

Scalable Online Training with Conjunctive Features

Naoki Yoshinaga and Masaru Kitsuregawa

IIS, The University of Tokyo

Proposal

- Kernel slicing for online training with conjunctive features
 - explicitly consider conjunctions among frequent features, while implicitly considering the others by polynomial kernel
 - reuse temporal margins of partial feature vectors
- Performance evaluation on two NLP tasks
(dependency parsing and hyponymy relation extraction)
 - orders of magnitudes faster than kernel-based online training, while retaining its space efficiency
 - model accuracy: comparable to batch SVM

Overview

- Research Backgrounds
 - Space-time trade-off in training with conjunctive features
 - Kernel splitting [Goldberg+ '08] for testing
- Methods
 - Online learning with kernel splitting
 - Online learning with kernel slicing
- Experiments
- Conclusion

Conjunctive features in NLP

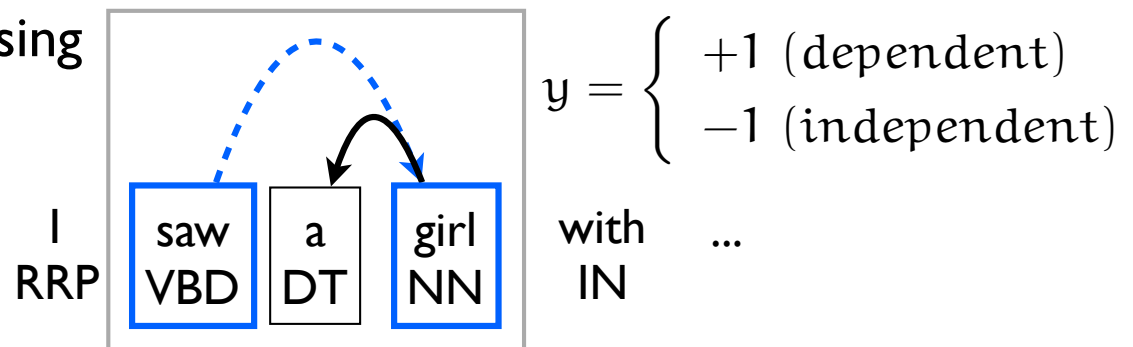
- **Conjunctive features** play a key role to obtain a high degree of accuracy in NLP classification problems
- dependency parsing [Koo+ '08], pronoun resolution [Nguyen+, 08], semantic role labeling [Liu+, '07], relation extraction [Sumida+ '08]

ex. | dependency parsing

Linear model
[LLM, Perceptron, etc.]

$$y = \text{sgn}(\underline{\mathbf{w}}^T \phi_d(\mathbf{x}))$$

high-dimensional
weight vector



$\mathbf{x} = \langle \underline{f_1, f_2, f_3, f_4} \rangle$ active (primitive) features

$\phi_2(\mathbf{x}) = \langle f_1, f_2, f_3, f_4, \underline{f_1 \wedge 2, f_1 \wedge 3, \dots, f_3 \wedge 4} \rangle$

conjunctive features

$f_{i \wedge j} \neq 0 \iff f_i \neq 0 \cup f_j \neq 0$

Conjunctive features in NLP

- **Conjunctive features** play a key role to obtain a high degree of accuracy in NLP classification problems
- dependency parsing [Koo+ '08], pronoun resolution [Nguyen+, 08], semantic role labeling [Liu+, '07], relation extraction [Sumida+ '08]

Linear model
[LLM, Perceptron, etc.]

$$y = \text{sgn}(\underline{\mathbf{w}}^T \phi_d(\mathbf{x}))$$

*high-dimensional
weight vector*

Kernel-based model
[SVM, Kernel perceptron, etc.]

$$y = \text{sgn}\left(\sum_{\mathbf{s}_i \in \underline{\mathcal{S}}} \alpha_i k_d(\mathbf{s}_i, \mathbf{x})\right)$$

*support set
(subset of training examples)*

polynomial kernel

$$\begin{aligned} k_d(\mathbf{s}, \mathbf{x}) \\ &= \phi_d(\mathbf{s})^T \phi_d(\mathbf{x}) \\ &= (\mathbf{s}^T \mathbf{x} + 1)^d \end{aligned}$$

