
Structured output models for image segmentation

Aurelien Lucchi

Machine Learning Workshop (MLWS) IDIAP - EPFL
Monday November 19th, 2012

Collaborators: Yunpeng Li, Kevin Smith, Raphael
Sznitman, Bohumil Maco, Graham Knott, Pascal Fua.

Outline

1. Review Conditional Random Fields (CRF)
2. Maximum likelihood training for CRFs
3. Maximum Margin Training for CRFs
 1. Cutting plane (Structured SVM)
 2. Online subgradient descent

1. Review CRF

Structured prediction

- Non structured output

$$f : X \rightarrow R$$

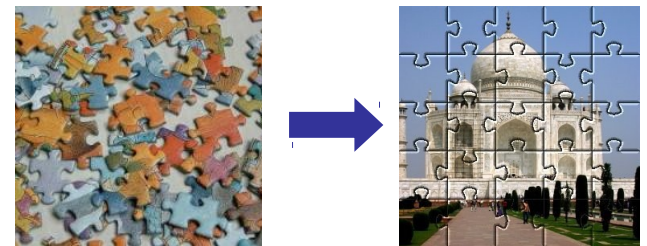
- inputs X can be any kind of objects
- output y is a real number

$$y = \{-1, +1\}$$

- Prediction of complex outputs

$$f : X \rightarrow Y$$

- Structured output y is complex (images, text, audio...)
- Ad hoc definition of structured data: data that consists of several parts, and not only the parts themselves contain information, but also the way in which the parts belong together



Structured prediction for image segmentation

X

$(x_1, \dots, x_i, \dots, x_n)$

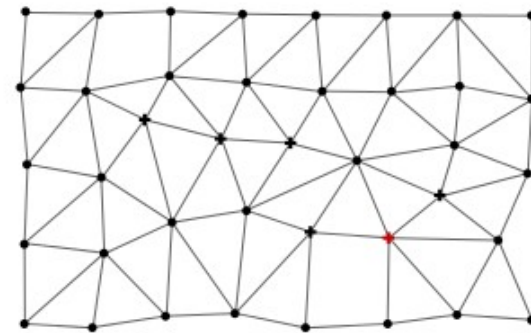


$$x_i \in \mathbb{R}^F$$

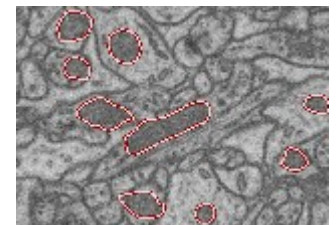
Histograms, Filter responses, ...

Y

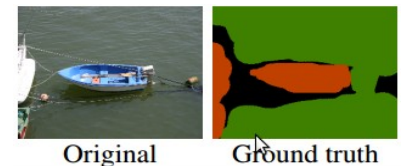
$(y_1, \dots, y_i, \dots, y_n)$



$$y_i = \{-1, +1\}$$



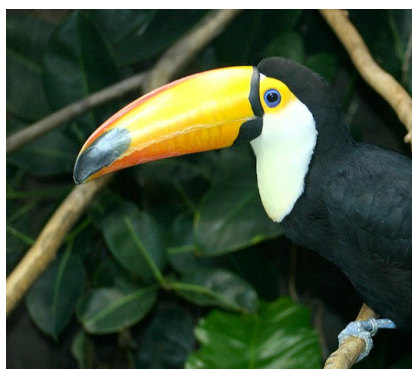
$$y_i = \{1, \dots, 21\}$$



CRF for image segmentation

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} V(y_i, y_j)$$

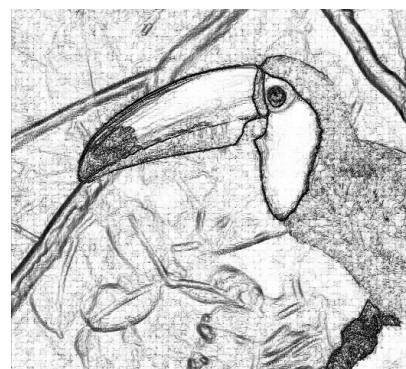
Maximum-a-posteriori (MAP) solution : $Y^* = \arg \min_{Y \in \mathcal{Y}} E_w(X, Y)$



Data (**D**)



Unary likelihood



Pair-wise Terms

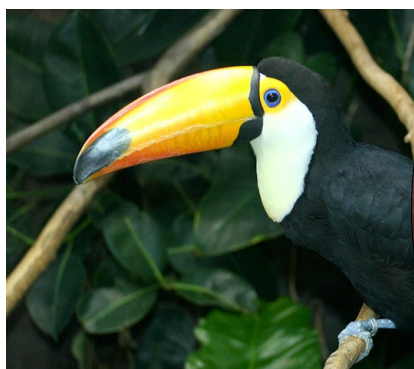


MAP Solution

CRF for image segmentation

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} V(y_i, y_j)$$

Maximum-a-posteriori (MAP) solution : $Y^* = \arg \min_{Y \in \mathcal{Y}} E_w(X, Y)$



Data (**D**)



Unary likelihood



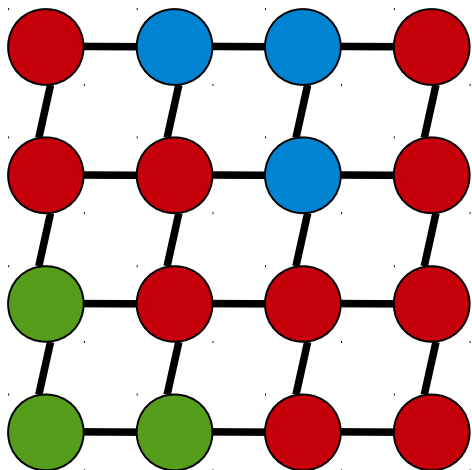
Pair-wise Terms



MAP Solution

CRF for image segmentation

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} V(y_i, y_j)$$



Pair-wise Terms

Favors the same label for neighboring nodes.

