

Linear-chain CRF for NLP in MADlib

KUN LI[§], CHRISTAN GRANT[§], DAISY ZHE WANG[§]

[§]Database Research Center, CISE, University of Florida

{kli, cgrant, daisyw} @cise.ufl.edu

ABSTRACT

MADlib is an open-source library for scalable in-database analytics developed by EMC, UC Berkeley, University of Wisconsin and our group in University of Florida. It provides an evolving suite of data-parallel implementations of mathematical, statistical and machine-learning methods for structured and unstructured data. We are contributing a linear-chain conditional random field for natural language processing such as part of speech tagging (POS), named entity resolution (NER). Motivation: in database analysis, support interactive query, time critical, leverage parallel database, e.g. Greenplum. Instead of extracting features for each token on fly, we extract features for each distinct token and materialized it in the database.

PARALLEL TEXT FEATURE EXTRACTION

Text feature extraction is a step in most statistical text analysis methods, and it can be an expensive operation. To achieve high quality, CRF methods often assign hundreds of features to each token in the document. Examples of such features include: (1) dictionary features: does this token exist in a provided dictionary? (2) regex features: does this token match a provided regular expression? (3) edge features: is the label of a token correlated with the label of a previous token? (4) word features: does this token appear in the training data? and (5) position features: is this token the first or last in the token sequence? The right combination of features depends on the application, and so our support for feature extraction heavily leverages the microprogramming interface provided by MADlib.

PARALLEL LINEAR-CHAIN CRF TRAINING

CRF training is time consuming due to the expensive forward-backward computation to evaluate the log-likelihood function and its gradient vector for each iteration.

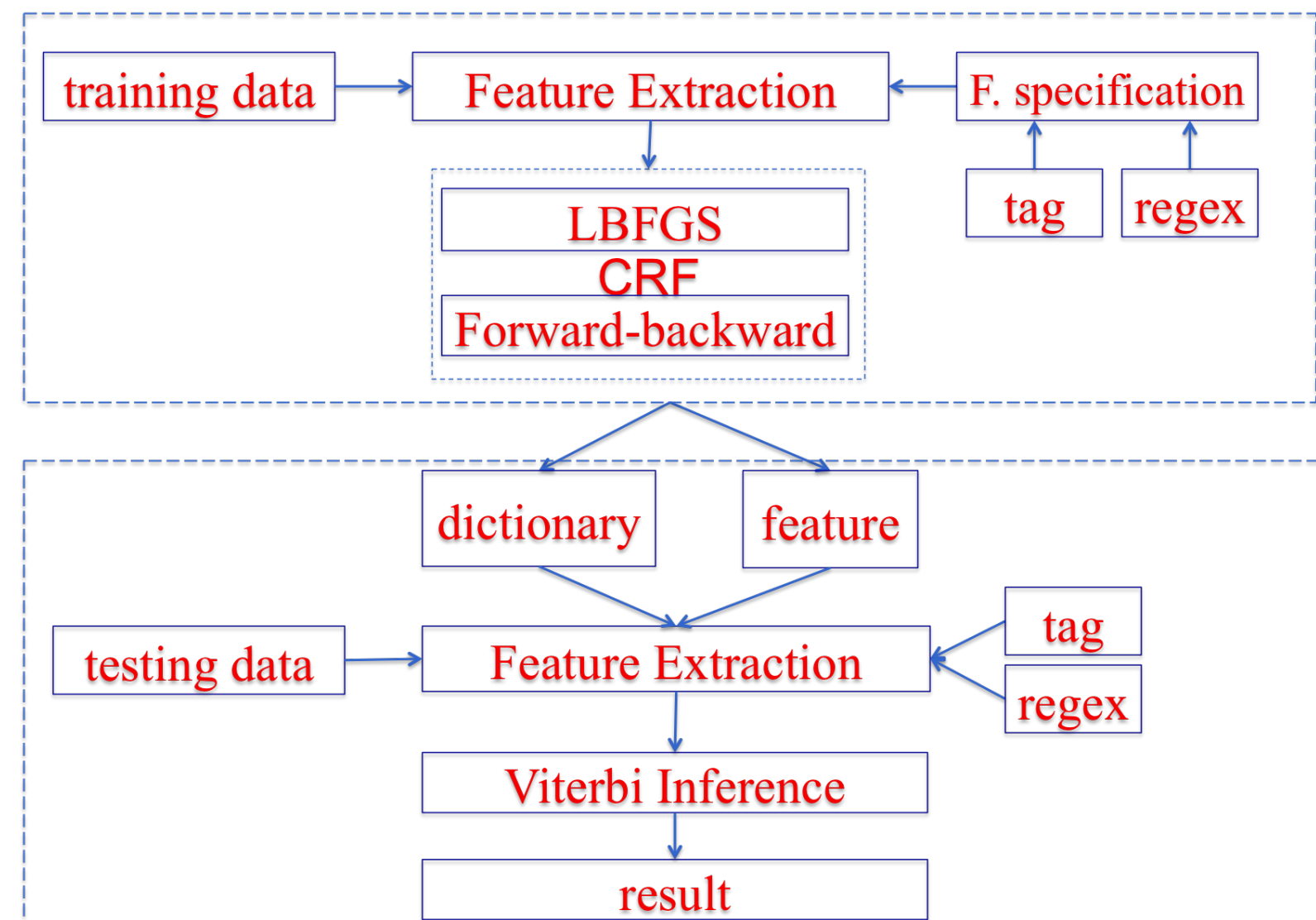
$$\ell'_\lambda = \sum_k [\lambda F(y_k, x_k) - \log Z_\lambda(x_k)] - \frac{\lambda^2}{2\sigma^2} + \text{const}$$

