

# Training algorithms for Structured Learning

Sunita Sarawagi  
IIT Bombay

<http://www.cse.iitb.ac.in/~sunita>

# Training

Given

- $N$  input output pairs  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$
- Features  $\mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_c \mathbf{f}(\mathbf{x}, \mathbf{y}_c, c)$
- Error of output :  $E(\mathbf{y}_i, \mathbf{y})$ 
  - ▶ (Use short form:  $E(\mathbf{y}_i, \mathbf{y}) = E_i(\mathbf{y})$ )
  - ▶ Also decomposes over smaller parts:  
$$E_i(\mathbf{y}) = \sum_{c \in C} E_{i,c}(\mathbf{y}_c)$$

Find  $\mathbf{w}$

- Small training error
- Generalizes to unseen instances
- Efficient for structured models

# Outline

- 1 Likelihood based Training
- 2 Max-margin training
  - Decomposition-based approaches
  - Cutting-plane approaches

# Probability distribution from scores

- Convert scores into a probability distribution

$$\Pr(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

where  $Z_{\mathbf{w}}(\mathbf{x}) = \sum_{\mathbf{y}'} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}'))$

# Probability distribution from scores

- Convert scores into a probability distribution

$$\Pr(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

where  $Z_{\mathbf{w}}(\mathbf{x}) = \sum_{\mathbf{y}'} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}'))$

- When  $\mathbf{y}$  vector of variables, say  $y_1, \dots, y_n$ , and decomposition parts  $c$  are subsets of variables we get a graphical model.

$$\Pr(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp\left(\sum_c \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}_c, c)\right) = \frac{1}{Z} \prod_c \psi_c(\mathbf{y}_c)$$

with clique potential  $\psi_c(\mathbf{y}_c) = \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}_c, c))$

# Probability distribution from scores

- Convert scores into a probability distribution

$$\Pr(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

where  $Z_{\mathbf{w}}(\mathbf{x}) = \sum_{\mathbf{y}'} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}'))$

- When  $\mathbf{y}$  vector of variables, say  $y_1, \dots, y_n$ , and decomposition parts  $c$  are subsets of variables we get a graphical model.

$$\Pr(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp\left(\sum_c \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}_c, c)\right) = \frac{1}{Z} \prod_c \psi_c(\mathbf{y}_c)$$

with clique potential  $\psi_c(\mathbf{y}_c) = \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}_c, c))$

- These are called **Conditional Random Fields (CRFs)**.

# Training via gradient descent

$$L(\mathbf{w}) = \sum_{\ell} \log \Pr(\mathbf{y}_{\ell} | \mathbf{x}_{\ell}, \mathbf{w}) = \sum_{\ell} (\mathbf{w} \cdot \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}_{\ell}) - \log Z_{\mathbf{w}}(\mathbf{x}_{\ell}))$$

Add a regularizer to prevent over-fitting.

$$\max_{\mathbf{w}} \sum_{\ell} (\mathbf{w} \cdot \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}_{\ell}) - \log Z_{\mathbf{w}}(\mathbf{x}_{\ell})) - \|\mathbf{w}\|^2 / C$$

Concave in  $\mathbf{w} \implies$  gradient descent methods will work.

Gradient:

$$\begin{aligned} \nabla L(\mathbf{w}) &= \sum_{\ell} \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}_{\ell}) - \frac{\sum_{\mathbf{y}'} \mathbf{f}(\mathbf{y}', \mathbf{x}_{\ell}) \exp \mathbf{w} \cdot \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}')}{Z_{\mathbf{w}}(\mathbf{x}_{\ell})} - 2\mathbf{w} / C \\ &= \sum_{\ell} \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}_{\ell}) - E_{\Pr(\mathbf{y}' | \mathbf{w})} \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}') - 2\mathbf{w} / C \end{aligned}$$

# Training algorithm

1: Initialize  $\mathbf{w}^0 = \mathbf{0}$



# Training algorithm

- 1: Initialize  $\mathbf{w}^0 = \mathbf{0}$
- 2: **for**  $t = 1 \dots T$  **do**
- 3:     **for**  $\ell = 1 \dots N$  **do**
- 4:          $\mathbf{g}_{k,\ell} = f_k(\mathbf{x}_\ell, \mathbf{y}_\ell) - E_{\text{Pr}(\mathbf{y}'|\mathbf{w})} f_k(\mathbf{x}_\ell, \mathbf{y}') \quad k = 1 \dots K$
- 5:     **end for**
- 6:      $\mathbf{g}_k = \sum_\ell \mathbf{g}_{k,\ell} \quad k = 1 \dots K$

