

Part 2: Structured Learning

Kai-Wei Chang

Microsoft Research → Univ. of Virginia

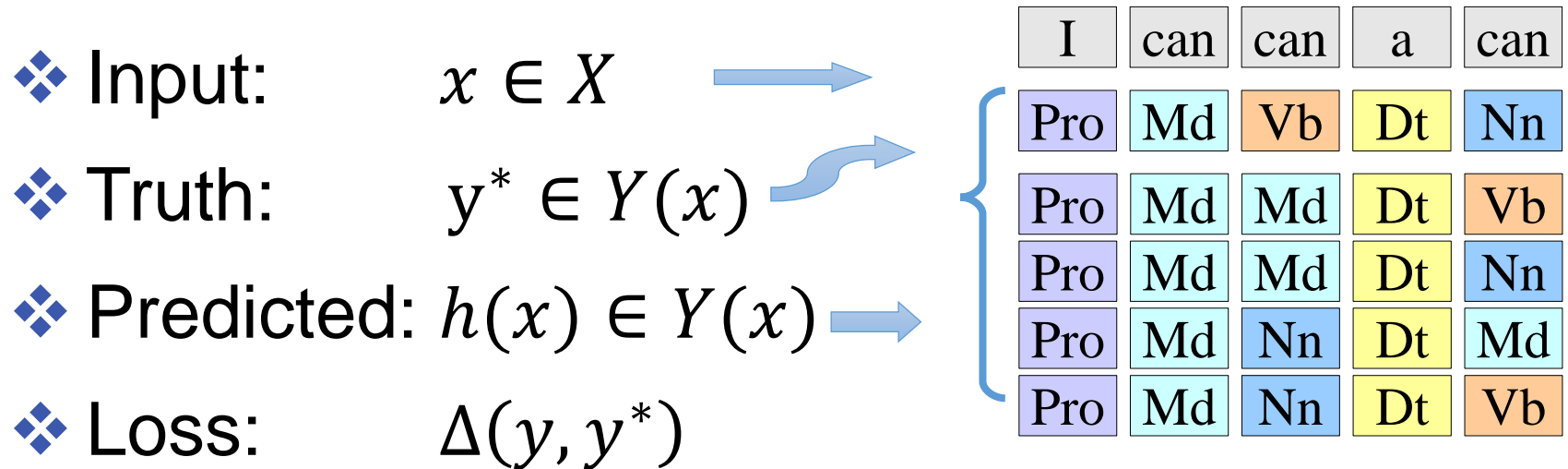
- **Part 2: Learning a Structured Prediction Model**
(45min)
 - **Definition of structured learning**
 - **Local Learning v.s. Global Learning**
 - **Global Learning Algorithms**
 - Online learning: Structured Perceptron
 - Batch learning: Structured SVM
 - **Optimization methods for Structured SVM**
 - Stochastic Gradient Decent
 - Dual Coordinate Descent
 - Learning on a multi-core machine

CCM Formulations

$$y = \operatorname{argmax}_{y \in Y} \mathbf{w}^T \phi(x, y) + \mathbf{u}^T \mathbf{C}(x, y)$$

This part of the tutorial focuses on learning \mathbf{w} (and \mathbf{u})

Learning a Structured Prediction Model



Goal: find $h \in H$ such that $h(x) \in Y(x)$

minimizing $E_{(x,y) \sim D} [\Delta(y, h(x))]$ based on N samples $(x_n, y_n) \sim D$

Learning Paradigms [\[Punyakanok+ 05\]](#)

❖ Local Learning

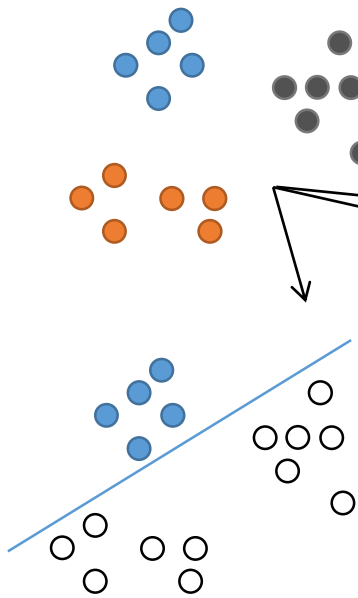
- ❖ Learning local models independently
- ❖ Ensure output is coherent at test time
- ❖ E.g., One-against-all multiclass classification

❖ Global Learning

- ❖ Learning with inference
- ❖ Training and testing are consistent
- ❖ Constrained classification for multiclass
[\[Har-Peled et. al 2002; Crammer et. al 2002\]](#)

