Towards Accurate Structured Output Learning and Prediction



Patrick Pletscher

ETH Zurich

October 25, 2012



🖾 pat@pletscher.org



Patrick Pletscher

Towards Accurate Structured Output Learning and Prediction

Object Recognition in Computer Vision



• Given an image. Task: locate the objects in it.



Patrick Pletscher

Object Recognition in Computer Vision



- Given an image. Task: locate the objects in it.
- Strong dependencies among labels between closeby pixels. "Pixel is likely grass if neighboring pixels are grass"



Object Recognition in Computer Vision



- Given an image. Task: locate the objects in it.
- Strong dependencies among labels between closeby pixels.
 "Pixel is likely grass if neighboring pixels are grass"
- Key idea: Impose some structural constraints. Predict labels of the pixels jointly.



Object Recognition in Computer Vision



- Given an image. Task: locate the objects in it.
- Strong dependencies among labels between closeby pixels.
 "pixel is likely grass if neighboring pixels are grass"
- Key idea: Impose some structural constraints. Predict labels of the pixels jointly.
- This thesis: learning & prediction with structured data.

Additional Examples of Structured Data

• Computer Vision: Image denoising or stereo.



• Natural language processing: Parsing or part-of-speech tagging.



• Biology: Protein side-chain prediction and design.



Structured Output Prediction

Setting

- Given observed input variables $x \in \mathcal{X}$; usually $\mathcal{X} = \mathbb{R}^{D}$.
- Predict a *multivariate* discrete output variable $y \in \mathcal{Y}$.
- Learning: Find good predictor $f_{w}(x) : \mathcal{X} \to \mathcal{Y}$. Parameterized by w.
- Energy *E*(*y*, *x*, *w*)
 - Cost function to score the different outputs.
 - Models the dependencies.

Structured Output Prediction

Setting

- Given observed input variables $x \in \mathcal{X}$; usually $\mathcal{X} = \mathbb{R}^{D}$.
- Predict a *multivariate* discrete output variable $y \in \mathcal{Y}$.
- Learning: Find good predictor $f_{w}(x) : \mathcal{X} \to \mathcal{Y}$. Parameterized by w.
- Energy *E*(*y*, *x*, *w*)
 - Cost function to score the different outputs.
 - Models the dependencies.

Standard Binary Classification

- Binary classification as a special case: $\mathcal{Y} = \{-1, 1\}$.
- Linear prediction function: $f_{\boldsymbol{w}}(\boldsymbol{x}) = \operatorname{sign} \langle \boldsymbol{w}, \boldsymbol{x} \rangle$.

• Energy:
$$E(\mathbf{y}, \mathbf{x}, \mathbf{w}) = -y \langle \mathbf{w}, \mathbf{x} \rangle$$

Structured Models: Why Is It Difficult?

Prediction for an Input x

Choose the best output: $\mathbf{y}^* = f_{\mathbf{w}}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}} E(\mathbf{y}, \mathbf{x}, \mathbf{w}).$

Due to dependencies no obvious way how to do this



Structured Models: Why Is It Difficult?

Prediction for an Input x

Choose the best output: $\mathbf{y}^* = f_{\mathbf{w}}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}} E(\mathbf{y}, \mathbf{x}, \mathbf{w}).$

Due to dependencies no obvious way how to do this

Exhaustive Enumeration? Here: Object Recognition

- *M* pixels, *K* different object classes. $|\mathcal{Y}| = K^M$ possibilities.
- Usually: K > 10, $M \gg 10 \times 10$.
- Compare to: 10^{80} atoms in the universe.
- Prediction as a computational problem.

Need for Clever Algorithms and Approximations

- Some problems exactly tractable. But often only approximate.
- Learning the energy: Even harder as prediction a subprocedure.

SOP Overview – Energy

• "Blueprint of a model" $\mathcal{G}=(\mathcal{V},\mathcal{E}).$



- Linear dependence on parameters: $E(m{y},m{x},m{w}) = -\langlem{w},\phi(m{x},m{y})
 angle$
- Explicit way to write the energy:

$$\mathsf{E}(oldsymbol{y}) = \sum_{i \in \mathcal{V}} heta_i(y_i) + \sum_{(i,j) \in \mathcal{E}} heta_{ij}(y_i, y_j)$$

construct θ from w and x.





SOP Overview – Loss

• Loss: the smaller the better!



