

Efficient Classification with Conjunctive Features

Naoki Yoshinaga and Masaru Kitsuregawa / Institute of Industrial Science, University of Tokyo

Efficiency-accuracy dilemma in feature-based classifiers

Conjunctive features (e.g., n-grams) play an important role in obtaining accurate classifiers in many NLP tasks.

Ex) Log-linear model (LLM) w/ binary features

$$p(y|x) = \frac{1}{Z(x)} \exp \left(\sum_i w_i f_i(x, y) \right), \quad w_i \in \mathbb{R}$$

Weight

$$f_F(x, y) = \bigwedge_{f_i \in F} f_i(x, y): \text{conjunctive feature}$$

But if we use many conjunctive features, the classification becomes very slow.

$$O(|x|) \Rightarrow O(|x^d|) \quad d \text{ (or less) conjunctions of primitive features in } x$$

Current approaches

There are two ways to speeding up the classification with conjunctive features.

- ✓ **Polynomial kernel** implicitly expresses conjunctive features, but its computation is heavy in practice $O(|SV||x|)$; $|SV| \gg |x|$ [1].
- ✓ **L1-regularized LLM** shrinks the feature space [2], but frequent features are likely to survive.

References

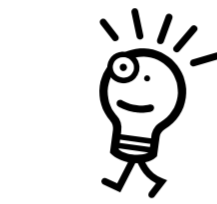
[1] Kudo and Matsumoto: Fast method for kernel-based text analysis. ACL '03.

[2] Gao et al. A comparative study of parameter estimation methods for statistical natural language processing. ACL '07

[3] Chuang et al. Power-law relationship and self-similarity in the itemset support distribution: analysis and applications. The VLDB Journal.

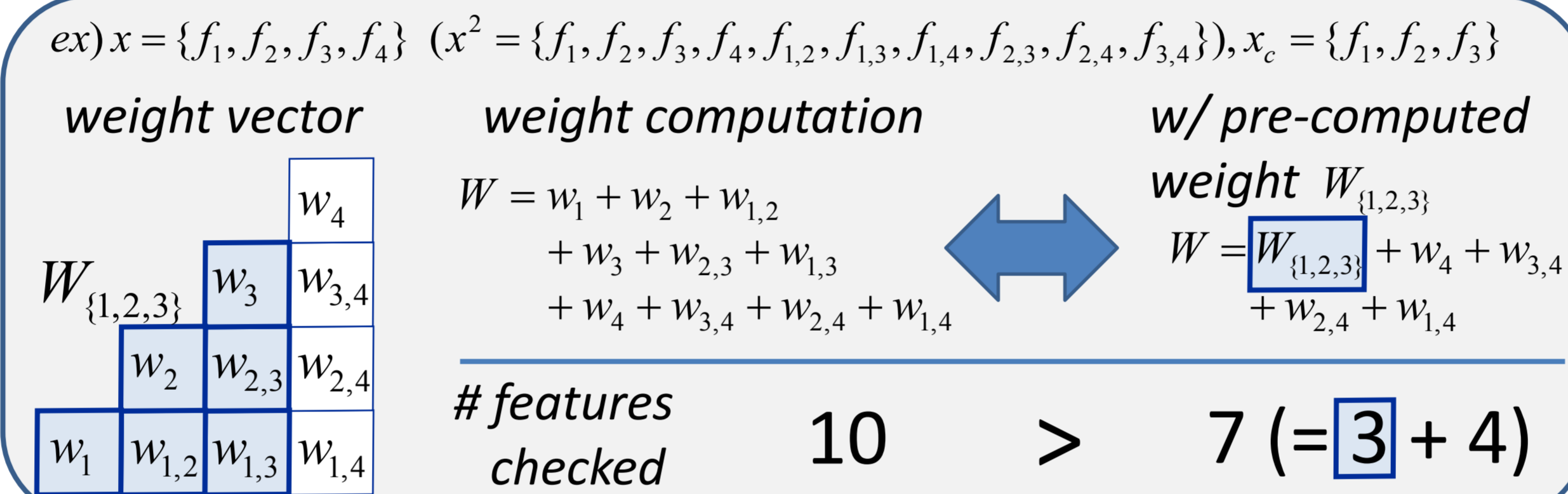
Efficient classifier with feature sequence trie

We speed up a classifier trained with many conjunctive features.