CRF for image segmentation

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} V(y_i, y_j)$$

Maximum-a-posteriori (MAP) solution :

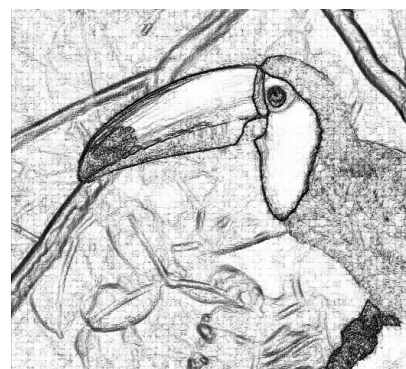
$$Y^* = \arg \min_{Y \in \mathcal{Y}} E_w(X, Y)$$



Data (**D**)



Unary likelihood



Pair-wise Terms



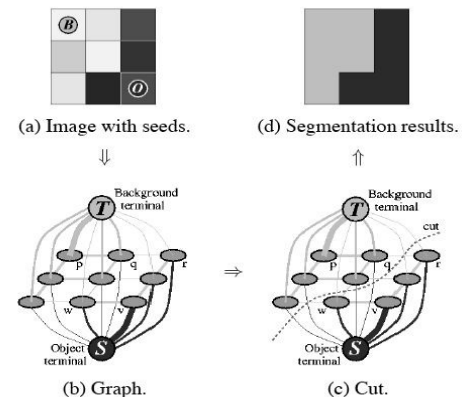
MAP Solution

Energy minimization

- MAP inference for discrete graphical models:

$$Y^* = \arg \min_{Y \in \mathcal{Y}} E_w(X, Y)$$

- Dynamic programming
 - Exact on non loopy graphs
- Graph-cuts (Boykov, 2001)
 - Optimal solution if energy function is submodular
- Belief propagation (Pearl, 1982)
 - No theoretical guarantees on loopy graphs but seems to work well in practice.
- Mean field (root in statistical physics)



Training a structured model ?

- First rewrite the energy function as:

$$\begin{aligned} E_w(X, Y) &= \sum_{i \in \mathcal{V}} D(y_i) + \sum_{i, j \in \mathcal{E}} V(y_i, y_j) \\ &= w^T \psi(X, Y) \end{aligned}$$

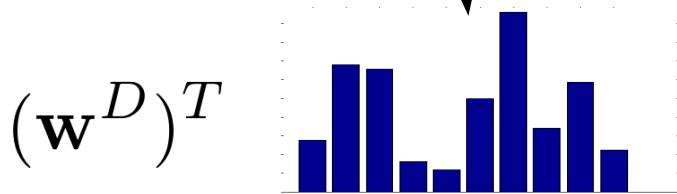
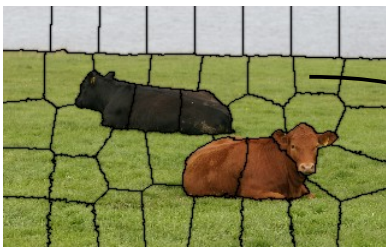
Log-linear model

- Efficient Learning/Training – need to efficiently learn parameters \mathbf{w} from training data ?

Training a structured model ?

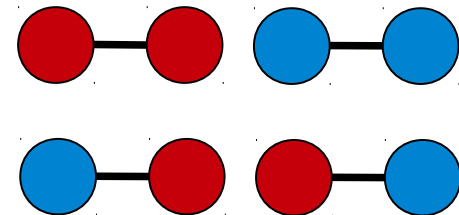
- Energy function is parametrized by vector \mathbf{w}

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{i, j \in \mathcal{E}} V(y_i, y_j) = \mathbf{w}^T \psi(X, Y)$$



+

\mathbf{w}^P	-1	1
-1	?	?
1	?	?



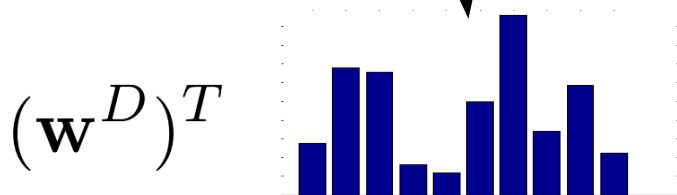
$$D(y_i) = (\mathbf{w}^D)^T x_i$$

$$V(y_i, y_j) = \mathbf{w}^V(y_i, y_j)$$

Training a structured model ?

- Energy function is parametrized by vector \mathbf{w}

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{i, j \in \mathcal{E}} V(y_i, y_j) = \mathbf{w}^T \psi(X, Y)$$



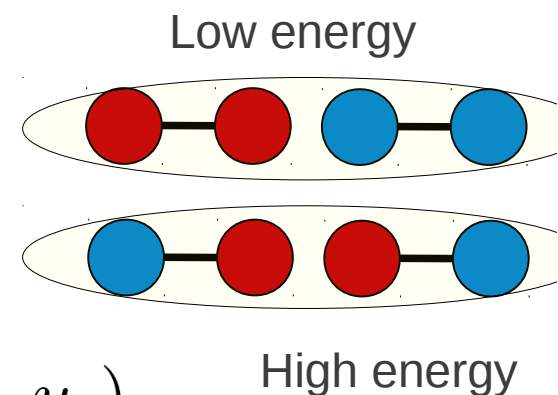
$$D(y_i) = (\mathbf{w}^D)^T x_i$$

$$\mathbf{w} = ((\mathbf{w}^D)^T, (\mathbf{w}^V)^T)^T$$

+

\mathbf{w}^P	-1	1
-1	0	1
1	1	0

$$V(y_i, y_j) = \mathbf{w}^V(y_i, y_j)$$



2. Maximum likelihood training



Maximum likelihood

$$\begin{aligned}w^* &= \arg \max_w L(w) = \arg \max_w \log p(Y|X, w) \\&= \arg \max_w \prod_n \log p(Y^n|X^n, w) \\&= \arg \max_w \sum_n \log p(Y^n|X^n, w) \\&= \arg \max_w \sum_n \log \frac{1}{Z(w)} \exp^{-E(Y^n; X^n)} \\&= \arg \max_w \sum_n -E(Y^n; X^n) - \log Z(w)\end{aligned}$$

Maximum likelihood

$$L(w) = \sum_i w^T \psi(x_i, y_i) - \log \sum_y \exp^{w^T \psi(x_i, y)}$$

- $L(w)$ is differentiable and convex (it has a positive definite Hessian) so gradient descent can find the global optimum.

Maximum likelihood

$$\begin{aligned}\nabla_w L(w) &= \sum_i \left[\psi(x_i, y_i) - \frac{\sum_y \exp^{w^T \psi(x_i, y)} \psi(x_i, y)}{\sum_y \exp^{w^T \psi(x_i, y)}} \right] \\ &= \sum_i \left[\psi(x_i, y_i) - \sum_{\textcolor{red}{y}} p(x_i, y|w) \psi(x_i, y) \right]\end{aligned}$$

- For general CRFs, there is still a problem with the computation of the derivative because the number of possible configurations for y is typically (exponentially) large.

$$y_i = \{-1, +1\} \rightarrow \left| \sum_y \right| = 2^n$$

Training a structured model ?