with gradient

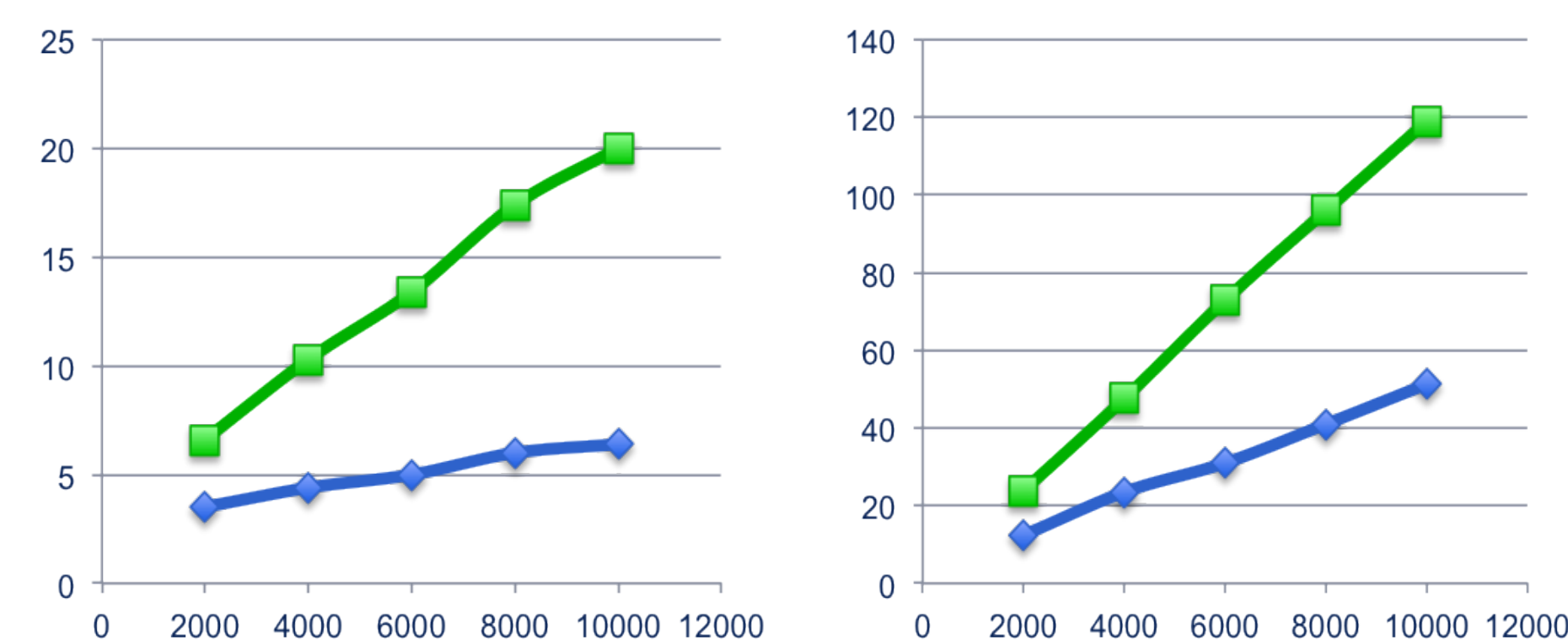
$$\nabla \ell'_\lambda = \sum_k [\lambda F(y_k, x_k) - E_{p_{\lambda(Y|x_k)}} F(Y, x_k)] - \frac{\lambda}{\sigma^2}$$