• Example: Missclassified pixels.

$$\Delta_{\boldsymbol{y}^{\star}}(\boldsymbol{y}) = \sum_{i \in \mathcal{V}} y_i \neq y_i^{\star}$$

- More complex losses (wrong direction):
 - Overlap of bounding boxes
 - F₁ score
 - Area-under-curve



Towards Accurate Structured Output Learning and Prediction

SOP Overview – Learning

- Goal: Learn a good energy.
 Ground-truth has small energy
- Learning = Estimation of w^* .

$$\min_{\boldsymbol{w}} \underbrace{\frac{\lambda}{2} \|\boldsymbol{w}\|^2}_{\text{regularizer}} + \frac{1}{N} \sum_{n=1}^{N} \underbrace{\ell(\boldsymbol{w}, \boldsymbol{x}^n, \boldsymbol{y}^n)}_{\text{surrogate loss}}$$

I. Max-margin loss for (x, y^*) :

$$E(\mathbf{y}^{\star}, \mathbf{x}, \mathbf{w}) - \min_{\mathbf{y} \in \mathcal{Y}} [E(\mathbf{y}, \mathbf{x}, \mathbf{w}) - \Delta_{\mathbf{y}^{\star}}(\mathbf{y})]$$

2. Log-loss for (x, y^*) :

$$E(\mathbf{y}^{\star}, \mathbf{x}, \mathbf{w}) + \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(-E(\mathbf{y}, \mathbf{x}, \mathbf{w}))$$



SOP Overview – Prediction

- Given input x and weights w^* .
- Construct potentials: w^* , $x \mapsto \theta$. •
- Minimize energy.

$$\min_{\mathbf{y}\in\mathcal{V}}\sum_{i\in\mathcal{V}}\theta_i(y_i) + \sum_{(i,j)\in\mathcal{E}}\theta_{ij}(y_i, y_j)$$



SOP Overview – Evaluation

Compare prediction y to ground-truth y^* using the loss function.



SOP Overview – Tractability

Exact algorithms for loopy graphs only if

I. $E(\mathbf{y})$ submodular and binary:

 $heta_{ij}(0,0) + heta_{ij}(1,1) \leq heta_{ij}(0,1) + heta_{ij}(1,0)$

- 2. Loss function "easy".
- 3. Max-margin for learning.



My PhD Thesis

Learning:

- Pletscher & Kohli, AISTATS 2012.
- Pletscher & Ong, AISTATS 2012.
- Lacoste-Julien, Jaggi, Schmidt & Pletscher, 2012.
- Pletscher, Ong & Buhmann, ECML 2010.
- Pletscher, Ong & Buhmann, AISTATS 2009.

Prediction:

• Pletscher & Wulff, ICML 2012.

Evaluation:

Pletscher, Nowozin, Kohli & Rother, DAGM 2011.



LPQP for Energy Minimization Pletscher & Wulff, ICML 2012

Novel Algorithm for Approximate Energy Minimization

• Compute minimum energy assignment for general pairwise energies:

$$\min_{\mathbf{y}} E(\mathbf{y}) = \min_{\mathbf{y}} \sum_{i \in \mathcal{V}} \theta_i(y_i) + \sum_{(i,j) \in \mathcal{E}} \theta_{ij}(y_i, y_j).$$

- Combination of two relaxations:
 - Linear Programming (Schlesinger 1976).
 - Quadratic Programming (Ravikumar and Lafferty 2006).
- Efficient Message Passing Algorithms.



Linear and Quadratic Programming Relaxations

LP for Marginal Polytope

$$\min_{oldsymbol{\mu}\in\mathcal{M}} \quad \sum_{i\in\mathcal{V}}oldsymbol{ heta}_i^{\mathsf{T}}oldsymbol{\mu}_i + \sum_{(i,j)\in\mathcal{E}}oldsymbol{ heta}_{ij}^{\mathsf{T}}oldsymbol{\mu}_{ij}$$

But: \mathcal{M} is exponentially large!





Linear and Quadratic Programming Relaxations

LP for Marginal Polytope

$$\min_{\boldsymbol{\mu}\in\mathcal{M}} \quad \sum_{i\in\mathcal{V}} \boldsymbol{\theta}_i^\mathsf{T} \boldsymbol{\mu}_i + \sum_{(i,j)\in\mathcal{E}} \boldsymbol{\theta}_{ij}^\mathsf{T} \boldsymbol{\mu}_{ij}$$

But: \mathcal{M} is exponentially large!

LP for Local Marginal Polytope

$$\min_{\boldsymbol{\mu} \in \mathcal{L}_{\mathcal{G}}} \sum_{i \in \mathcal{V}} \boldsymbol{\theta}_i^\mathsf{T} \boldsymbol{\mu}_i + \sum_{(i,j) \in \mathcal{E}} \boldsymbol{\theta}_{ij}^\mathsf{T} \boldsymbol{\mu}_{ij}$$