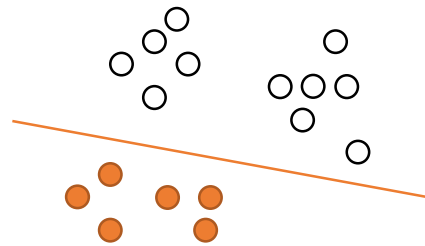
Visualizing One-vs-all

From the full dataset, construct three binary classifiers, one for each class

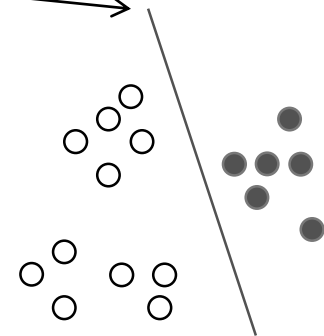


$\mathbf{w}_{\text{blue}}^T \mathbf{x} > 0$ for
blue inputs

Notation: Score
for blue label



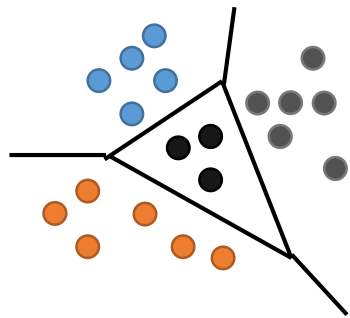
$\mathbf{w}_{\text{red}}^T \mathbf{x} > 0$ for
red inputs



$\mathbf{w}_{\text{green}}^T \mathbf{x} > 0$ for
green inputs

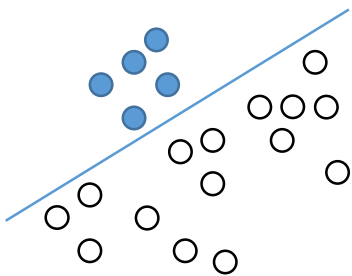
*Winner Take All will predict the right answer.
Only the correct label will have a positive score*

One-vs-all may not always work

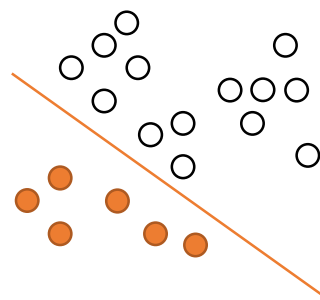


Black points are not separable with a single binary classifier

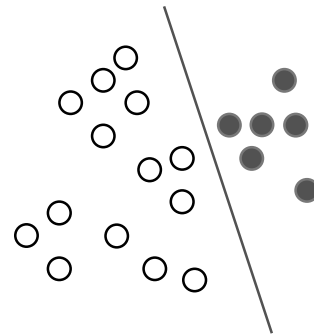
The decomposition will not work for these cases!



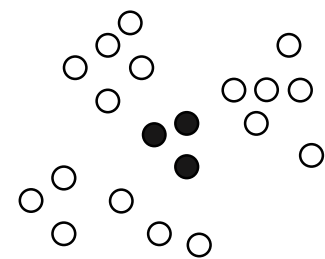
$\mathbf{w}_{\text{blue}}^T \mathbf{x} > 0$
for **blue**
inputs



$\mathbf{w}_{\text{red}}^T \mathbf{x} > 0$
for **red**
inputs



$\mathbf{w}_{\text{green}}^T \mathbf{x} > 0$
for **green**
inputs



???

Local Learning: One-vs-all classification

- ❖ Easy to learn

- ❖ Use any binary classifier learning algorithm

- ❖ **Problems**

- ❖ Calibration issues

- ❖ We are comparing scores produced by K classifiers trained independently. No reason for the scores to be in the same numerical range!

- ❖ Might not always work

- ❖ Yet, works fairly well in many cases, especially if the underlying binary classifiers are tuned, regularized

Global Learning Motivation

- ❖ Decomposition methods
 - ❖ Do not account for how the final predictor will be used
 - ❖ Do not optimize any global measure of correctness
- ❖ **Goal:** To train a multiclass classifier that is “global”

Global “One-vs-all” Approach

- ❖ *Idea*: Create K classifiers $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$. For examples with label i , we want

$$\mathbf{w}_i^T \mathbf{x} > \mathbf{w}_j^T \mathbf{x} \text{ for all } j$$

- ❖ **Prediction**: $\operatorname{argmax}_i \mathbf{w}_i^T \mathbf{x}$

- ❖ **Training**: For each training example (\mathbf{x}_i, y_i) :

$$\hat{y} \leftarrow \operatorname{arg} \max_j \mathbf{w}_j^T \phi(\mathbf{x}_i, y_i)$$

if $\hat{y} \neq y_i$

η : learning rate

$$\mathbf{w}_{y_i} \leftarrow \mathbf{w}_{y_i} + \eta \mathbf{x}_i \quad (\text{promote})$$

$$\mathbf{w}_{\hat{y}} \leftarrow \mathbf{w}_{\hat{y}} - \eta \mathbf{x}_i \quad (\text{demote})$$

Joint Inference with General Constraint Structure

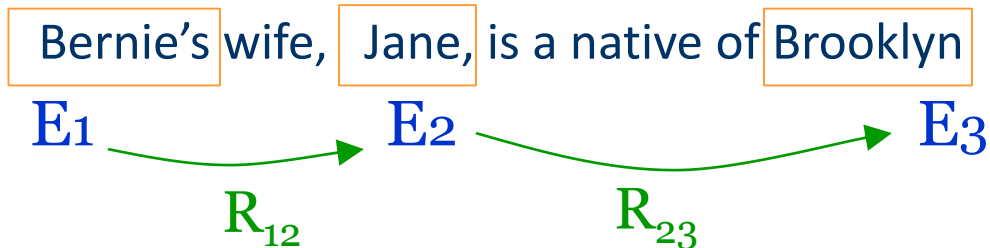
[Roth&Yih'04,07,....]