Space-time tradeoff in training with conjunctive features

- Training with conjunctive features involves **space-time tradeoff** in the way conjunctive features are handled

linear training
(perceptron)

```

 $w \leftarrow \mathbf{0}$ 
for  $t = 1$  to  $T$  do
   $m_t \leftarrow w^T \phi_d(x_t)$ 
  if  $y_t m_t < 0$ 
     $w$   $\leftarrow w + y_t \phi_d(x_t)$ 
  endif needs huge memory
end
  
```

} ← initialize → {
 } ← compute margin → {
 } ← update → {

kernel-based training
(kernel perceptron)

```

 $\mathcal{S}_0 \leftarrow \emptyset, \alpha \leftarrow \mathbf{0}$  linearly increase as
for  $i = 1$  to  $T$  do training proceeds
   $m_t \leftarrow \sum_{s_i \in \mathcal{S}_{t-1}} \alpha_i k_d(s_i, x_t)$ 
  if  $y_t m_t < 0$ 
     $\alpha_t$   $\leftarrow y_t, \mathcal{S}_t = \mathcal{S}_{t-1} \cup \{x_t\}$ 
  endif
end
  
```

Linear training: **polynomial space** in the number of primitive features
 Kernel-based training: **quadratic time** in the number of examples

Kernel splitting [Goldberg+ 2008]

(for testing)

- Split features into common ones \mathcal{F}_C and rare ones $\mathcal{F} \setminus \mathcal{F}_C$ and divide margin computation: *according to frequency in \mathcal{S}*
 - explicitly consider conjunctions among common features
 - implicitly consider remaining conjunctions by kernel

$$\begin{aligned}\sum_{\mathbf{s}_i \in \mathcal{S}} \alpha_i k_d(\mathbf{s}_i, \mathbf{x}) &= \sum_{\mathbf{s}_i \in \mathcal{S}} \alpha_i k_d(\mathbf{s}_i, \underline{\mathbf{x}_C}) + \sum_{\mathbf{s}_i \in \mathcal{S}} \alpha_i \{k_d(\mathbf{s}_i, \mathbf{x}) - k_d(\mathbf{s}_i, \mathbf{x}_C)\} \\ &= \mathbf{w}_C^T \phi_d(\mathbf{x}_C) + \sum_{\mathbf{s}_i \in \mathcal{S}_R} \alpha_i \{k_d(\mathbf{s}_i, \mathbf{x}) - k_d(\mathbf{s}_i, \mathbf{x}_C)\}\end{aligned}$$

Kernel splitting [Goldberg+ 2008]

(for testing)

- Split features into common ones \mathcal{F}_C and rare ones $\mathcal{F} \setminus \mathcal{F}_C$ and divide margin computation: *according to frequency in \mathcal{S}*
- explicitly consider conjunctions among common features
- implicitly consider remaining conjunctions by kernel

$$\sum_{s_i \in \mathcal{S}} \alpha_i k_d(s_i, x) = \sum_{s_i \in \mathcal{S}} \alpha_i k_d(s_i, \underline{x_C}) + \sum_{s_i \in \mathcal{S}} \alpha_i \{k_d(s_i, x) - k_d(s_i, x_C)\}$$

$$= w_C^T \phi_d(x_C) + \sum_{s_i \in \mathcal{S}_R} \alpha_i \{k_d(s_i, x) - k_d(s_i, x_C)\}$$

explicit weights for common conjunctive features $|w_C| \ll |w|$

$$w_C = \sum_{s_j \in \mathcal{S}} \alpha_j \phi_d(s_j \cap \mathcal{F}_C)$$

space-efficient linear classification

Kernel splitting [Goldberg+ 2008]

(for testing)

