.pl.dir.0 चल\VM.0.0.0.0.0.0.ifn.n रही\VAUX.fem.sg.3.0.prg.dcl.nfn.n  
हैं\VAUX.0.sg.3.prs.sim.dcl.fin.n | \PU

The tag follows the word separated by a ‘\’ (back slash) immediately after the word. There are no blank spaces in between. After the whole POS tag there should be at least one blank (white space) before the next word or a sentinel. In the above example, the first string of 2 or 3 uppercase characters denotes the Category and Type. For example, in the above sentence the word राज्यअ is marked as NC which stands for Noun Common (N denotes Category Noun and C denotes type Common).

The attributes are denoted as numbers or letters, as the case may be, after the tag for the lexical category separated by ‘.’ (dot). The order of the attributes is fixed and cannot be arbitrarily swapped. To illustrate this, consider the category *proper noun* (NC) whose attribute set is {Gender, Number, Case, and Case-marker}. Gender can take values from {Masculine(mas), Feminine(fem), Not-applicable(0)}; Number can take values from the set {Singular (sg), Plural (pl), Not-applicable (0)}; Case can take values from set {Direct(dir), Oblique(obl), Not-applicable(0)}; and Case-marker can take values from the set {Accusative (acc), Genitive (gen), Locative (loc), Not-applicable (0)}. Therefore, for the Common Noun राज्यअ, in the above example sentence, which is feminine, plural, oblique and have no case-marker, the complete tag should be:

\NC.fem.pl.obl.0

- Corpus size
  - o Hindi – Manually annotated 4,859 sentences ( 98,450 words)

- Format of the Data
  - The annotated data is available in both XML and TEXT format for both the language
  - Each data file contains approximately 1- 5,000 words and kept in sentence level.
  - The XML file contains the metadata about the language, encoding, data size etc.
- This data was created under the supervision of
  - **Multilingual Systems Group, Microsoft Research Labs India**
  - **Dr. Girish Nath Jha, Jawaharlal Nehru University, New Delhi.**

### **SOURCE DATA**

The Hindi annotated data targets to cover written modern standard Hindi. Hindi raw corpus is collected randomly from the *Web Duniya* corpus.

### **ANNOTATION PROCEDURE**

The detail of the annotation procedure is downloadable along with the annotated data. Please go through the *annotation guideline* (specific to the language) for clarification and further annotation.

### **DIRECTORY STRUCTURE**

In *MSRI-Data\Annotated\_data\Hindi*:

\*.xml files are in the *Annotated\_data\XML\_files*

\*text files are in *Annotated\_data\text\_files*

In the *docs\* directory:

More detailed information about the part-of-speech tagset and annotation process.

### **CONTACT**

[ilpost@microsoft.com](mailto:ilpost@microsoft.com)