Recognizing Entities and Relations

other	0.05
per	0.85
loc	0.10

other	0.10
per	0.60
loc	0.30

other	0.05
per	0.50
loc	0.45



irrelevant	0.05
spouse_of	0.45
born_in	0.50

irrelevant	0.10
spouse_of	0.05
born_in	0.85

Models could be learned separately/jointly; constraints may come up only at decision time.

Training Methods

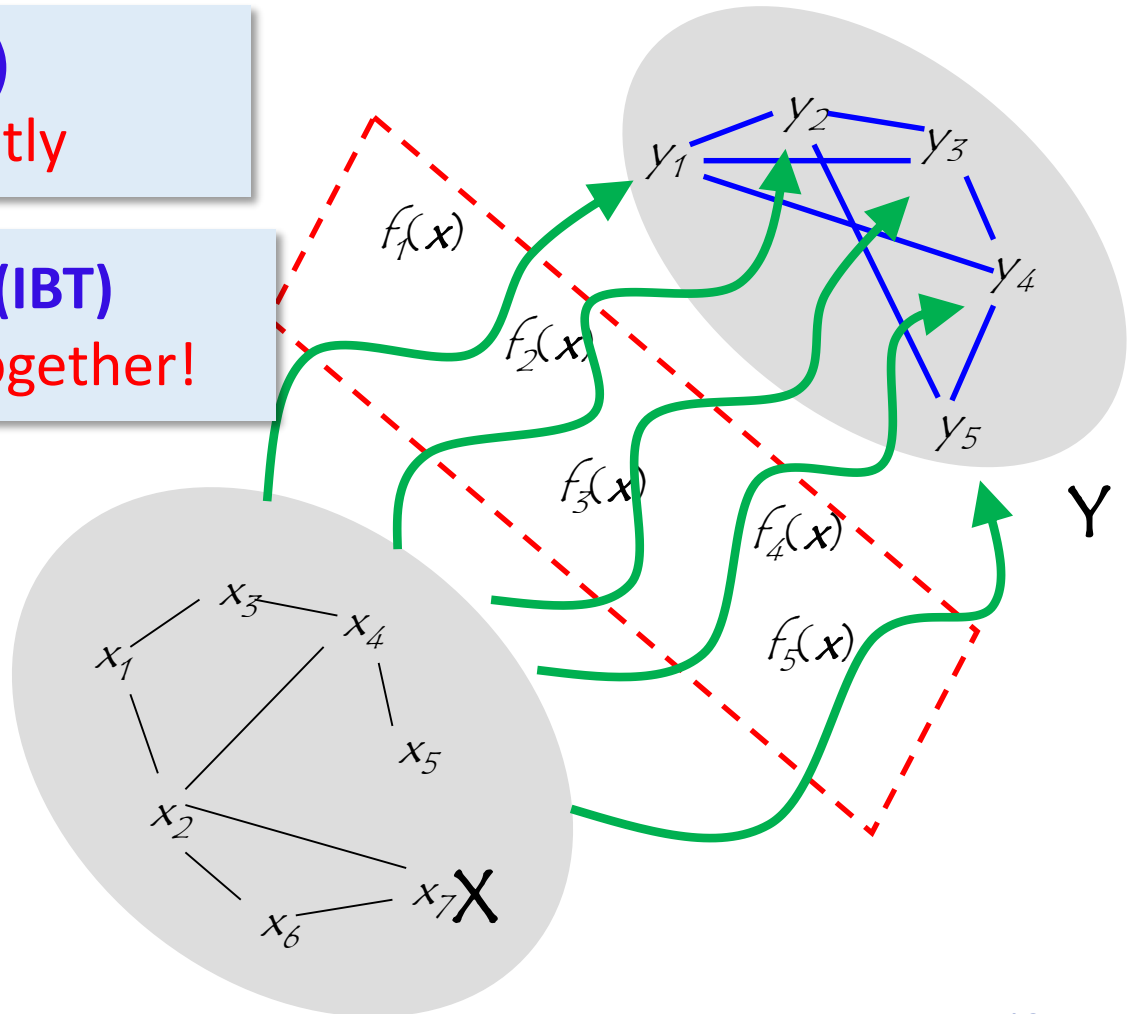
Learning + Inference (L+I)

Learn models **independently**

Inference Based Training (IBT)

Learn one model, all y 's **together!**

Intuition: Learning with constraints may make learning more difficult



Local Learning v.s. Global Learning

- ❖ Advantages of local training:

- ❖ Modular, very simple to implement
- ❖ Often efficient

- ❖ Advantages of global training:

- ❖ Models directly capture correlations between outputs (have better theoretical guarantees)
- ❖ Might achieve better performance when data is sufficient