- Split features into common ones \mathcal{F}_C and rare ones $\mathcal{F} \setminus \mathcal{F}_C$ and divide margin computation: *according to frequency in \mathcal{S}*
- explicitly consider conjunctions among common features
- implicitly consider remaining conjunctions by kernel

$$\begin{aligned}
 \sum_{s_i \in \mathcal{S}} \alpha_i k_d(s_i, x) &= \sum_{s_i \in \mathcal{S}} \alpha_i k_d(s_i, \underline{x_C}) + \sum_{s_i \in \mathcal{S}} \alpha_i \{ \underbrace{k_d(s_i, x) - k_d(s_i, x_C)}_{= 0 \text{ when } s_j^T x = s_j^T x_C} \} \\
 &= w_C^T \phi_d(x_C) + \sum_{s_i \in \mathcal{S}_R} \alpha_i \{ k_d(s_i, x) - k_d(s_i, x_C) \}
 \end{aligned}$$

consider a few support vectors $|\mathcal{S}_R| \ll |\mathcal{S}|$
that have rare feature $f_R \in x_R = x \setminus x_C$

space-efficient linear classification + efficient kernel-based testing

Online learning with kernel splitting

- Replace margin computation part in kernel-based online learning with **kernel splitting**

Kernel perceptron

```

$$\mathcal{S}_0 \leftarrow \emptyset, \alpha \leftarrow \mathbf{o},$$
  
for  $t = 1$  to  $T$  do  
  
$$m_t \leftarrow \sum_{s_i \in \mathcal{S}_{t-1}} \alpha_i k_d(s_i, x_t)$$
  
  if  $y_t m_t \leq 0$   
     $\alpha_t \leftarrow y_t, \mathcal{S}_t = \mathcal{S}_{t-1} \cup \{x_t\}$   
  endif  
end
```

replace this with
kernel splitting

Online learning with kernel splitting

- Replace margin computation part in kernel-based online learning with **kernel splitting**

Kernel perceptron

kernel splitting

```

$$\mathcal{S}_0 \leftarrow \emptyset, \alpha \leftarrow \mathbf{o},$$
  
  
for  $t = 1$  to  $T$  do  
   $\mathbf{x}_C \leftarrow \mathbf{x}_t \cap \mathcal{F}_C$     Q. 1: how to determine  $\mathcal{F}_C$  ?  
   $m_t \leftarrow \mathbf{w}_C \phi_d(\mathbf{x}_C)$     Q. 2: how to maintain  $\mathbf{w}_C$  ?  
     $+ \sum_{\mathbf{s}_i \in \mathcal{S}_R} \alpha_i \{k_d(\mathbf{s}_i, \mathbf{x}_t) - k_d(\mathbf{s}_i, \mathbf{x}_C)\}$   
  if  $y_t m_t \leq 0$   
     $\alpha_t \leftarrow y_t, \mathcal{S}_t = \mathcal{S}_{t-1} \cup \{\mathbf{x}_t\}$   
  endif  
end
```

Online learning with kernel splitting

- Replace margin computation part in kernel-based online learning with **kernel splitting**

Kernel perceptron with **kernel splitting**

A.1: Choose top-N frequent features in the training examples as \mathcal{F}_C

kernel splitting

A.2 Online-update w_C to correspond with $\langle \mathcal{S}_t, \alpha_t \rangle$

```

 $\mathcal{S}_0 \leftarrow \emptyset, \alpha \leftarrow \mathbf{o},$ 
 $\mathcal{F}_C \leftarrow \{f \mid \text{RANK}(f) \leq N\}, w_C \leftarrow \mathbf{o}$ 
for  $t = 1$  to  $T$  do
     $x_C \leftarrow x_t \cap \mathcal{F}_C$     Q.1: how to determine  $\mathcal{F}_C$ ?
     $m_t \leftarrow w_C \phi_d(x_C)$   Q.2: how to maintain  $w_C$ ?
         $+ \sum_{s_i \in \mathcal{S}_R} \alpha_i \{k_d(s_i, x_t) - k_d(s_i, x_C)\}$ 
    if  $y_t m_t \leq 0$ 
         $\alpha_t \leftarrow y_t, \mathcal{S}_t = \mathcal{S}_{t-1} \cup \{x_t\}$ 
         $w_C \leftarrow w_C + y_t \phi_d(x_C)$ 
    endif
end
```