- Other solutions exist:
 - Pseudo-likelihood
 - Variational approximation
 - Contrastive divergence
 - **Maximum-margin framework (e.g. Structured SVM)**

3.1. Maximum Margin Training of Structured Models: cutting plane (structured SVM)



Structured SVM

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} V(y_i, y_j)$$

- Given a set of N training examples with ground truth labels $(Y^{(1)}, \dots, Y^{(N)})$, we can write

$$Y^* = \arg \min_{Y \in \mathcal{Y}} E_w(X, Y)$$

\equiv

$$\forall n, Y \in \mathcal{Y}_n \setminus Y^{(n)}$$

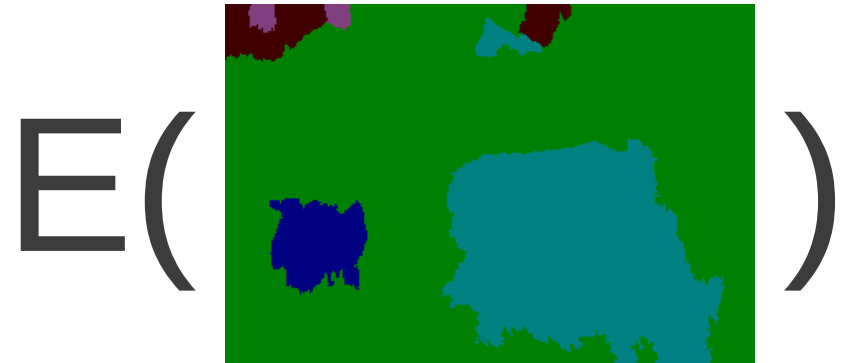
$$E_{\mathbf{w}}(Y^{(n)}) \leq E_{\mathbf{w}}(Y)$$

Energy for the correct labeling at least as low as energy of any incorrect labeling..

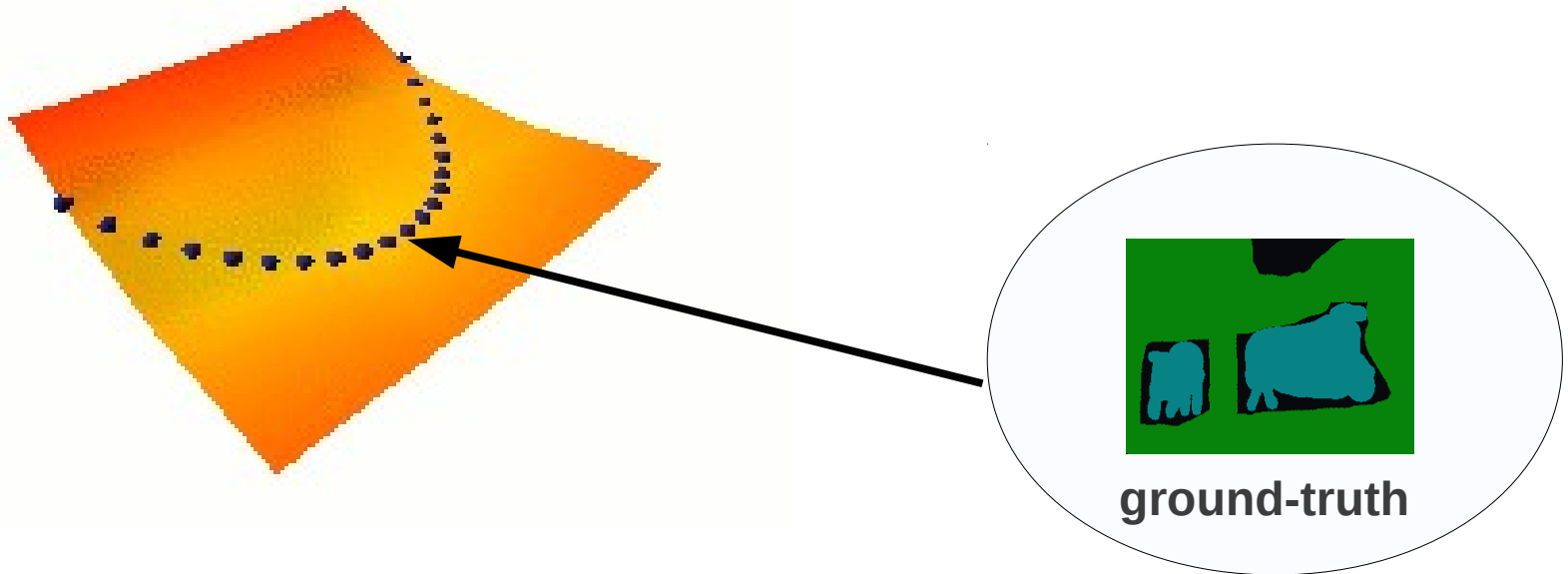
Structured SVM



ground-truth



Energy based landscape



See <http://www.cs.nyu.edu/~yann/research/ebm/>

Structured SVM

- Given a set of N training examples with ground truth labellings $(Y^{(1)}, \dots, Y^{(N)})$ we optimize :

$$\min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{N} \sum_{n=1}^N \xi_n$$

$$\text{s.t. } \forall n, Y \in \mathcal{Y}_n \setminus Y^{(n)} : \delta E_{\mathbf{w}}(Y) \geq \Delta(Y^{(n)}, Y) - \xi_n$$

$$\delta E_{\mathbf{w}}(Y) = E_{\mathbf{w}}(Y) - E_{\mathbf{w}}(Y^{(n)}) \quad \Delta(Y^{(n)}, Y) = \sum_{i \in \mathcal{V}} I(y_i \neq y_i^{(n)})$$

Structured SVM

- Since the SSVM operates by solving a quadratic program (QP), all the constraints must be linear.
- Energy function must be expressible as an inner product between the parameter vector \mathbf{w} and a feature map.

$$E_{\mathbf{w}}(X, Y) = \sum_{i \in \mathcal{V}} D(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} V(y_i, y_j)$$

$$E_{\mathbf{w}}(Y) = \langle \mathbf{w}, \Psi(Y) \rangle .$$

$$\mathbf{w} = ((\mathbf{w}^D)^T, (\mathbf{w}^V)^T)^T$$

Structured SVM

$$\min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{N} \sum_{n=1}^N \xi_n$$

$$\text{s.t. } \forall n, Y \in \mathcal{Y}_n \setminus Y^{(n)} : \delta E_{\mathbf{w}}(Y) \geq \Delta(Y^{(n)}, Y) - \xi_n$$

Exponential number of constraints: $N \times 2^{|\mathcal{V}|}$

Electron microscopy (EM) dataset

See <http://cvlab.epfl.ch/research/medical/em/mitochondria/>

$|\mathcal{V}| = 133555$ nodes



Structured SVM

- In order to deal with the exponential number of constraints in the QP, Tsochantaridis proposed a **cutting plane algorithm**.