parallel forward-backward computation

SYSTEM OVERVIEW



PERFORMANCE EVALUATION



ACKNOWLEDGEMENTS

This work is supported by EMC/Greenplum. Thanks to EMC/Greenplum employee Caleb Welton and Hitoshi Harada's precious feedback and code review.

PARALLEL VITERBI INFERENCE

The Viterbi dynamic programming algorithm [14] is the popular algorithm to find the top-k most likely labelings of a document for (linear chain) CRF models. Like any dynamic programming algorithm, the Viterbi algorithm is recursive. We experimented with two direct implementations of macro-coordination over time. First, we chose to implement it using a combination of recursive SQL and window aggregate functions. We discussed this implementation at some length in earlier work [42]. Our initial recursive SQL implementation only runs over PostgreSQL versions 8.4 and later; it does not run in Greenplum. Second, we chose to implement a Python UDF that uses iterations to drive the recursion in Viterbi. This iterative implementation runs over both PostgreSQL and Greenplum. In Greenplum, Viterbi can be run in parallel over different subsets of the document on a multi-core machine. We want to find the sequence of tags that maximizes the formula.

$$P(T_1, T_2, \dots, T_n | W_1, W_2, \dots, W_n) \text{ which can be estimated as } \prod_{1 \leq i \leq n} (P(T_i | T_{i-1}) * P(W_i | T_i))$$

$P(T_i | T_{i-1})$ is the transition weight

$P(W_i | T_i)$ is the emitting weight

-Create tables and import model data to the database `SELECT madlib.load_crf_model('/path/to/modeldata')`

-Feature extraction `SELECT madlib.text_feature_extraction`

-Viterbi inference

`SELECT madlib.vcrf_label(...)` Sample queries

`Q1 : Tagging for sentence with id = 501`

`SELECT doc_id, madlib.vcrf_top1_label(mfactors.score, rfactors.score)`

`FROM m_factors m_factors, rfactors rfactor;`

`Q2 : Tagging for all sentences containing the word 'professor'`

`SELECT DISTINCT ON doc_id, madlib.vcrf_top1_label(mfactors.score, rfactors.score)`

`FROM m_factor, rfactor, seg_tbl, segmenttbl`

`WHERE segtbl.doc_id = segmenttbl.doc_id AND segmenttbl.seg_text = 'professor'`

REFERENCE

[1] <http://crf.sourceforge.net/>

<http://www.cs.berkeley.edu/~daisyw/ViterbiCRF.html>

Xuan-Hieu Phan, Le-Minh Nguyen, and Cam-Tu Nguyen. FlexCRFs: Flexible Conditional Random Fields, 2004.