Online learning with kernel splitting

- Replace margin computation part in kernel-based online learning with **kernel splitting**

Kernel perceptron with **kernel splitting**

```

$$\begin{aligned} \mathcal{S}_0 &\leftarrow \emptyset, \alpha \leftarrow \mathbf{o}, \\ \mathcal{F}_C &\leftarrow \{f \mid \text{RANK}(f) \leq N\}, \mathbf{w}_C \leftarrow \mathbf{o} \\ \text{for } t = 1 \text{ to } T \text{ do} \\ &\quad \mathbf{x}_C \leftarrow \mathbf{x}_t \cap \mathcal{F}_C \\ &\quad \mathbf{m}_t \leftarrow \mathbf{w}_C \phi_d(\mathbf{x}_C) \\ &\quad \quad + \sum_{\mathbf{s}_i \in \mathcal{S}_R} \alpha_i \{k_d(\mathbf{s}_i, \mathbf{x}_t) - k_d(\mathbf{s}_i, \mathbf{x}_C)\} \\ &\quad \text{if } y_t \mathbf{m}_t \leq 0 \\ &\quad \quad \alpha_t \leftarrow y_t, \mathcal{S}_t = \mathcal{S}_{t-1} \cup \{\mathbf{x}_t\} \\ &\quad \quad \mathbf{w}_C \leftarrow \mathbf{w}_C + y_t \phi_d(\mathbf{x}_C) \\ &\quad \text{endif} \\ \text{end} \end{aligned}$$

```

Assumption:
additive updates

$\forall t' > t$
 $\langle \alpha_t, \mathcal{S}_t \rangle \subseteq \langle \alpha_{t'}, \mathcal{S}_{t'} \rangle$

Intricacy in setting Parameter N

- Kernel splitting can **control** space-time trade-off in training with conjunctive features, but it **does not resolve** it

$$\begin{aligned}x_C &\leftarrow x_t \cap \mathcal{F}_C \\ m_t &\leftarrow w_C \phi_d(x_C) + \sum_{s_i \in \mathcal{S}_R} \alpha_i \{k_d(s_i, x_t) - k_d(s_i, x_C)\}\end{aligned}$$

Time complexity : $\mathcal{O}(|x_C|^d + |\mathcal{S}_R||x_C|)$

- N ($= |\mathcal{F}_C|$) should be **smaller** for higher-order conjunctive features (to keep $|x_C|^d$ and $|\mathcal{F}_C|^d$ small)
- N ($= |\mathcal{F}_C|$) should be **larger** when we handle a larger number training examples (to keep $|\mathcal{S}_R|$ small)

Kernel slicing | basic idea

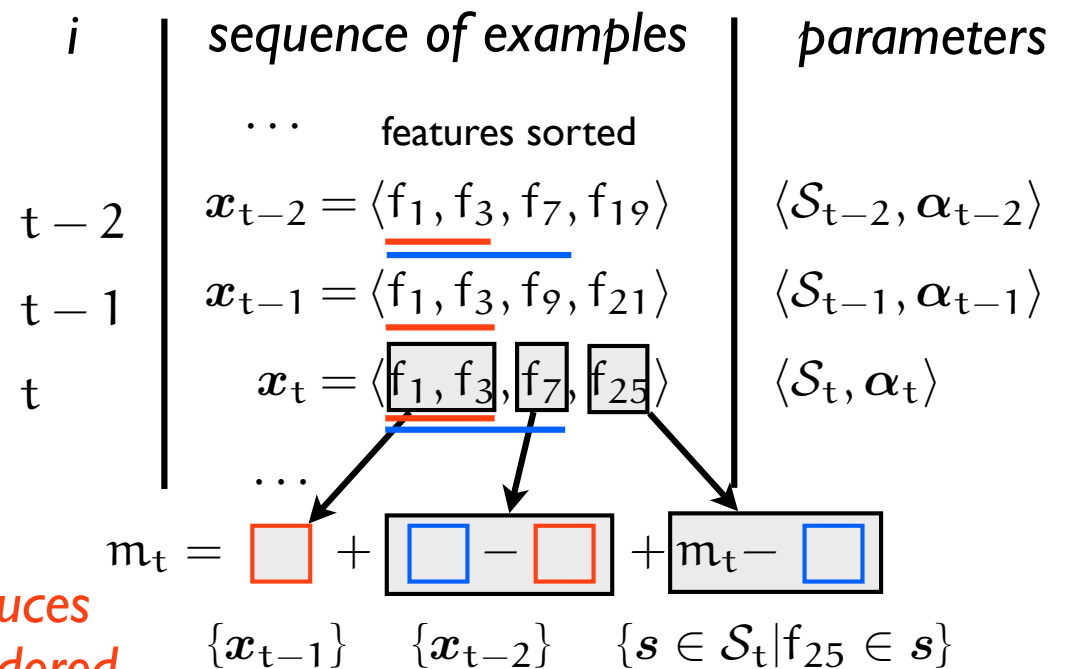
- Examples in real-world data are redundant [Yoshinaga+ '09]
- Online learner will repeatedly compute margins of **common partial feature vectors**