Patrick Pletscher

Towards Accurate Structured Output Learning and Prediction

Linear and Quadratic Programming Relaxations



LP for Local Marginal Polytope

$$\min_{\boldsymbol{\mu} \in \mathcal{L}_{\mathcal{G}}} \sum_{i \in \mathcal{V}} \boldsymbol{\theta}_{i}^{\mathsf{T}} \boldsymbol{\mu}_{i} + \sum_{(i,j) \in \mathcal{E}} \boldsymbol{\theta}_{ij}^{\mathsf{T}} \boldsymbol{\mu}_{ij}$$

$$\begin{split} \min_{\boldsymbol{\mu} \in \mathcal{L}_{\mathcal{G}}} & \sum_{i \in \mathcal{V}} \boldsymbol{\theta}_{i}^{\mathsf{T}} \boldsymbol{\mu}_{i} + \sum_{(i,j) \in \mathcal{E}} \boldsymbol{\theta}_{ij}^{\mathsf{T}} \boldsymbol{\mu}_{ij} \\ \text{s.t.} & \boldsymbol{\mu}_{ij} = \mathsf{vec}(\boldsymbol{\mu}_{i} \boldsymbol{\mu}_{i}^{\mathsf{T}}) \quad \forall (i,j) \in \mathcal{E} \end{split}$$

LPQP: Combine LP and QP relaxations

Joint LP and QP Objective

$$\min_{\mu \in \mathcal{L}_{\mathcal{G}}} \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{\mu} + \underbrace{\rho \sum_{(i,j) \in \mathcal{E}} D_{\mathcal{KL}}(\mu_{ij}, \mu_i \mu_j^{\mathsf{T}})}_{\text{encourages consistency}}.$$

Numerical Solution

- Non-convex KL divergence: use the Concave-Convex Procedure.
- Iteratively solve convex optimization problems.
- Efficient message-passing algorithms.
- Gradual increase of ρ .

Decision Tree Fields

LPQP Results in Low Energy Solutions



Towards Accurate Structured Output Learning and Prediction

High-order Loss Functions

Pletscher & Kohli, AISTATS 2012

Exact Max-margin Learning for High-order Losses

• High-order loss $\Delta_{y^*}(y)$: does not factorize into unaries



Maximum margin ⇔ Energy Minimization:

 $\ell_{mm}(\boldsymbol{w}, \boldsymbol{x}, \boldsymbol{y}^{\star}) = E(\boldsymbol{y}^{\star}, \boldsymbol{x}, \boldsymbol{w}) - \min_{\boldsymbol{y} \in \mathcal{Y}} [E(\boldsymbol{y}, \boldsymbol{x}, \boldsymbol{w}) - \Delta_{\boldsymbol{y}^{\star}}(\boldsymbol{y})]$

• We characterize a family of tractable high-order loss functions.



Label-count Loss: Introduce Auxiliary Variable

Label-count loss

Pairwise graphical model



Patrick Pletscher

Towards Accurate Structured Output Learning and Prediction



Patrick Pletscher

Towards Accurate Structured Output Learning and Prediction

Summary & Future Work

Summary

- Work in most aspects of structured output prediction and learning.
- Goal: fast and accurate approximations for training and prediction.

Future Work

- Smooth-max for direct loss minimization.
- Max-margin learning on a relaxed polytope. Use similar ideas as in LPQP for enforcing consistency.
- Applications. Currently working with ImageNet, a 2 TB dataset.



Thanks To My Collaborators





thanks also to: Committee, ETHZ ML Groups (JB+AK), Friends/Family



Towards Accurate Structured Output Learning and Prediction

15/16



References I

- Lacoste-Julien, Simon et al. (2012). "Stochastic Block-Coordinate Frank-Wolfe Optimization for Structural SVMs".
- Pletscher, Patrick and Pushmeet Kohli (2012). "Learning low-order models for enforcing high-order statistics". In: Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS). JMLR: W&CP 22, pp. 886–894.
- Pletscher, Patrick and Cheng Soon Ong (2012). "Part & Clamp: An efficient algorithm for structured output learning". In: Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS). JMLR: W&CP 22, pp. 877–885.
- Pletscher, Patrick, Cheng Soon Ong, and Joachim M. Buhmann (2009). "Spanning Tree Approximations for Conditional Random Fields". In: Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics (AISTATS). JMLR: W&CP 5, pp. 408–415.

References II

Pletscher, Patrick, Cheng Soon Ong, and Joachim M. Buhmann (2010).
"Entropy and Margin Maximization for Structured Output Learning".
In: Proceedings of the 20th European Conference on Machine Learning (ECML). Vol. 6321. Lecture Notes in Computer Science, pp. 83–98.
Pletscher, Patrick and Sharon Wulff (2012). "LPQP for MAP: Putting LP Solvers to Better Use". In: Proceedings of the 29th International Conference on Machine Learning (ICML).

- Pletscher, Patrick et al. (2011). "Putting MAP back on the Map". In: 33rd Annual Symposium of the German Association for Pattern Recognition (DAGM). Vol. 6835. Lecture Notes in Computer Science. Springer, pp. 111–121.
- Ravikumar, Pradeep and John Lafferty (2006). "Quadratic Programming Relaxations for Metric Labeling and Markov Random Field MAP Estimation". In: Proceedings of the 23rd international conference on Machine learning (ICML), pp. 737–744.

References III

Schlesinger, M I (1976). "Syntactic analysis of two-dimensional visual signals in noisy conditions". In: *Kibernetika* 4, pp. 113–130.