- Iteratively finds the most violated constraint and adds it to the working set of constraints.

$$\hat{Y} = \operatorname{argmin}_{Y \in \mathcal{Y}_n} E_{\mathbf{w}}(Y) - \Delta(Y^{(n)}, Y)$$

- See Tsochantaridis, I., Hofmann, T., Joachims, T., Altun, Y.: Support vector machine learning for interdependent and structured output spaces. In: ICML. (2004)



Structured SVM

- In order to deal with the exponential number of constraints in the QP, Tsochantaridis proposed a **cutting plane algorithm**.

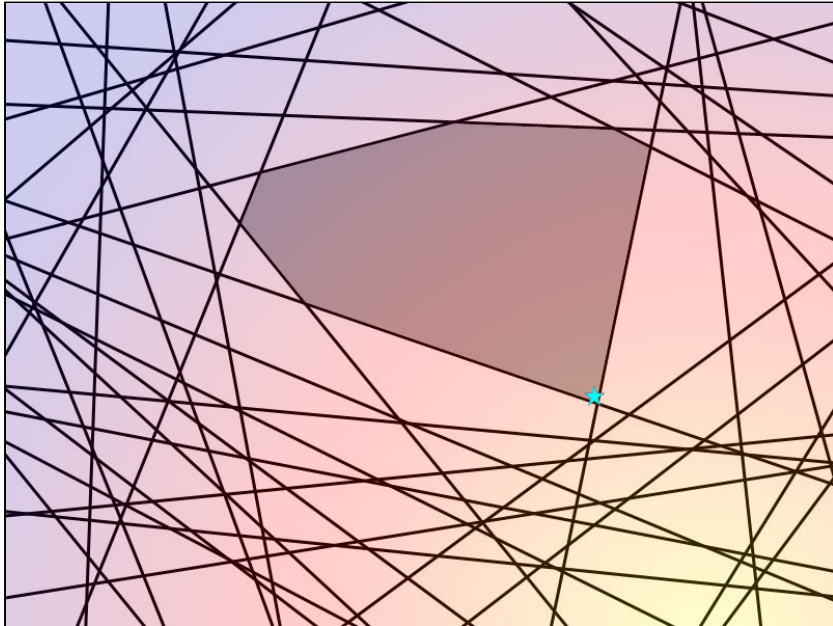
- Iteratively finds the most violated constraint and adds it to the working set of constraints.

$$\hat{Y} = \operatorname{argmin}_{Y \in \mathcal{Y}_n} E_{\mathbf{w}}(Y) - \Delta(Y^{(n)}, Y)$$

- See Tsochantaridis, I., Hofmann, T., Joachims, T., Altun, Y.: Support vector machine learning for interdependent and structured output spaces. In: ICML. (2004)

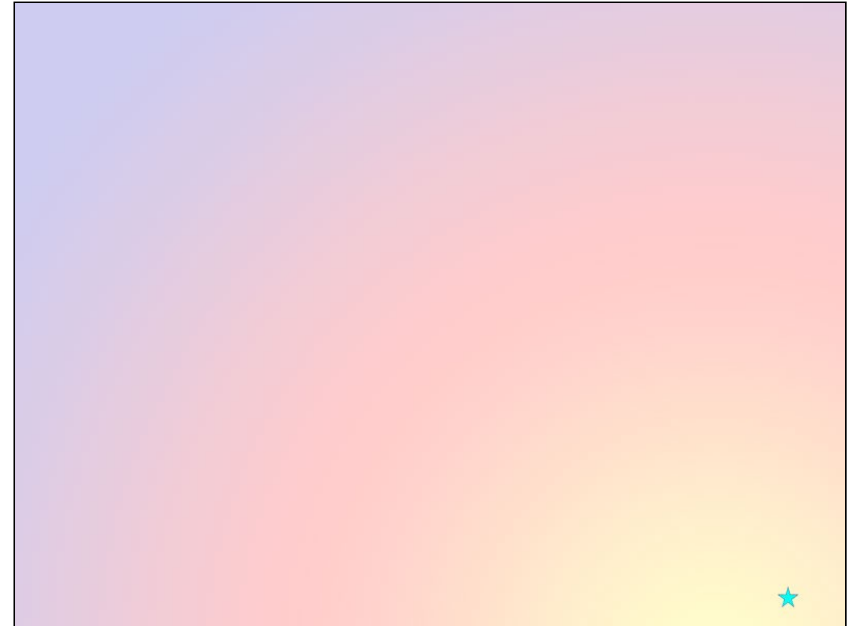


Illustrative Example



SSVM Problem

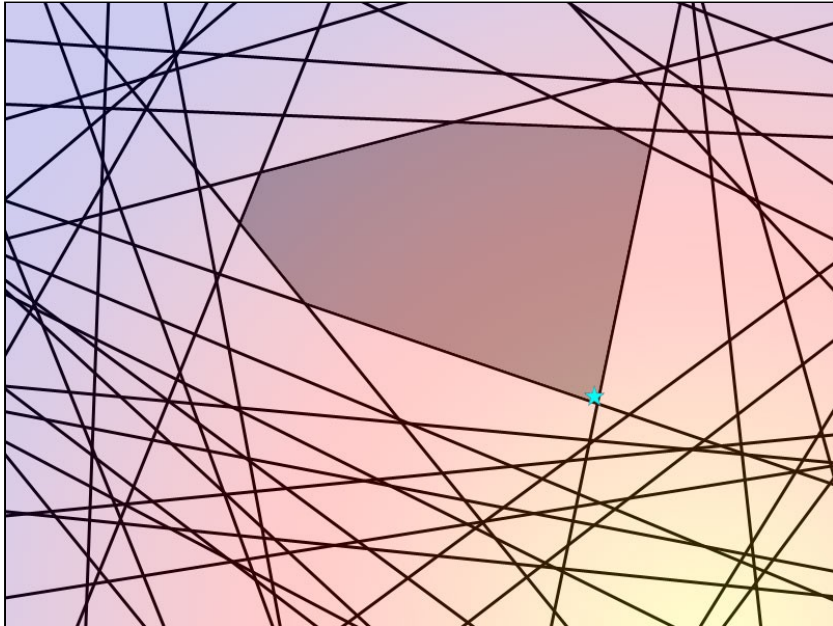
- Exponential constraints
- Most are dominated by a small set of “important” constraints