kernel perceptron

```

 $\mathcal{S}_0 \leftarrow \emptyset, \alpha \leftarrow \mathbf{o}$ 
for  $i = 1$  to  $T$  do
   $m_t \leftarrow \sum_{s_i \in \mathcal{S}_{t-1}} \alpha_i k_d(s_i, \mathbf{x}_t)$ 
  if  $y_t m_t < 0$ 
     $\alpha_t \leftarrow y_t, \mathcal{S}_t = \mathcal{S}_{t-1} \cup \{\mathbf{x}_t\}$ 
  endif
end
    
```

*reusing partial margins reduces
support vectors to be considered*



Kernel slicing | feature-wise splitting

- Kernel slicing: incrementally compute a **partial margin** of x_t when adding features to $x_t^0 (= \emptyset)$ from frequent to rare

$$m_t = m_t^0 + \sum_{j=1}^{|x_t|} \underline{m_t^j} \quad \begin{array}{l} \text{margin change when we add } j\text{-th frequent feature} \\ \text{temporal} \\ \text{partial margin} \end{array} \quad m_t^j = \sum_{s_i \in \mathcal{S}_t} \alpha_i (k_d(s, x_t^j) - k_d(s, x_t^{j-1}))$$

- retrieve / update partial margins (with time index t) in a trie

$$m_t^j = \underline{m_{t'}^j} + \sum_{s_i \in \mathcal{S}_j} \alpha_i \{k_d(s_i, x_j) - k_d(s_i, x_{j-1})\}$$

*partial margin computed
for x_t^j at past round $t' (< t)$*

$\mathcal{S}_j = \{s \in \mathcal{S}_t \setminus \mathcal{S}_{t'} \mid f_j \in s\}$
newly added support vectors since we finally see x_t^j

- when common feature $f_j \in \mathcal{F}_C$ is added and the retrieved margin was too old, use w_C to compute the partial margin

$$|\phi_d(x_j) - \phi_d(x_{j-1})| < |\mathcal{S}_j| |x_{j-1}| \quad m_t^j = w_C^T \{\phi_d(x_t^j) - \phi(x_t^{j-1})\}$$

Experiments

- Implement online passive aggressive I (PA-I) [Crammer+ '06] with kernel slicing
- Compare our learner with
 - Support vector machine (SVM) [TinySVM by T. Kudo]
 - kernel-based PA-I with inverted indices [Okanohara+ '07]
 - SGD-training of ℓ_1 -regularized log-linear model [Tsuruoka, '09]
- Evaluate on two NLP tasks:
Japanese dependency parsing and hyponymy relation extraction

Task settings

- Japanese dependency parsing
 - Classifier judges whether a given head/dependent candidate has a dependency relation (in shift-reduce parser [Sassano, '04])
 - Features: POS(-subcategory), inflection form of head / dependent, and surrounding contexts (distance etc.)
- Hyponymy relation extraction
 - Classifier judges whether a given pair of entities extracted from Wikipedia articles forms a hyponymy relation [Sumida+ '08]
 - Features: POS, surface string, morpheme, listing type of each entity, and surrounding contexts (distance etc.)