Structured Prediction: Learning

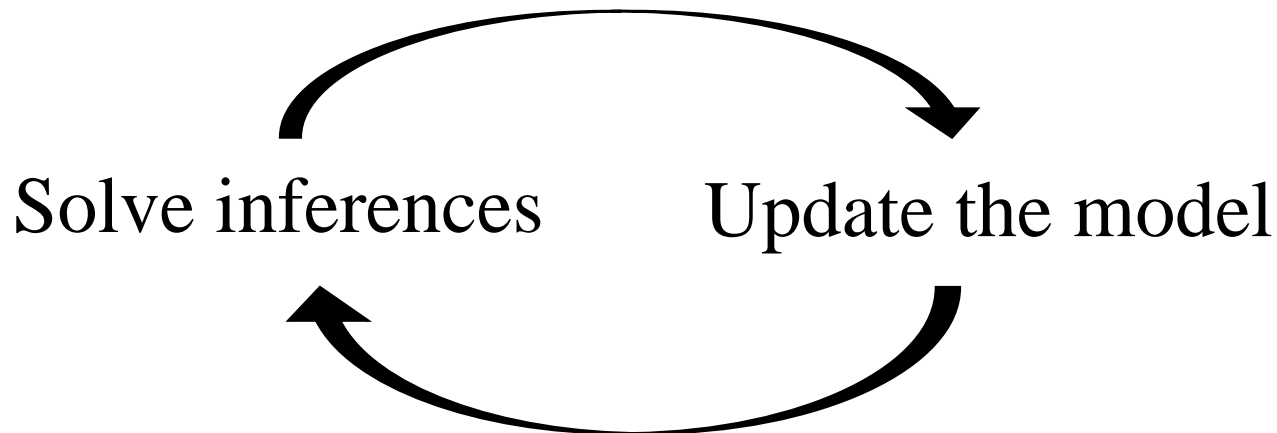
- ❖ Learning is thus driven by the attempt to find a weight vector w such that for each given annotated example (x_i, y_i) :

$$\text{Score of annotated structure} \geq \text{Score of any other structure} + \text{Penalty for predicting other structure}$$

- ❖ The update of the weight vector w can be done in an **on-line** or a **batch** fashion

Global Learning Algorithms

- ❖ **Online** e.g., Structured Perceptron [Collins 02]
 - ❖ Receive an instance; and update
- ❖ **Batch** e.g., Structured SVM [Taskar+05]
 - ❖ Collect a set of data; formulate learning as an optimization problem



Structured Perceptron

❖ Perceptron (binary classification):

For each training example (\mathbf{x}_i, y_i) :

$$\hat{y} \leftarrow \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$\text{if } y_i \neq \hat{y}, \text{ then } \mathbf{w} \leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i$$

❖ Structured Perceptron

For each training example $(\mathbf{x}_i, \mathbf{y}_i)$:

$$\hat{\mathbf{y}} \leftarrow \arg \max_{\mathbf{y} \in Y} \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i)$$

$$\mathbf{w} \leftarrow \mathbf{w} + \eta' [\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}})]$$

Perceptron is a special case when

$$\phi(\mathbf{x}, y) = y\mathbf{x}$$

$$\eta' = \frac{1}{2}\eta$$

Structured Perceptron Variations

❖ (Marginal) Structured Perceptron

For each training example $(\mathbf{x}_i, \mathbf{y}_i)$:

$$\hat{\mathbf{y}} \leftarrow \arg \max_{\mathbf{y} \in Y} \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) + \Delta(\mathbf{y}, \mathbf{y}_i)$$

$$\mathbf{w} \leftarrow \mathbf{w} + \eta' [\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}})]$$

❖ Parallel Structured Perceptron [McDonald et al 10]:

1. Split data into p parts.
2. Train Structured Perceptron on each data block in parallel.
3. Mixed the models using a linear combination.
4. Repeat Step 2 and use the mixed model as the initial model.

Recap

- ❖ Learning is thus driven by the attempt to find a weight vector \mathbf{w} such that for each given annotated example $(\mathbf{x}_i, \mathbf{y}_i)$:

$$\text{Score of annotated structure} \geq \text{Score of any other structure} + \text{Penalty for predicting other structure}$$

- ❖ Introduce slack variables $\{\xi_i\}$

$$\mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) \geq \Delta(\mathbf{y}_i, \mathbf{y}) - \xi_i$$

Structured SVM (Batch)

(Tsochantaridis et al. 05)

For all samples and feasible structures

Given a set of training examples $D = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^l$

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i$$

$$\text{s. t. } \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) \geq \Delta(\mathbf{y}_i, \mathbf{y}) - \xi_i \quad \forall i, \mathbf{y} \in Y_i$$

Note:

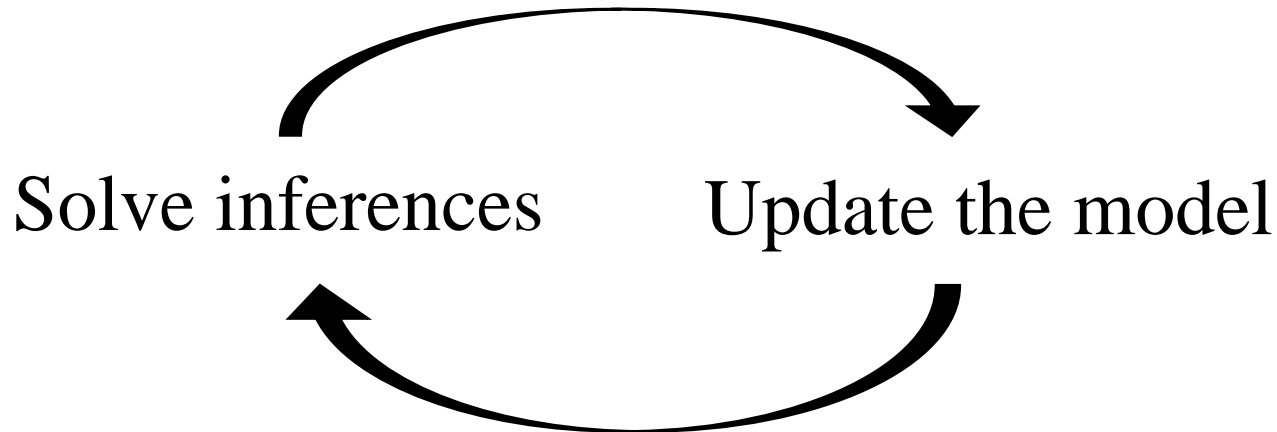
1. $\mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$ is the **scoring function** used in inference.
2. Equivalent to the following optimization problem:

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \left(\max_{\mathbf{y}} [\Delta(\mathbf{y}_i, \mathbf{y}) + \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y})] - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \right)$$

Loss-augmented inference

Optimization Methods for Structured SVM

- ❖ **Online** e.g., Structured Perceptron [Collins 02]
- ❖ **Batch** e.g., Structured SVM [Taskar+05]
 - Cutting plane: [Tsochantaridis+ 05, Joachims+ 09]
 - Dual Coordinate Descent: [Shevade+ 11, Chang+ 13]
 - Block-Coordinate Frank-Wolfe: [Lacoste-Julien+ 13]
 - Parallel Dual Coordinate Descent: [Chang+ 13a]



Stochastic (sub-)gradient descent

To minimize a function $g(z)$ that has the form $\sum_i g_i(z)$

- ❖ Initialize z_0
- ❖ Iterate until convergence
 - ❖ Pick a random g_i and compute its (sub)gradient at z_t : $\nabla g_i(z_t)$
 - ❖ Update: $z_{t+1} \leftarrow z_t - \gamma_t \nabla g_i(z_t)$ γ_t : learning rate

General idea: Replace the gradient with a noisy estimate

Sub-gradient computation

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \left(\max_y [\Delta(\mathbf{y}_i, \mathbf{y}) + \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y})] - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \right)$$

❖ Solve the max. Suppose solution is \mathbf{y}'

❖ The **loss-augmented/loss-sensitive/cost-augmented inference** step

❖ Compute gradient of

$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \left([\Delta(\mathbf{y}_i, \mathbf{y}) + \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y})] - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \right)$$

❖ Subgradient is

$$\mathbf{w} + C (\phi(\mathbf{x}_i, \mathbf{y}') - \phi(\mathbf{x}_i, \mathbf{y}_i))$$

SGD for structural SVM: The update

- ❖ Solve inference and compute sub-gradient:

$$\mathbf{w} + C (\phi(\mathbf{x}_i, \mathbf{y}') - \phi(\mathbf{x}_i, \mathbf{y}_i))$$

- ❖ At each step, go down the gradient:

$$\mathbf{w} \leftarrow \mathbf{w} - \gamma_t (\mathbf{w} + C (\phi(\mathbf{x}_i, \mathbf{y}') - \phi(\mathbf{x}_i, \mathbf{y}_i)))$$

- ❖ Compared to Structured Perception update:

$$\mathbf{w} \leftarrow \mathbf{w} + \eta' [\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}})]$$

(γ_t, η : learning rate)

L2-loss Structured SVM

Given a set of training examples $D = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^l$

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i^2$$

$$\text{s. t. } \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) \geq \Delta(\mathbf{y}_i, \mathbf{y}) - \xi_i \quad \forall i, \mathbf{y} \in Y_i$$

- Solve problem in the dual space
 - Recall that in dual SVM: one dual variable for each instance
 - In structured case: one for (instance, output assignment)

Dual Problem of Structural SVM

$\min_{\alpha > 0} D(\alpha)$, and

$$D(\alpha) \equiv \frac{1}{2} \left\| \sum_{\alpha_{i,y}} \alpha_{i,y} \phi(\mathbf{y}, \mathbf{y}_i, \mathbf{x}_i) \right\|^2 + \frac{1}{4C} \sum_i \left(\sum_{\mathbf{y}} \alpha_{i,y} \right)^2 - \sum_{i,y} \Delta(\mathbf{y}, \mathbf{y}_i) \alpha_{i,y}$$

where $\phi(\mathbf{y}, \mathbf{y}_i, \mathbf{x}_i) = \phi(\mathbf{y}_i, \mathbf{x}_i) - \phi(\mathbf{y}, \mathbf{x}_i)$.