Cutting plane approach

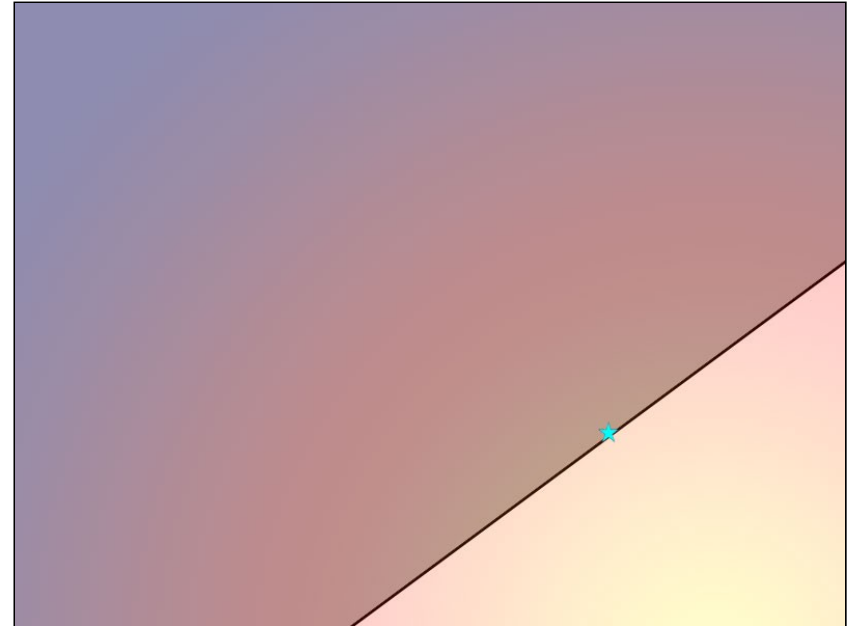
- Repeatedly finds the next most violated constraint...
- ...until set of constraints is a good approximation.

Illustrative Example



SSVM Problem

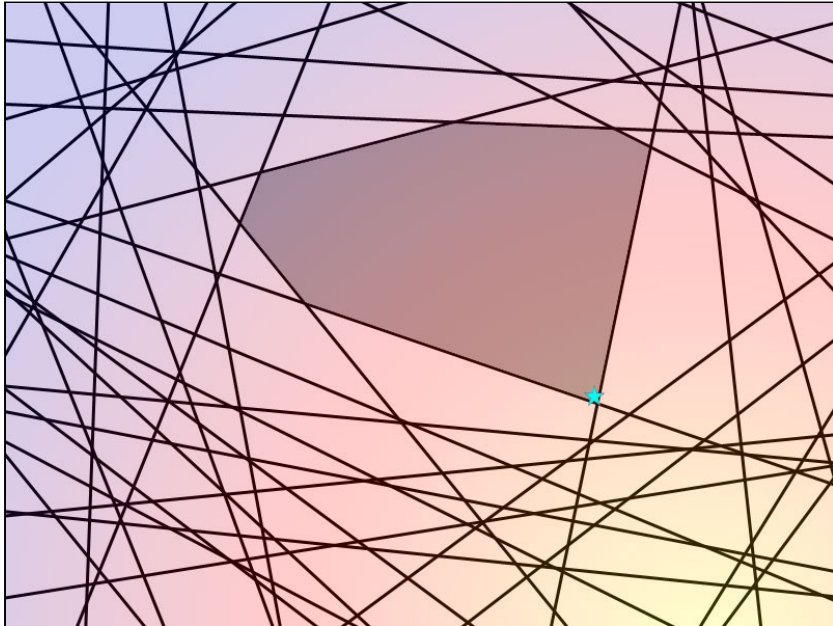
- Exponential constraints
- Most are dominated by a small set of “important” constraints



Cutting plane approach

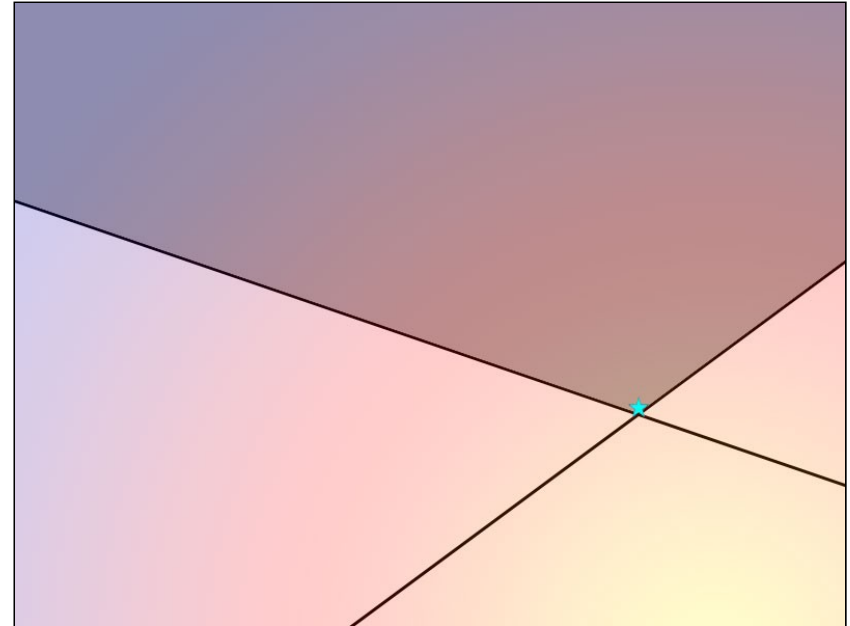
- Repeatedly finds the next most violated constraint...
- ...until set of constraints is a good approximation.

Illustrative Example



SSVM Problem

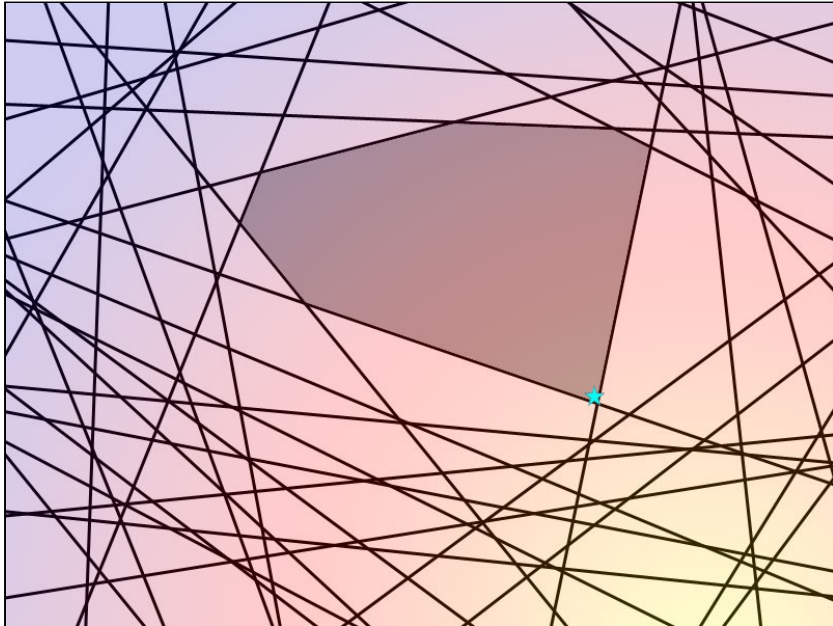
- Exponential constraints
- Most are dominated by a small set of “important” constraints