# Training algorithm

- 1: Initialize  $\mathbf{w}^0 = \mathbf{0}$
- 2: **for**  $t = 1 \dots T$  **do**
- 3:     **for**  $\ell = 1 \dots N$  **do**
- 4:          $\mathbf{g}_{k,\ell} = f_k(\mathbf{x}_\ell, \mathbf{y}_\ell) - E_{\text{Pr}(\mathbf{y}'|\mathbf{w})} f_k(\mathbf{x}_\ell, \mathbf{y}')$       $k = 1 \dots K$
- 5:     **end for**
- 6:      $\mathbf{g}_k = \sum_\ell \mathbf{g}_{k,\ell}$       $k = 1 \dots K$
- 7:      $\mathbf{w}_k^t = \mathbf{w}_k^{t-1} + \gamma_t (\mathbf{g}_k - 2\mathbf{w}_k^{t-1} / C)$
- 8:     **Exit** if  $\|\mathbf{g}\| \approx \text{zero}$
- 9: **end for**

# Training algorithm

- 1: Initialize  $\mathbf{w}^0 = \mathbf{0}$
- 2: **for**  $t = 1 \dots T$  **do**
- 3:     **for**  $\ell = 1 \dots N$  **do**
- 4:          $\mathbf{g}_{k,\ell} = f_k(\mathbf{x}_\ell, \mathbf{y}_\ell) - E_{\text{Pr}(\mathbf{y}'|\mathbf{w})} f_k(\mathbf{x}_\ell, \mathbf{y}')$       $k = 1 \dots K$
- 5:     **end for**
- 6:      $\mathbf{g}_k = \sum_{\ell} \mathbf{g}_{k,\ell}$       $k = 1 \dots K$
- 7:      $\mathbf{w}_k^t = \mathbf{w}_k^{t-1} + \gamma_t(\mathbf{g}_k - 2\mathbf{w}_k^{t-1}/C)$
- 8:     **Exit** if  $\|\mathbf{g}\| \approx \text{zero}$
- 9: **end for**

Running time of the algorithm is  $O(INn(m^2 + K))$  where  $I$  is the total number of iterations.

Calculating  $E_{\text{Pr}(\mathbf{y}'|\mathbf{w})} f_k(\mathbf{x}_\ell, \mathbf{y}')$  using inference.

# Likelihood-based trainer

- 1 Penalizes all wrong  $\mathbf{y}$ s the same way, does not exploit  $E_i(\mathbf{y})$
- 2 Requires the computation of sum-marginals, not possible in all kinds of structured learning.
  - 1 Collective extraction
  - 2 Sentence Alignment
  - 3 Ranking

# Outline

- 1 Likelihood based Training
- 2 Max-margin training
  - Decomposition-based approaches
  - Cutting-plane approaches

# Two formulations

## 1 Margin scaling

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i$$

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

# Two formulations

## 1 Margin scaling

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \geq E_i(\mathbf{y}) - \xi_i \quad \forall \mathbf{y}, i \end{aligned}$$

## 2 Slack scaling

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \geq 1 - \frac{\xi_i}{E_i(\mathbf{y})} \quad \forall \mathbf{y}, i \end{aligned}$$



# Max-margin loss surrogates

True error  $E_i(\operatorname{argmax}_{\mathbf{y}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}_i, \mathbf{y}))$

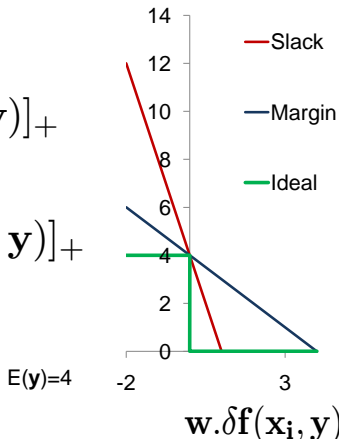
Let  $\mathbf{w} \cdot \delta \mathbf{f}(\mathbf{x}_i, \mathbf{y}) = \mathbf{w} \cdot \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w} \cdot \mathbf{f}(\mathbf{x}_i, \mathbf{y})$