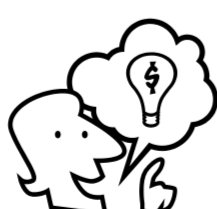
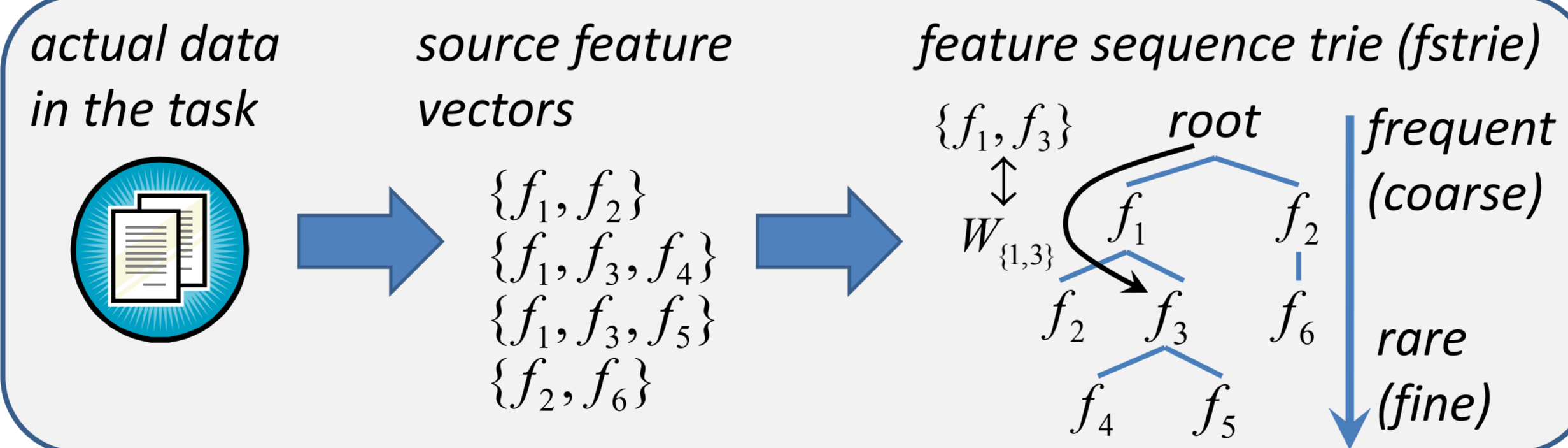
Idea:

Pre-compute weight W_c for some feature vector x_c and use it to obtain weight W of $x \supset x_c$



Issue:

- ✓ Which x_c should we pre-compute weight for?
- ✓ How to find an optimal x_c for input x ?



solution:

Obtain source feature vectors from actual data, sort features in them by frequency, and store weights of the prefix feature vectors into a trie.

Time complexity: $O(|x_c| + |x^d| - |x_c^d|)$ remaining weights (sparse)