Cutting plane approach

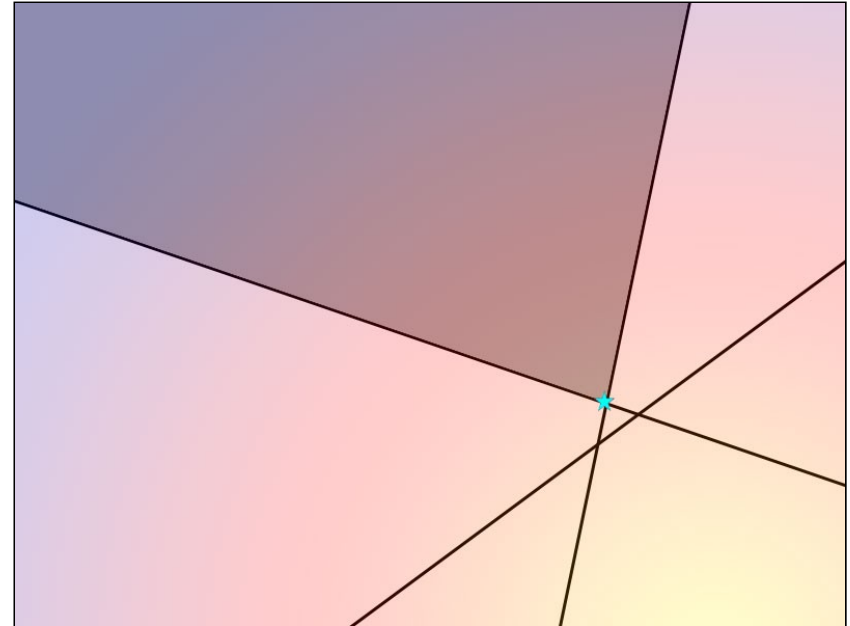
- Repeatedly finds the next most violated constraint...
- ...until set of constraints is a good approximation.

Illustrative Example



SSVM Problem

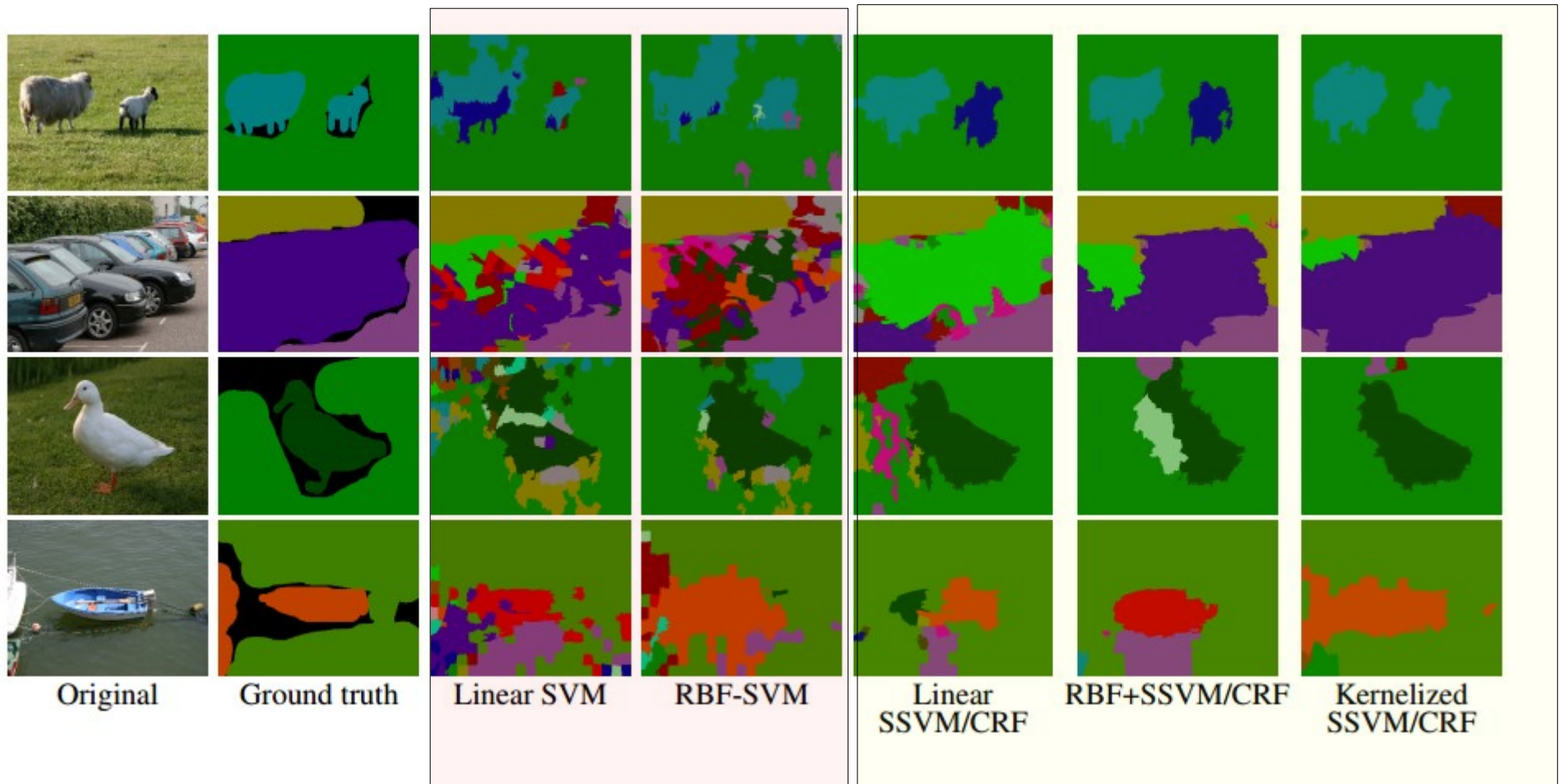
- Exponential constraints
- Most are dominated by a small set of “important” constraints



Cutting plane approach

- Repeatedly finds the next most violated constraint...
- ...until set of constraints is a good approximation.