1. Margin Loss

$$\max_{\mathbf{y}} [E_i(\mathbf{y}) - \mathbf{w} \cdot \delta \mathbf{f}(\mathbf{x}_i, \mathbf{y})]_+$$

2. Slack Loss

$$\max_{\mathbf{y}} E_i(\mathbf{y}) [1 - \mathbf{w} \cdot \delta \mathbf{f}(\mathbf{x}_i, \mathbf{y})]_+$$



# Max-margin training: margin-scaling

The Primal (P):

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \mathbf{w}^T \delta \mathbf{f}_i(\mathbf{y}) \geq E_i(\mathbf{y}) - \xi_i \quad \forall \mathbf{y}, i : 1 \dots N \end{aligned}$$

# Max-margin training: margin-scaling

The Primal (P):

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i$$

$$\text{s.t. } \mathbf{w}^T \delta \mathbf{f}_i(\mathbf{y}) \geq E_i(\mathbf{y}) - \xi_i \quad \forall \mathbf{y}, i : 1 \dots N$$

- Good news: Convex in  $\mathbf{w}, \xi$
- Bad news: exponential number of constraints
- Two main lines of attacks
  - 1 Decomposition: polynomial-sized rewrite of objective in terms of parts of  $\mathbf{y}$
  - 2 Cutting-plane: generate constraints on the fly.

1 The Primal (P):

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i$$
$$\text{s.t. } \mathbf{w}^T \delta \mathbf{f}_i(\mathbf{y}) \geq E_i(\mathbf{y}) - \xi_i \quad \forall \mathbf{y}, i : 1 \dots N$$

1 The Primal (P):

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i$$
$$\text{s.t. } \mathbf{w}^T \delta \mathbf{f}_i(\mathbf{y}) \geq E_i(\mathbf{y}) - \xi_i \quad \forall \mathbf{y}, i : 1 \dots N$$

2 The Dual (D) of (P)

$$\max_{\alpha_i(\mathbf{y})} -\frac{1}{2} \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \delta \mathbf{f}_i(\mathbf{y}) \sum_{j, \mathbf{y}'} \alpha_j(\mathbf{y}') \delta \mathbf{f}_j(\mathbf{y}') + \sum E_i(\mathbf{y}) \alpha_i$$
$$\text{s.t. } \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = \frac{C}{N}$$
$$\alpha_i(\mathbf{y}) \geq 0 \quad i : 1 \dots N$$

# Properties of Dual

- 1 Strong duality holds: Primal (P) solution = Dual (D) solution.
- 2  $\mathbf{w} = \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \delta \mathbf{f}_i(\mathbf{y})$
- 3 Dual (D) is concave in  $\alpha$ , constraints are simpler.
- 4 Size of  $\alpha$  is still intractably large  $\implies$  cannot solve via standard libraries.

# Decomposition-based approaches

$$\textcircled{1} \quad \delta \mathbf{f}_i(\mathbf{y}) = \sum_c \delta \mathbf{f}_{i,c}(\mathbf{y}_c)$$

$$\textcircled{2} \quad E_i(\mathbf{y}) = \sum_c E_{i,c}(\mathbf{y}_c)$$

# Decomposition-based approaches

- 1  $\delta \mathbf{f}_i(\mathbf{y}) = \sum_c \delta \mathbf{f}_{i,c}(\mathbf{y}_c)$
- 2  $E_i(\mathbf{y}) = \sum_c E_{i,c}(\mathbf{y}_c)$

Rewrite the dual as

$$\begin{aligned} \max_{\mu_{i,c}(\mathbf{y}_c)} & -\frac{1}{2} \sum_{i,c,\mathbf{y}_c} \delta \mathbf{f}_{i,c}(\mathbf{y}_c) \mu_{i,c}(\mathbf{y}_c) \sum_{j,d,\mathbf{y}'_d} \delta \mathbf{f}_{j,d}(\mathbf{y}'_d) \mu_{j,d}(\mathbf{y}'_d) \\ & + \sum_{i,c,\mathbf{y}_c} E_{i,c}(\mathbf{y}_c) \mu_{i,c}(\mathbf{y}_c) \\ \text{s.t.} & \sum_{\mathbf{y}} \mu_{i,c}(\mathbf{y}_c) = \sum_{\mathbf{y} \sim \mathbf{y}_c} \alpha_i(\mathbf{y}) \\ & \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = \frac{C}{N}, \alpha_i(\mathbf{y}) \geq 0 \quad i : 1 \dots N \end{aligned}$$