We considered third-order conjunctive features in training

Example / Feature Statistics

- Feature conjunctions dramatically increase
 - the average number of active features
 - the feature space

DATA SET	dependency parsing	hyponymy extraction
T (# examples)	296,776	201,664
Ave. of $ \mathbf{x} $	27.6	15.4
Ave. of $ \phi_3(\mathbf{x}) $	3558.3	798.7
$ \mathcal{F} $ (# features)	64,493	306,036
$ \mathcal{F}^3 $ (# conj. features)	58,361,669	64,249,234

x130

x50

x900

x210

Labeled examples are available from: <http://www.tkl.iis.u-tokyo.ac.jp/~ynaga/pecco/>
<http://nlpwww.nict.go.jp/hyponymy/>

Results | dependency parsing

- PA-I with kernel slicing was the fastest, while retaining space-efficiency of kernel-based training
- hyper-parameters are tuned to maximize model accuracy on development set

	METHOD	ACC.	TIME	MEMORY
kernel-based training	SVM (batch)	90.93%	25912s	243MB
	PA-I kernel	90.90%	8704s	83MB
	PA-I splitting	90.90%	351s	149MB
	PA-I slicing	90.89%	262s	175MB
linear training	PA-I linear	90.90%	465s	993MB
	ℓ_1 -LLM (SGD)	90.76%	4057s	21499MB

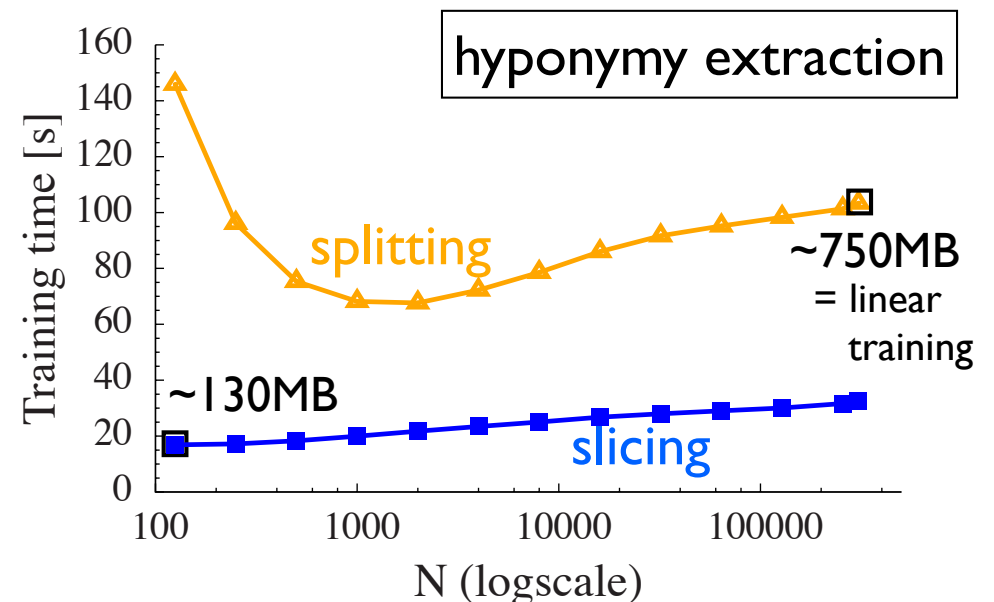
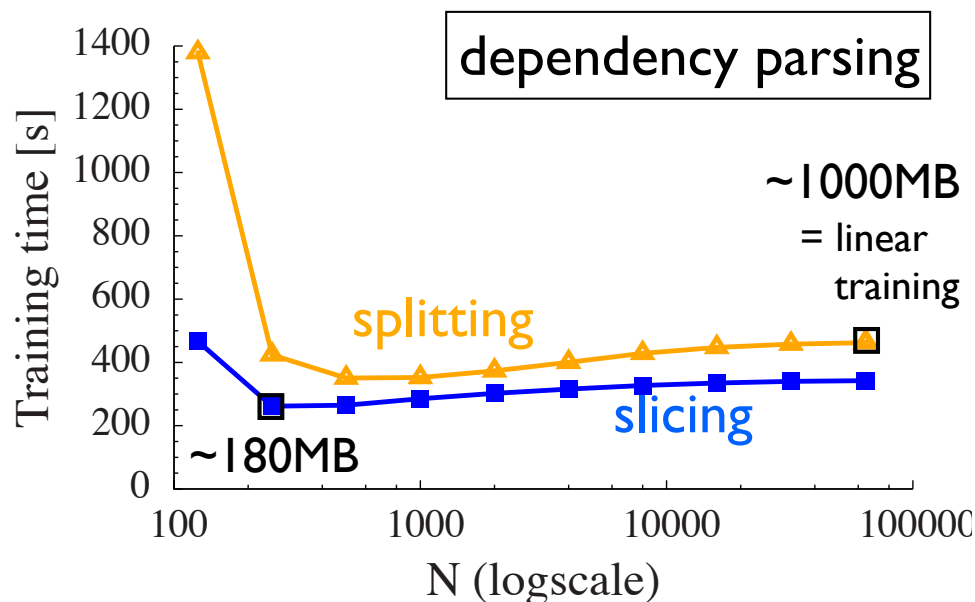
Results | hyponymy extraction