- Number of α variables can be **exponentially large**.
- Relationship between w^* and α^*

Corresponds to different y .

$$w^* = \sum_{i,y} \alpha_{i,y}^* \phi(\mathbf{y}, \mathbf{y}_i, \mathbf{x}_i).$$

For linear model: maintain the relationship between w and α though out the learning process [Hsieh et.al. 08].

Structured Learning by Dual Coordinate Descent

- Number of dual variables can be **exponentially large**
- Maintain an **active set A** of dual variables:
 - Identify dual variables that will be likely non-zero
- Single-thread implementation:
 - Select and maintain A (active set selection step).

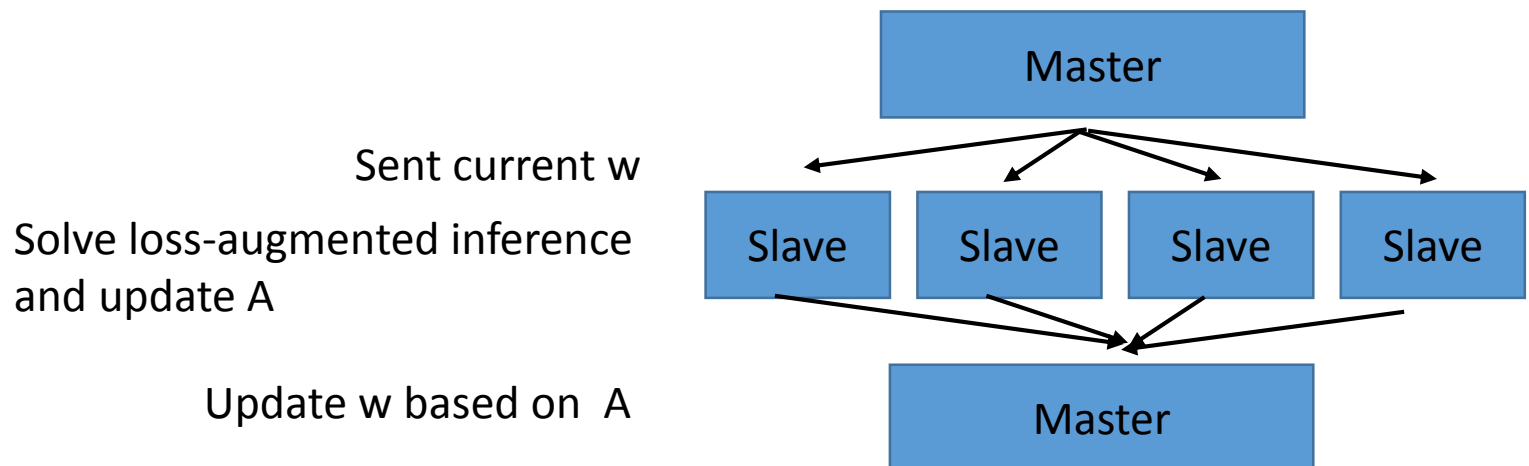
$$\max_{\mathbf{y} \in \mathcal{Y}_i} \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}_i, \mathbf{y})$$

- **Update the values of $\alpha_{i,y} \in A$ (learning step).**
 - (Approximately) solving a sub-problem.

A parallel Dual Coordinate Descent Algorithm

A **Master-Slave** architecture (MS-DCD):

- Given p processors: split data into p parts.
- At each iterations:
 - **Master sends current model to slave threads.**
 - **Each slave thread solves loss-augmented inference problems associated with a data block and updates the active set.**
 - **After all slave threads finish, master thread updates the model w according to the active set.**

