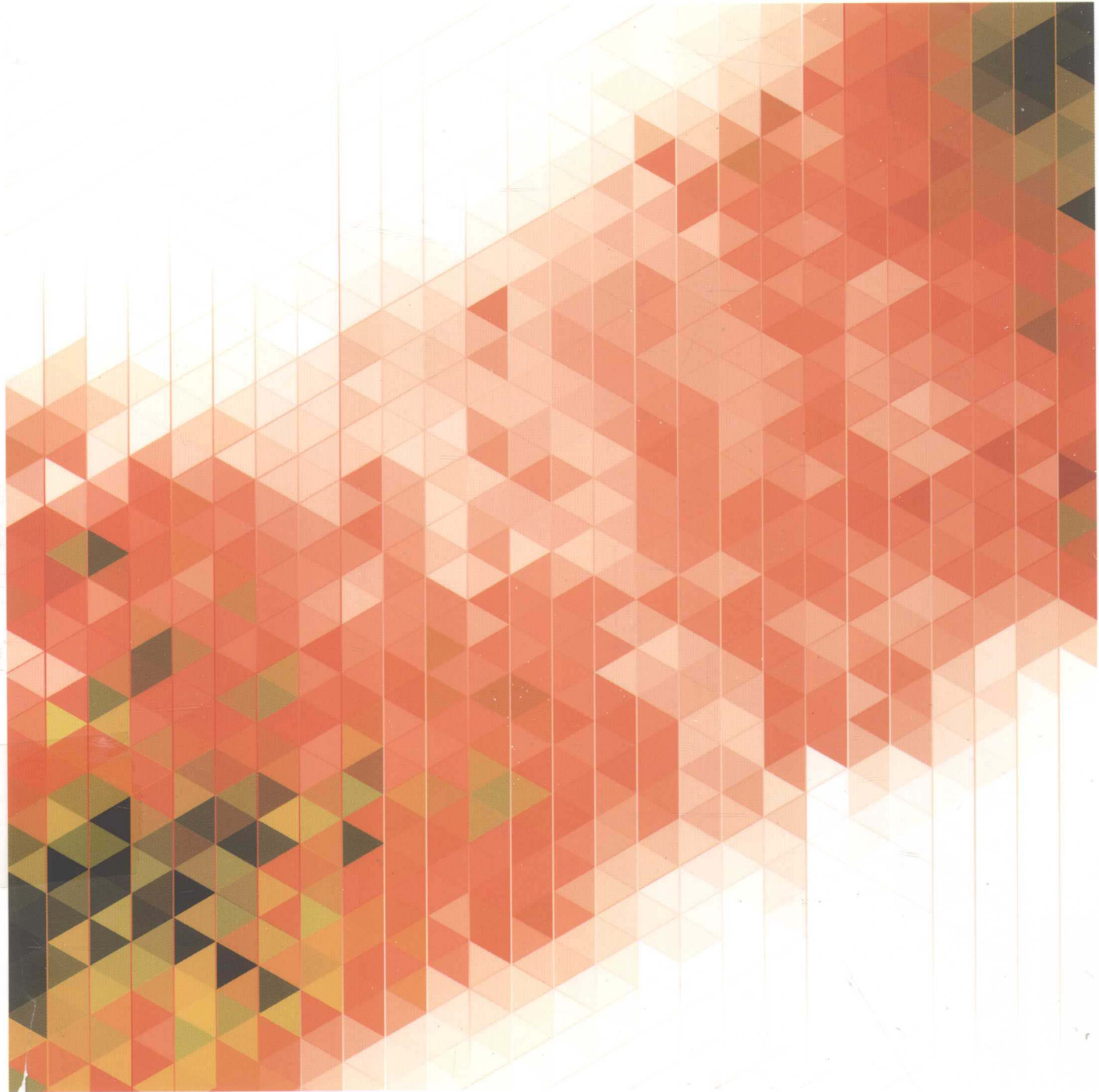


Advanced Structured Prediction

EDITED BY

Sebastian Nowozin, Peter V. Gehler,
Jeremy Jancsary, and Christoph H. Lampert



Advanced Structured Prediction

Edited by

Sebastian Nowozin, Peter V. Gehler, Jeremy Jancsary, and
Christoph H. Lampert

The MIT Press
Cambridge, Massachusetts
London, England

© 2014 Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

MIT Press books may be purchased at special quantity discounts for business or sales promotional use. For information, please email special_sales@mitpress.mit.edu

This book was set in L^AT_EX by the authors and editors. Printed and bound in the United States of America.

Library of Congress Cataloging-in-Publication Data

Advanced structured prediction / edited by Sebastian Nowozin, Peter V. Gehler, Jeremy Jancsary, and Christoph H. Lampert.

pages cm. - (Neural information processing series)

Includes bibliographical references and index.

ISBN 978-0-262-02837-0 (hardcover : alk. paper) 1. Machine learning. 2. Computer algorithms. 3. Data structures (Computer science) I. Nowozin, Sebastian, 1980– editor.

Q325.5.A295 2014

006.3'1-dc23

2014013235

10 9 8 7 6 5 4 3 2 1

Series Foreword

The yearly Neural Information Processing Systems (NIPS) workshops bring together scientists with broadly varying backgrounds in statistics, mathematics, computer science, physics, electrical engineering, neuroscience, and cognitive science, unified by a common desire to develop novel computational and statistical strategies for information processing and to understand the mechanisms for information processing in the brain. In contrast to conferences, these workshops maintain a flexible format that both allows and encourages the presentation and discussion of work in progress. They thus serve as an incubator for the development of important new ideas in this rapidly evolving field. The series editors, in consultation with workshop organizers and members of the NIPS Foundation Board, select specific workshop topics on the basis of scientific excellence, intellectual breadth, and technical impact. Collections of papers chosen and edited by the organizers of specific workshops are built around pedagogical introductory chapters, while research monographs provide comprehensive descriptions of workshop-related topics, to create a series of books that provides a timely, authoritative account of the latest developments in the exciting field of neural computation.