- PA-I with kernel slicing was the fastest, while retaining space-efficiency of kernel-based training
- hyper-parameters are tuned to maximize model accuracy on development set

	METHOD	ACC.	TIME	MEMORY
kernel-based training	SVM (batch)	93.09%	17354s	140MB
	PA-I kernel	93.14%	1074s	49MB
	PA-I splitting	93.10%	68s	108MB
	PA-I slicing	93.05%	17s	131MB
linear training	PA-I linear	93.11%	103s	751MB
	ℓ_1 -LLM (SGD)	92.86%	779s	14089MB

Splitting vs. Slicing | Parameter N

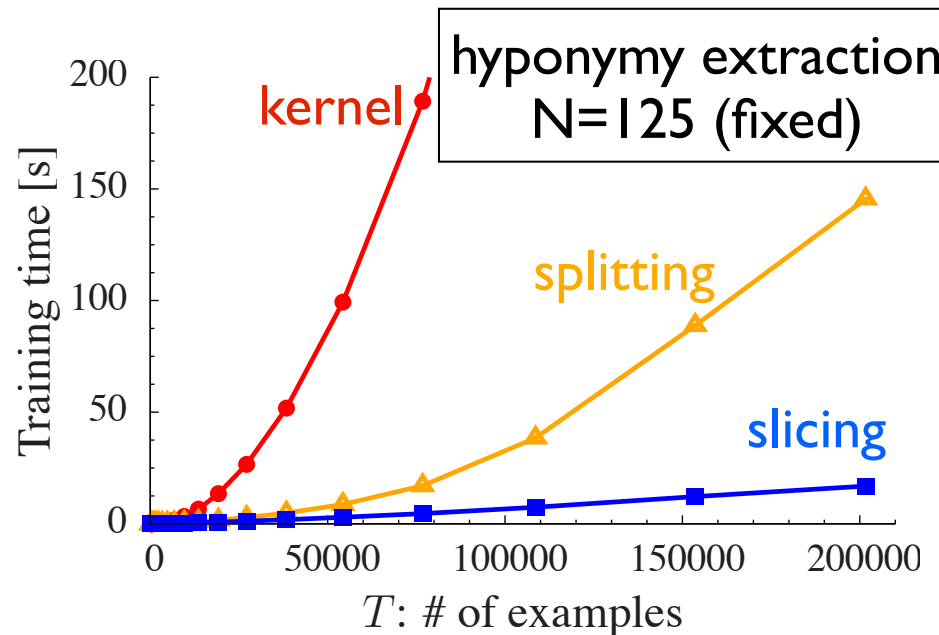
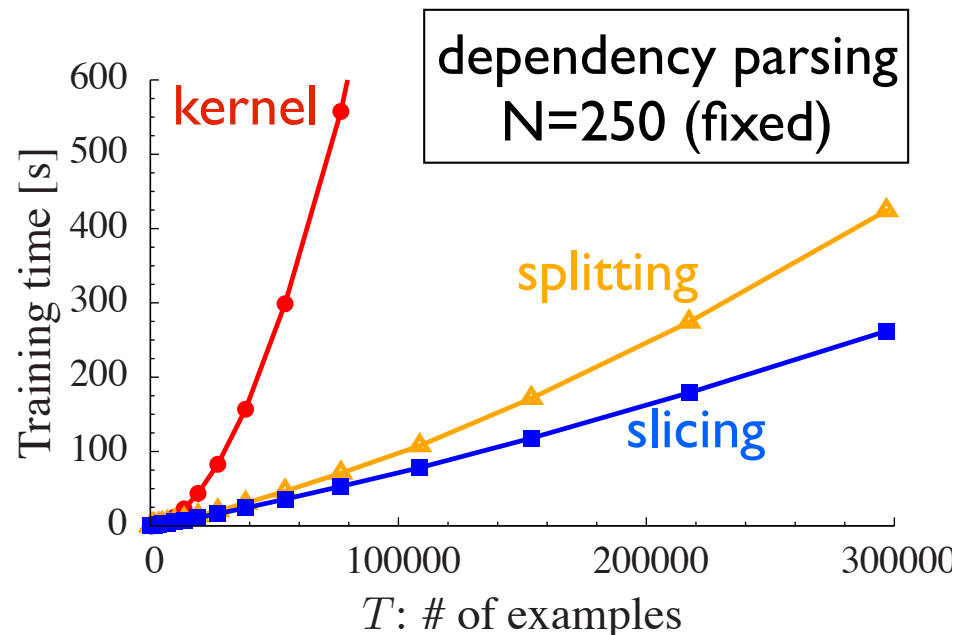
- Training time as a function of parameter N [mem. usage]
 - **kernel splitting**: parameter-sensitive
 - **kernel slicing**: parameter-insensitive



You don't need to *tune* N (it doesn't change the model, anyway)

Splitting vs. Slicing | # examples, T

- Training time as a function of the number of examples
 - **kernel splitting**: in between linear and quadratic
 - **kernel slicing**: almost linear



Reusing temporal margins → more scalable training

Related Work

- Feature selection in linear training [Wu+ '07, Okanohara+, '09]
(limit the number of conjunctive features)
 - Simpler model → faster, more space-efficient, *less accurate*
x17 but 94.19% → 93.71% (named entity recognition [Wu+ '07]),
x37 but 89.52% → 89.03% (dependency parsing [Okanohara+ '09])
- Bounded Kernel-based training [Dekel+ '06; Cavallanti+ '07]
(limit the number of support vectors)
 - These lightweight algorithms could not bound the number of support vectors, while retaining model accuracy [Orabona+ '09]

Our method exploits the data redundancy in evaluating the kernel to train the **same** model as the base learner

Conclusion

- Scalable online training method with kernel slicing
 - Kernel slicing generalizes kernel splitting [Goldberg+ '08], to reuse temporal partial margins for common partial feature vectors
 - orders of magnitude faster than kernel-based online training, while retaining its space efficiency
- Things I didn't mention in this talk (see our paper):
 - Efficient management of feature weights and partial margins (packed training examples) with a double-array trie [Yata+ '09]
 - Termination of margin computations that will never contribute to parameter updates (safely skipping rare features)

Future work

- ~~Release C++ implementation and dataset:~~ *done.*
<http://www.tkl.iis.u-tokyo.ac.jp/~ynaga/opal/>
- Fast testing? - you may want to try pecco [Yoshinaga+, EMNLP '09]
<http://www.tkl.iis.u-tokyo.ac.jp/~ynaga/pecco/>
- Implement kernel slicing for other online algorithms
- Generalize kernel slicing to accommodate other kernels

Thank you