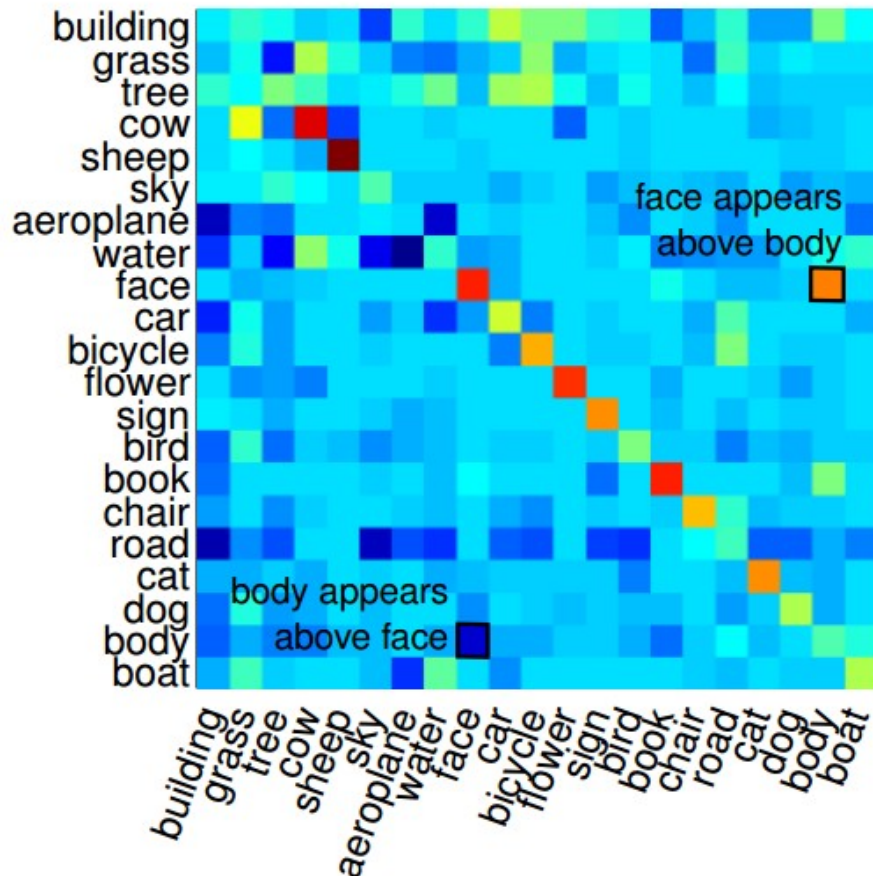
Results



No structure

CRF

Learned pairwise term



- The pairwise term in matrix form where columns indicate classes y_i belonging to superpixel i and rows indicate classes y_j belonging to neighboring j for the MSRC-21 dataset.

W^V

Drawbacks

- Finding the most violated constraint at each iteration of the cutting plane is intractable in loopy graphical models.
- Approximations can sometimes be imprecise enough to have a major impact on learning
 - An unsatisfactory constraint can cause the cutting plane algorithm to prematurely terminate.

3.2. Maximum Margin Training of Structured Models: online subgradient descent (SGD)



SGD approach

- Can reformulate the problem as an unconstrained optimization by plugging the constraints in the objective function.

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{N} \sum_{n=1}^N [E_{\mathbf{w}}(Y^n) + \Delta(Y^n, Y^*) - E_{\mathbf{w}}(Y^*)]_+$$

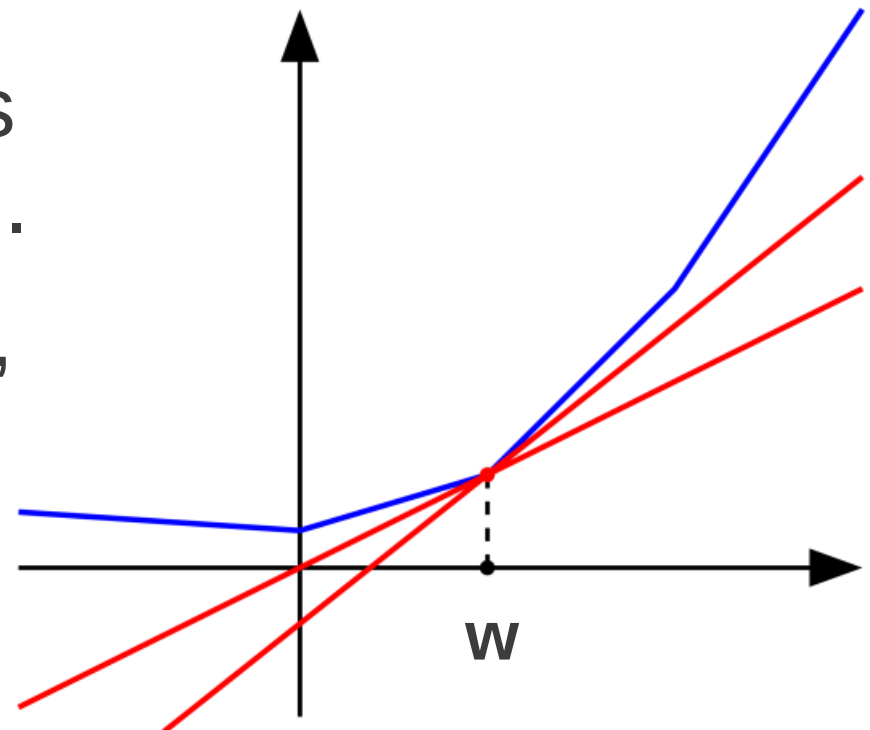
- SGD approach : compute and step in the negative direction of a sub-gradient of \mathcal{L}_w
- See N. Ratliff, J. A. Bagnell, and M. Zinkevich. (Online) Subgradient Methods for Structured Prediction. In AISTATS, 2007.

Subgradient

- A subgradient for the convex loss function \mathbf{L} at \mathbf{w} is defined as a vector \mathbf{g} , such that:

$$\forall \mathbf{w}' \in \mathcal{W}, \mathbf{g}^T (\mathbf{w}' - \mathbf{w}) \leq l(\mathbf{w}') - l(\mathbf{w})$$

- Set of all subgradients is called the subdifferential.
- If \mathbf{L} is differentiable at \mathbf{w} , then \mathbf{g} is the gradient .