## $\alpha$ s as probabilities

Scale  $\alpha$ s with  $\frac{C}{N}$ .

$$\begin{aligned} \max_{\mu_{i,c}(\mathbf{y}_c)} & -\frac{C}{2N} \sum_{i,c,\mathbf{y}_c} \delta \mathbf{f}_{i,c}(\mathbf{y}_c) \mu_{i,c}(\mathbf{y}_c) \sum_{j,d,\mathbf{y}'_d} \delta \mathbf{f}_{j,d}(\mathbf{y}'_d) \mu_{j,d}(\mathbf{y}'_d) \\ & + \sum_{i,c,\mathbf{y}_c} E_{i,c}(\mathbf{y}_c) \mu_{i,c}(\mathbf{y}_c) \end{aligned}$$

s.t.  $\mu_{i,c}(\mathbf{y}_c) \in$  Marginals of any valid distribution

## $\alpha$ s as probabilities

Scale  $\alpha$ s with  $\frac{C}{N}$ .

$$\begin{aligned} \max_{\mu_{i,c}(\mathbf{y}_c)} & -\frac{C}{2N} \sum_{i,c,\mathbf{y}_c} \delta \mathbf{f}_{i,c}(\mathbf{y}_c) \mu_{i,c}(\mathbf{y}_c) \sum_{j,d,\mathbf{y}'_d} \delta \mathbf{f}_{j,d}(\mathbf{y}'_d) \mu_{j,d}(\mathbf{y}'_d) \\ & + \sum_{i,c,\mathbf{y}_c} E_{i,c}(\mathbf{y}_c) \mu_{i,c}(\mathbf{y}_c) \end{aligned}$$

s.t.  $\mu_{i,c}(\mathbf{y}_c) \in$  Marginals of any valid distribution

Solve via the exponentiated gradient method.

# Exponentiated gradient algorithm

- 1 Initially  $\mu_{i,c}(\mathbf{y}_{i,c}) = 1$ , for  $\mathbf{y}_{i,c} \neq \mathbf{y}_c$ ,  $\mu_{i,c}(\mathbf{y}_c) = 0$
- 2 For  $t = 1, \dots, T$ 
  - 1 Choose a  $i$  from  $1, \dots, N$ .

# Exponentiated gradient algorithm

- 1 Initially  $\mu_{i,c}(\mathbf{y}_{i,c}) = 1$ , for  $\mathbf{y}_{i,c} \neq \mathbf{y}_c$ ,  $\mu_{i,c}(\mathbf{y}_c) = 0$
- 2 For  $t = 1, \dots, T$ 
  - 1 Choose a  $i$  from  $1, \dots, N$ .
  - 2 Ignore constraints and perform a gradient-based update:

$$s_{i,c} = \mu_{i,c}^t + \eta(E_{i,c} - \mathbf{w}^t \delta \mathbf{f}_i(\mathbf{y}_c))$$

$$\text{where } \mathbf{w}^t = \sum_{i,c,\mathbf{y}_c} \mu_{i,c}^t(\mathbf{y}_c) \delta \mathbf{f}_{i,c}(\mathbf{y}_c)$$

# Exponentiated gradient algorithm

- 1 Initially  $\mu_{i,c}(\mathbf{y}_{i,c}) = 1$ , for  $\mathbf{y}_{i,c} \neq \mathbf{y}_c$ ,  $\mu_{i,c}(\mathbf{y}_c) = 0$
- 2 For  $t = 1, \dots, T$ 
  - 1 Choose a  $i$  from  $1, \dots, N$ .
  - 2 Ignore constraints and perform a gradient-based update:

$$s_{i,c} = \mu_{i,c}^t + \eta(E_{i,c} - \mathbf{w}^t \delta \mathbf{f}_i(\mathbf{y}_c))$$

where  $\mathbf{w}^t = \sum_{i,c,\mathbf{y}_c} \mu_{i,c}^t(\mathbf{y}_c) \delta \mathbf{f}_{i,c}(\mathbf{y}_c)$