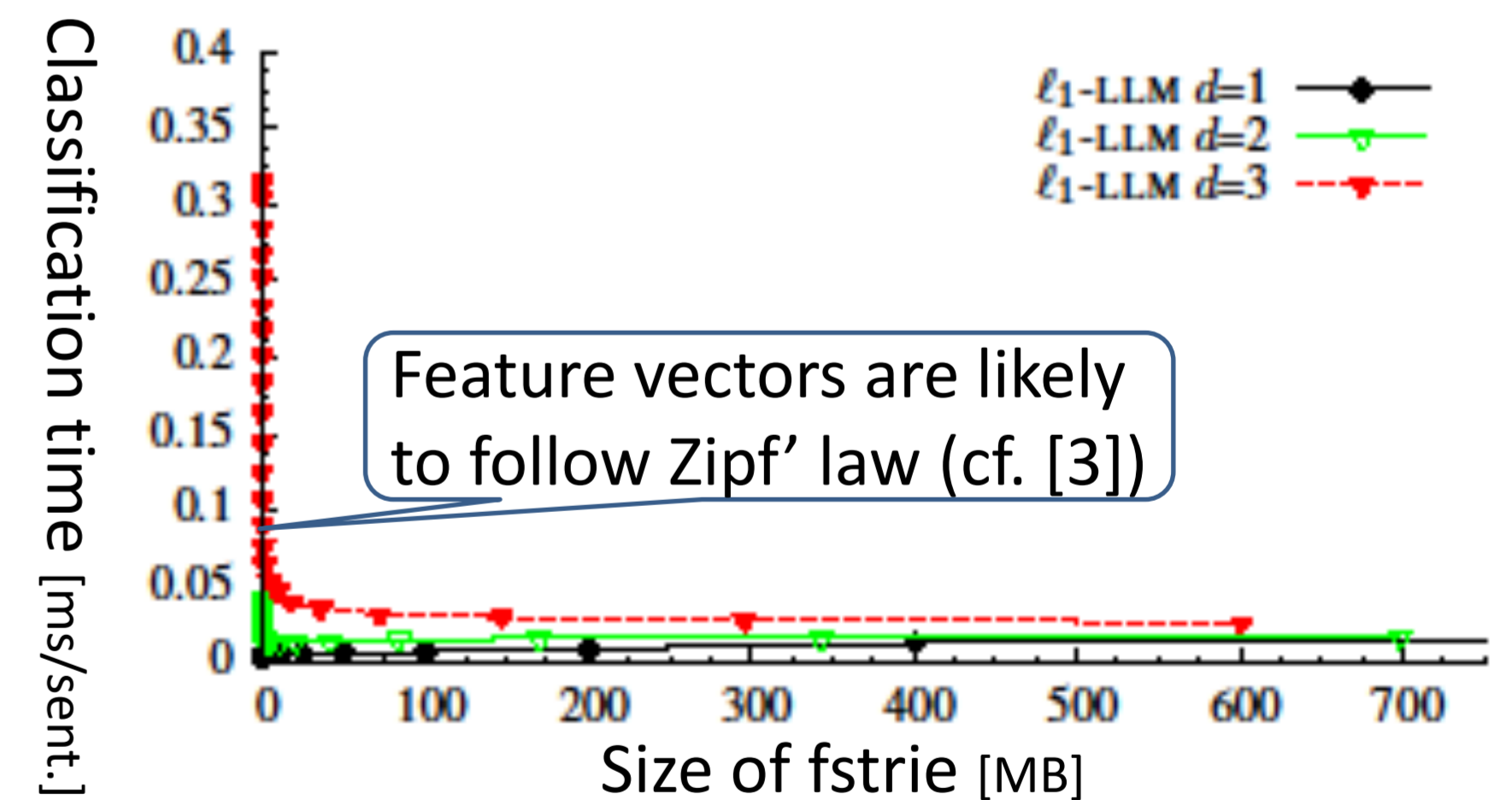
Evaluation: dependency parsing

Results of Japanese dependency parsing

Model	SVM		L1-LLM	
Conj. degree	1	3	1	3
$ F^d $ # feat. ($\times 10^3$)	39.7	26194.4	9.3	129.5
$ x^d $ ave. # feat.	27.3	3286.7	26.5	2088.3
Dep. Acc. (%)	88.29	90.93	88.22	90.71
Classify [ms./sent.]	kernel	13.480	10.945	NA
	baseline	0.003	0.345	0.004
	w/ fstrie	NA	0.079	NA
Parse [ms./sent.]	0.015	0.093	0.016	0.040

- ✓ Kyoto Corpus; Sassano's Shift/reduce parser
- ✓ 3,258,313 sentences (news article) were parsed to build fstrie (weight calculation took 1 hour)

Classification time as a function of fstrie size



- * The nodes in fstrie are pruned according to their probability and impact on computation reduction

What's next?

- ✓ Evaluation in other tasks (try the code at <http://www.tkl.iis.u-tokyo.ac.jp/~ynaga/>)
- ✓ Application to structured prediction (CRF)