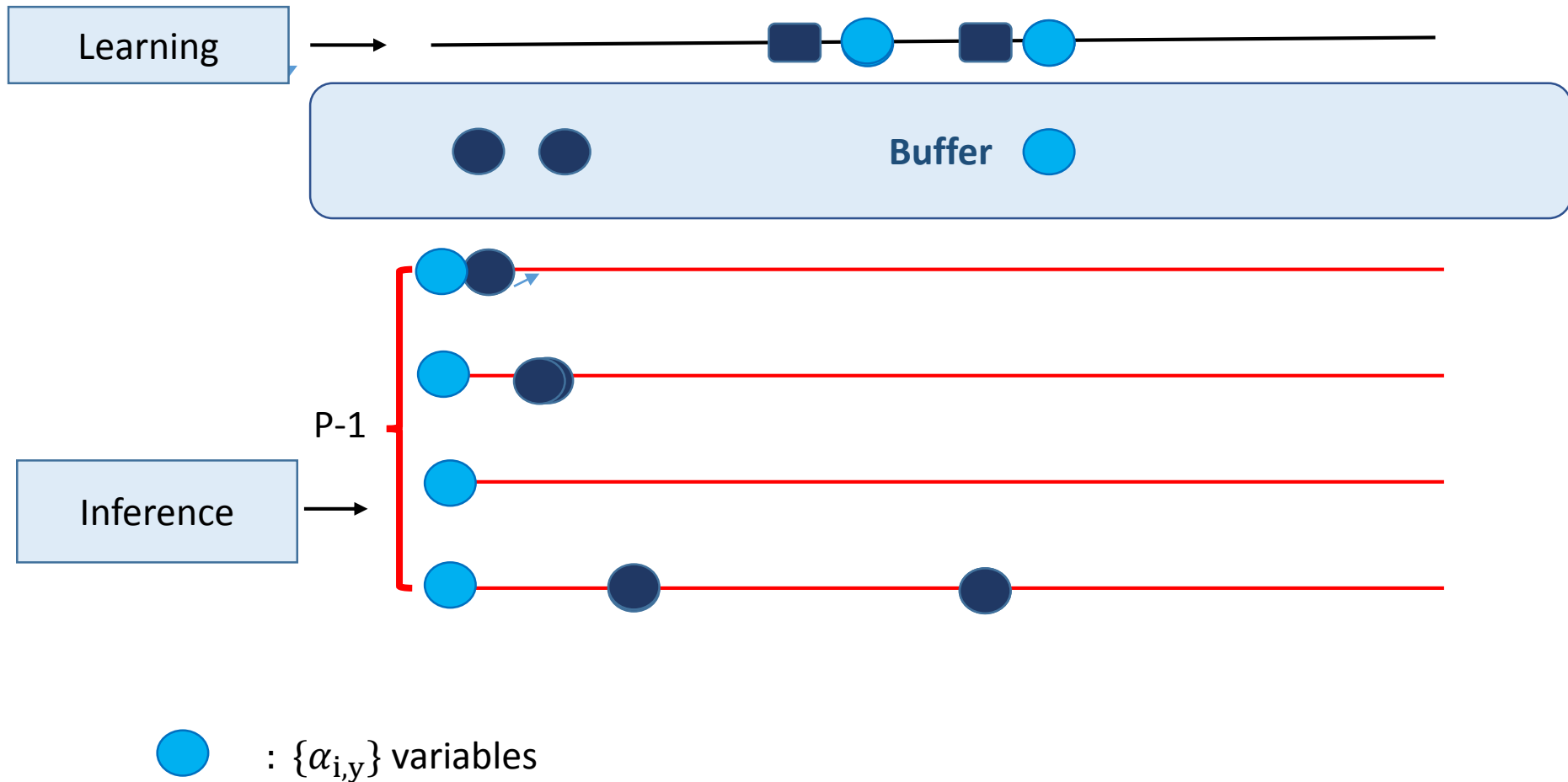


DEMI-DCD

DEMI-DCD: **Decouple Model-update and Inference** with Dual Coordinate Descent.

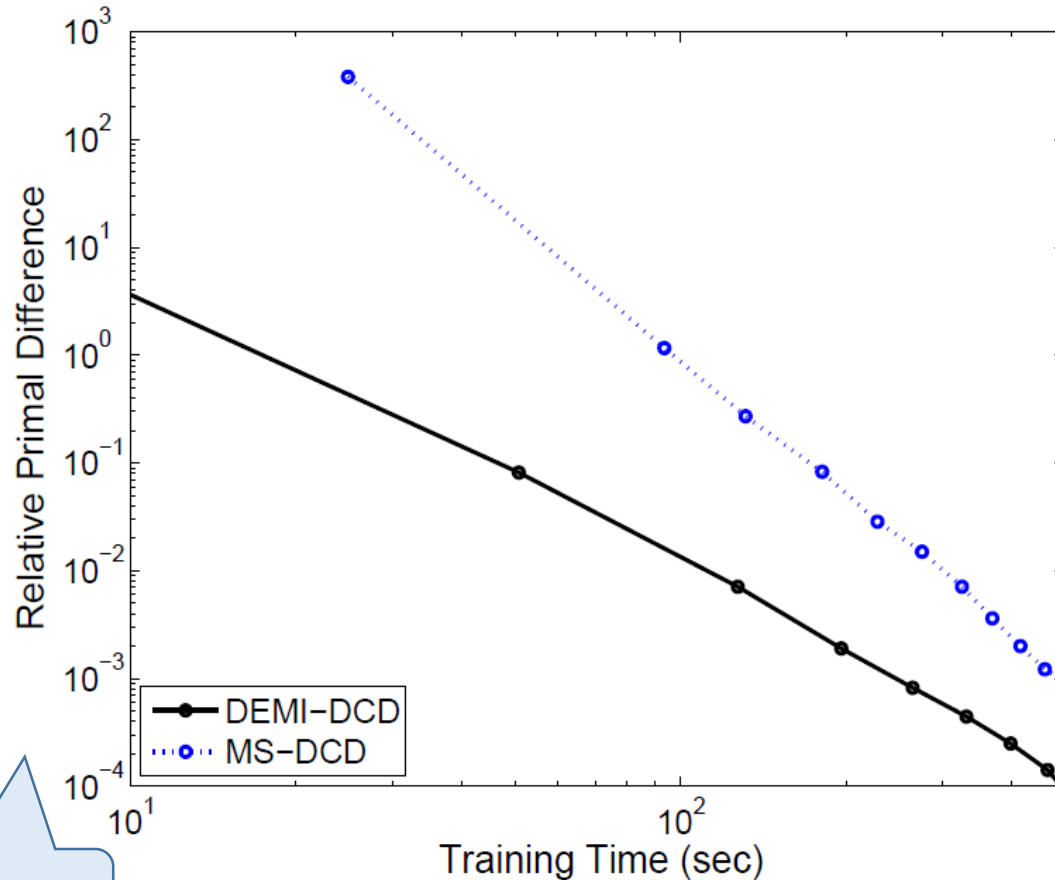
- Let p be #threads, and split training data into B_1, B_2, \dots, B_{p-1} .
- **Active set selection thread j** : select and maintain the active set A_i for each example i in B_j .
- **Learning thread**: loop over all examples and update model w .
- **A and w are shared** between threads using shared memory buffers.