SGD approach

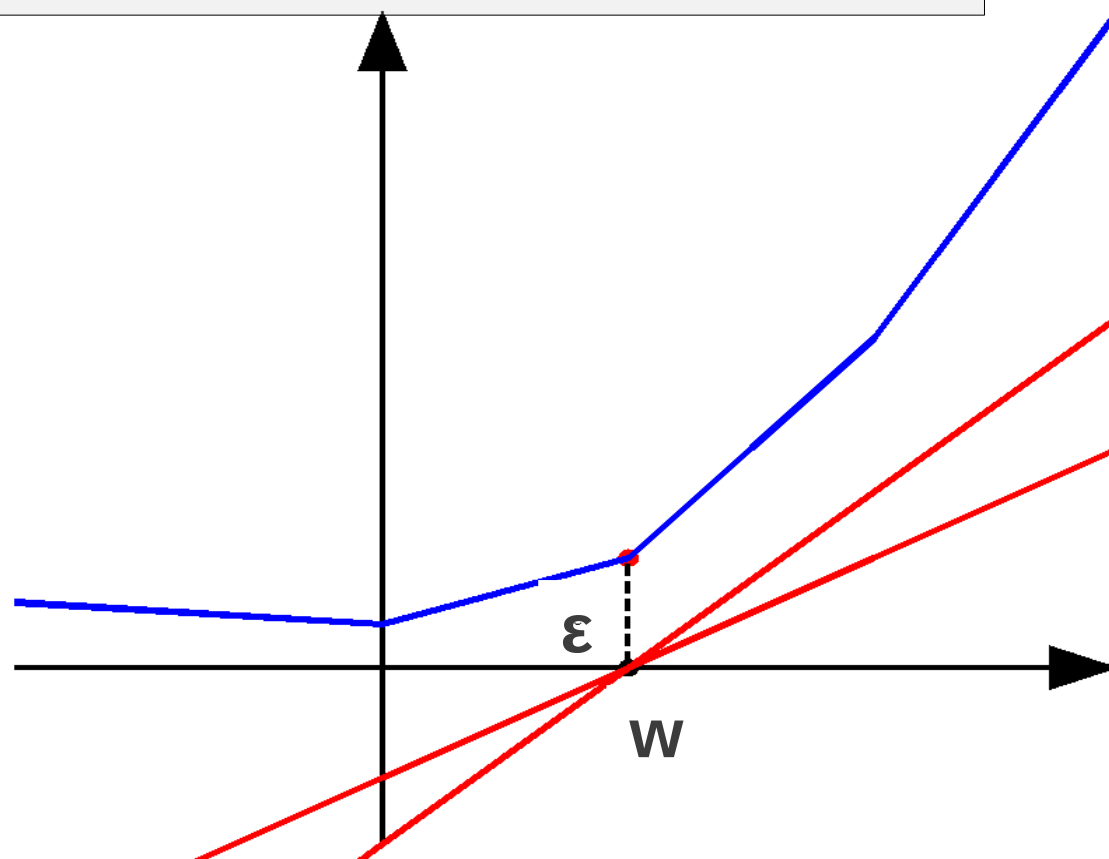
Algorithm 1

```
1: INPUTS :  
2:    $\mathcal{D}$  : training set.  
3:    $\lambda$  : multiplicative factor for step-size.  
4:    $\mathbf{w}_0$  : arbitrary initial values, e.g. 0.  
5: OUTPUT :  $\mathbf{w}_T$   
6: for  $t = 1 \dots T$  do  
7:   for all examples  $(X^n, Y^n)$  in  $\mathcal{D}$  do  
8:      $Y^* = \arg \min_{Y \in \mathcal{Y}_n} (E_{\mathbf{w}}(Y) - \Delta(Y^n, Y))$   
9:      $\eta_t \leftarrow \frac{\lambda}{t}$   
10:     $g_t \leftarrow \frac{\partial l(Y^n, Y, \mathbf{w}_{t-1})}{\partial \mathbf{w}_{t-1}}$   
11:     $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \eta_t g_t$   
12:   end for  
13: end for
```

What can we say about approximate subgradients ?

- Epsilon subgradients:

$$\forall \mathbf{w}' : \mathbf{g}^T (\mathbf{w} - \mathbf{w}') \geq l(\mathbf{w}) - l(\mathbf{w}') - \epsilon$$



What can we say about approximate subgradients ?

- Convergence guarantees (see S. M. Robinson. Linear convergence of epsilon-subgradient descent methods for a class of convex functions. Mathematical Programming, 86:41–50, 1999):

$$\mathcal{L}(\mathbf{w}^{best}) - \mathcal{L}(\mathbf{w}^*) \leq \frac{R^2 + \sum_i (\eta^{(i)})^2 G^2}{2 \sum_i \eta^{(i)}} + \epsilon$$

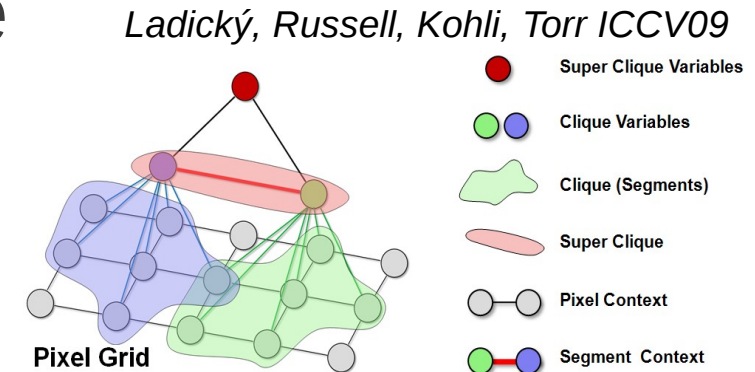
$$\mathbf{w}^{best} = \operatorname{argmin}_{\mathbf{w}^{(t)}} \mathcal{L}(\mathbf{w}^{(t)})$$

$$\left\| \mathbf{w}^{(1)} - \mathbf{w}^* \right\|^2 \leq R^2 \quad \left\| \mathbf{g} \right\|^2 \leq G^2$$

Goes to 0 with appropriate step size

What can we say about approximate subgradients ?

- Models in computer vision are getting more complex.
- Need better approximation algorithms.



$$\mathcal{L}(\mathbf{w}^{best}) - \mathcal{L}(\mathbf{w}^*) \leq \frac{R^2 + \sum_i (\eta^{(i)})^2 G^2}{2 \sum_i \eta^{(i)}} + \epsilon$$

Experimental results

	QP formulation	SGD + argmax	SGD + sampling	SGD + argmax CVPR13
Time for 1000 iterations	19628s	5315s	2481s	5842s
VOC score on EM dataset	80.5 %	79.9 %	77.5 %	84.5 %

Future challenges

- Better energy functions, fully connected CRFs...
- Higher order potentials.
- Better and faster training algorithms.



Increasing connectivity →

Questions



Resources

- <http://cvlab.epfl.ch/research/medical/em/mitochond>
- <http://cvlab.epfl.ch/research/medical/em/synapses/>
- <http://cvlab.epfl.ch/~lucchi/>

Credits

- Slides courtesy :
 - Christoph Lampert (Learning with Structured Inputs and Outputs)
 - Pushmeet Kohli (Efficiently Solving Dynamic Markov Random Fields using Graph Cuts)
 - Ben Taskar (Structured Prediction: A Large Margin Approach, NIPS tutorial)
 - Yisong Yue, Thorsten Joachims (An Introduction to Structured Output Learning Using Support Vector Machines)