- 3 Define a distribution  $\alpha$  by exponentiating the updates:

$$\alpha_i(\mathbf{y})^{t+1} = \frac{1}{Z} \exp\left(\sum_c s_{i,c}(\mathbf{y}_c)\right)$$

where  $Z = \sum_{\mathbf{y}} \exp(\sum_c s_{i,c}(\mathbf{y}_c))$

# Exponentiated gradient algorithm

- 1 Initially  $\mu_{i,c}(\mathbf{y}_{i,c}) = 1$ , for  $\mathbf{y}_{i,c} \neq \mathbf{y}_c$ ,  $\mu_{i,c}(\mathbf{y}_c) = 0$
- 2 For  $t = 1, \dots, T$

- 1 Choose a  $i$  from  $1, \dots, N$ .
- 2 Ignore constraints and perform a gradient-based update:

$$s_{i,c} = \mu_{i,c}^t + \eta(E_{i,c} - \mathbf{w}^t \delta \mathbf{f}_i(\mathbf{y}_c))$$

where  $\mathbf{w}^t = \sum_{i,c,\mathbf{y}_c} \mu_{i,c}^t(\mathbf{y}_c) \delta \mathbf{f}_{i,c}(\mathbf{y}_c)$

- 3 Define a distribution  $\alpha$  by exponentiating the updates:

$$\alpha_i(\mathbf{y})^{t+1} = \frac{1}{Z} \exp\left(\sum_c s_{i,c}(\mathbf{y}_c)\right)$$

where  $Z = \sum_{\mathbf{y}} \exp(\sum_c s_{i,c}(\mathbf{y}_c))$

- 4 New feasible values are marginals of  $\alpha$

$$\mu_{i,c}^{t+1}(\mathbf{y}_c) = \sum_{\mathbf{y} \sim \mathbf{y}_c} \alpha_i(\mathbf{y})^{t+1}$$

# Convergence results

## Theorem

$J(\alpha^{t+1}) - J(\alpha^t) \geq \frac{1}{\eta} KL(\alpha^t, \alpha^{t+1})$  where  $\eta \leq \frac{1}{nR^2}$  where  
 $R = \max \delta \mathbf{f}_i(\mathbf{y}) \delta \mathbf{f}_j(\mathbf{y}')$

# Convergence results

## Theorem

$J(\alpha^{t+1}) - J(\alpha^t) \geq \frac{1}{\eta} KL(\alpha^t, \alpha^{t+1})$  where  $\eta \leq \frac{1}{nR^2}$  where  
 $R = \max \delta \mathbf{f}_i(\mathbf{y}) \delta \mathbf{f}_j(\mathbf{y}')$

## Theorem

Let  $\alpha^*$  = dual optimal. Then at the  $T$ th iteration.

$$J(\alpha^*) - \frac{1}{T\eta} KL(\alpha^*; \alpha^0) \leq J(\alpha^{T+1}) \leq J(\alpha^*)$$



# Convergence results

## Theorem

$J(\alpha^{t+1}) - J(\alpha^t) \geq \frac{1}{\eta} \text{KL}(\alpha^t, \alpha^{t+1})$  where  $\eta \leq \frac{1}{nR^2}$  where  
 $R = \max \delta \mathbf{f}_i(\mathbf{y}) \delta \mathbf{f}_j(\mathbf{y}')$

## Theorem

Let  $\alpha^*$  = dual optimal. Then at the  $T$ th iteration.

$$J(\alpha^*) - \frac{1}{T\eta} \text{KL}(\alpha^*; \alpha^0) \leq J(\alpha^{T+1}) \leq J(\alpha^*)$$

## Theorem

The number of iterations of the algorithm is at most  
 $\frac{N^2}{\epsilon} R^2 \text{KL}(\alpha^*; \alpha^0)$

# Cutting plane method

- 1 Exponentiated gradient approach requires computation of sum-marginals and decomposable losses.
- 2 Cutting plane — a more general approach that just requires MAP

# Cutting-plane algorithm [TJHA05]

1: Initialize  $\mathbf{w}^0 = \mathbf{0}$ , Active constraints=Empty.