Animation



Convergence on Primal Function Value

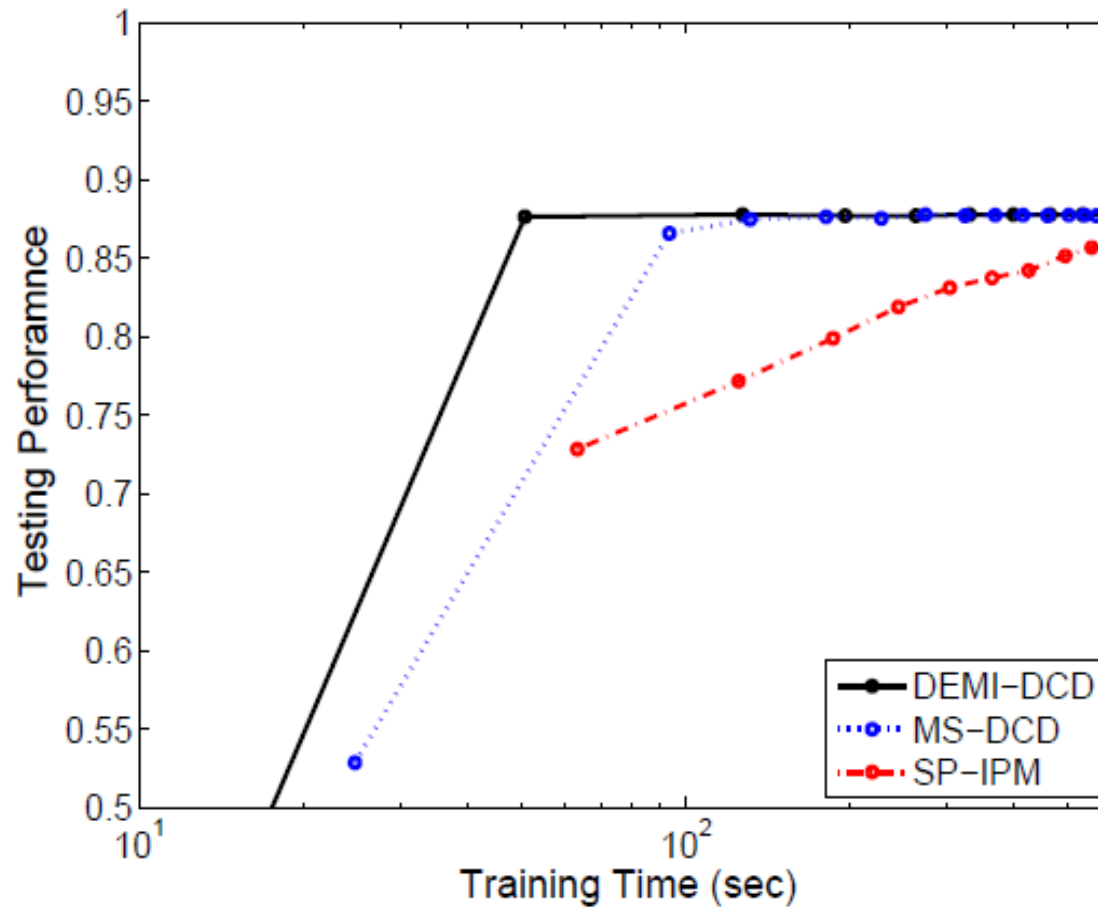
Relative primal function value difference along training time (Entity-Relation task)



3
Log-scale

Test Performance

Test Performance (F1 score) along training time



- **Part 2: Learning a Structured Prediction Model**
 - **Definition of structured learning**
 - **Local Learning v.s. Global Learning**
 - **Global Learning Algorithms**
 - Online learning: Structured Perceptron
 - Batch learning: Structured SVM
 - **Optimization methods for Structured SVM**
 - Stochastic Gradient Decent
 - Dual Coordinate Descent
 - Learning on a multi-core machine

In part 4, we will describe how to use them in practice