# Cutting-plane algorithm [TJHA05]

- 1: Initialize  $\mathbf{w}^0 = \mathbf{0}$ , Active constraints=Empty.
- 2: **for**  $t = 1 \dots T$  **do**
- 3:     **for**  $\ell = 1 \dots N$  **do**
- 4:          $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} (E_{\ell}(\mathbf{y}) + \mathbf{w}^t \cdot \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}))$
- 5:         **if**  $\mathbf{w}^t \cdot \delta \mathbf{f}_{\ell}(\hat{\mathbf{y}}) < E_{\ell}(\hat{\mathbf{y}}) - \xi_{\ell}^t - \epsilon$  **then**
- 6:             Add  $(\mathbf{x}_{\ell}, \hat{\mathbf{y}})$  to set of active constraints.
- 7:              $\mathbf{w}^t, \xi^t = \text{solve QP with active constraints.}$
- 8:         **end if**
- 9:     **end for**
- 10:     **Exit** if no new constraint added.
- 11: **end for**

# Efficient solution in the dual space

Solve QP in the dual space.

- 1 Initially  $\alpha_i^t(\mathbf{y}) = 0, \forall \mathbf{y} \neq \mathbf{y}_i, \alpha_i^t(\mathbf{y}_i) = \frac{C}{N}$
- 2 For  $t = 1, \dots, T$ 
  - 1 Choose a  $i$  from  $1, \dots, N$ .
  - 2  $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} (E_{\ell}(\mathbf{y}) + \mathbf{w}^t \cdot \mathbf{f}(\mathbf{x}_{\ell}, \mathbf{y}))$  where  $\mathbf{w}^t = \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \delta \mathbf{f}_i(\mathbf{y})$
  - 3  $\alpha_i(\hat{\mathbf{y}}) =$  coordinate with highest gradient.
  - 4 Optimize  $J(\alpha)$  over set of  $\mathbf{y}$ s in the active set (SMO applicable here).

# Convergence results

Let  $R^2 = \max \delta \mathbf{f}_i(\mathbf{y}) \delta \mathbf{f}_j(\mathbf{y}')$ ,  $\Delta = \max_{i, \mathbf{y}} E_i(\mathbf{y})$

## Theorem

$$J(\alpha^{t+1}) - J(\alpha^t) \geq \min\left(\frac{C\epsilon}{2N}, \frac{\epsilon^2}{8R^2}\right)$$

## Theorem

*The number of constraints that the cutting plane algorithm adds is at most  $\max\left(\frac{2N\Delta}{\epsilon}, \frac{8C\Delta R^2}{\epsilon^2}\right)$*

# Single slack formulation [JFY09]

## Theorem

*The number of constraints that the cutting plane algorithm adds in the single slack formulation is at most*

$$\max\left(\log \frac{C\Delta}{4R^2C^2}, \frac{16CR^2}{\epsilon}\right)$$

# Summary

- 1 Two very efficient algorithms for training structured models that avoids the problem of exponential output space.
- 2 Other alternatives
  - 1 Online training, example MIRA and Collins Trainer
  - 2 Stochastic trainers: LARank.
  - 3 Local training: SEARN
- 3 Extension to Slack-scaling and other loss functions.



 Peter L. Bartlett, Michael Collins, Ben Taskar, and David McAllester.

Exponentiated gradient algorithms for large-margin structured classification.

In Lawrence K. Saul, Yair Weiss, and Léon Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 113–120. MIT Press, Cambridge, MA, 2005.

 Michael Collins, Amir Globerson, Terry Koo, Xavier Carreras, and Peter L. Bartlett.

Exponentiated gradient algorithms for conditional random fields and max-margin markov networks.

*J. Mach. Learn. Res.*, 9:1775–1822, 2008.

 Thorsten Joachims, Thomas Finley, and Chun-Nam John Yu.

Cutting-plane training of structural svms.

*Machine Learning*, 77(1):27–59, 2009.



I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun.

Large margin methods for structured and interdependent output variables.

*Journal of Machine Learning Research (JMLR)*,  
6(Sep):1453–1484, 2005.