Michael I. Jordan and Thomas G. Dietterich

Preface

Machine learning is one of the fastest growing areas of computer science, and with good reason: predictive machine learning models trained on ever growing data sets provide relevant information to scientists and business decision makers alike, as well as enabling intelligent consumer applications.

Structured prediction refers to machine learning models that predict relational information that has structure such as being composed of multiple interrelated parts. For example, these models are used to predict a natural language sentence or segment an image into meaningful components. Structured prediction models are important in many application domains and have been used with great success in biology, computer vision, and natural language processing.

This volume is not the first on the topic of structured prediction; seven years ago, in 2007, MIT Press released the edited volume *Predicting Structured Data*. Since then structured prediction has blossomed into many application areas, but it has not settled down yet; there continues to be a stream of interesting and original work. In an introduction chapter, we summarize the state-of-the-art and recent developments. The remainder of the volume is a careful selection of contributed chapters.

We would like to thank all chapter contributors for their high-quality work, Marie Lufkin Lee from MIT Press for her support and patience, Suvrit Sra for help in preparing a L^AT_EX template for this volume, and Jasmin Pielorz for help with proofreading and copy-editing.

We dedicate this volume to the memory of Ben Taskar, a pioneer of the field.

Sebastian Nowozin, Peter V. Gehler,
Jeremy Jancsary, Christoph H. Lampert

Cambridge, Tübingen, Vienna, Klosterneuburg
January 2014

Contents

Series Foreword	xi
Preface	xiii
1 Introduction to Structured Prediction	
<i>S. Nowozin, P.V. Gehler, J. Jancsary, and C. H. Lampert</i>	1
1.1 Structured Prediction	2
1.2 Recent Developments	7
1.3 Summary of the Chapters	10
1.4 Conclusion	14
1.5 References	15
2 The Power of LP Relaxation for MAP Inference	
<i>S. Živný, T. Werner, and D. Průša</i>	19
2.1 Valued Constraint Satisfaction Problem	20
2.2 Basic LP Relaxation	22
2.3 Languages Solved by the BLP	22
2.4 Universality of the BLP	31
2.5 Conclusions	36
2.6 References	38
3 AD³: A Fast Decoder for Structured Prediction	
<i>A. Martins</i>	43
3.1 Introduction	43
3.2 Factor Graphs and MAP Decoding	45
3.3 Factor Graphs for NLP	48
3.4 LP-MAP Decoding	52
3.5 Alternating Directions Dual Decomposition (AD ³)	55
3.6 Local Subproblems in AD ³	61
3.7 Experiments	65
3.8 Related Work	68
3.9 Conclusions	69

3.10	References	70
4	Generalized Sequential Tree-Reweighted Message Passing	
	<i>T. Schoenemann and V. Kolmogorov</i>	75
4.1	Introduction	75
4.2	Background and Notation	77
4.3	TRW-S Algorithm	80
4.4	Algorithm's Analysis	86
4.5	Experimental Results	89
4.6	Conclusions	91
4.7	Appendices	94
4.8	References	100
5	Smoothed Coordinate Descent for MAP Inference	
	<i>O. Meshi, T. Jaakkola, and A. Globerson</i>	103
5.1	Introduction	104
5.2	MAP and LP Relaxations	106
5.3	Coordinate Minimization Algorithms	111
5.4	Dual Convergence Rate Analysis	113
5.5	Primal Convergence	116
5.6	The Augmented Dual LP Algorithm	118
5.7	Experiments	120
5.8	Discussion	123
5.9	Appendix: Primal Convergence Rate	125
5.10	References	129
6	Getting Feasible Variable Estimates from Infeasible Ones: MRF Local Polytope Study	
	<i>B. Savchynskyy and S. Schmidt</i>	133
6.1	Introduction	133
6.2	Optimizing Projection	138
6.3	MRF Inference and Optimizing Projections	141
6.4	Optimizing Projection in Algorithmic Schemes	146
6.5	Experimental Analysis and Evaluation	150
6.6	Conclusions	153
6.7	References	155
7	Perturb-and-MAP Random Fields: Reducing Random Sampling to Optimization, with Applications in Computer Vision	
	<i>G. Papandreou and A. Yuille</i>	159
7.1	Introduction	160

7.2	Energy-Based Modeling: Standard Deterministic and Probabilistic Approaches	162
7.3	Perturb-and-MAP for Gaussian and Sparse Continuous MRFs	166
7.4	Perturb-and-MAP for MRFs with Discrete Labels	171
7.5	Related Work and Recent Developments	180
7.6	Discussion	182
7.7	References	182
8	Herding for Structured Prediction	
	<i>Y. Chen, A.E. Gelfand, and M. Welling</i>	187
8.1	Introduction	188
8.2	Integrating Local Models Using Herding	194
8.3	Application: Image Segmentation	199
8.4	Application: Go Game Prediction	205
8.5	Conclusion	209
8.6	References	211
9	Training Structured Predictors Through Iterated Logistic Regression	
	<i>J. Domke</i>	213
9.1	Introduction	213
9.2	Linear vs. Nonlinear Learning	215
9.3	Overview	216
9.4	Loss Functions	216
9.5	Message-Passing Inference	220
9.6	Joint Learning and Inference	222
9.7	Logistic Regression	223
9.8	Reducing Structured Learning to Logistic Regression	224
9.9	Function Classes	225
9.10	Example	230
9.11	Conclusions	231
9.12	Appendix: Proofs	232
9.13	References	237
10	PAC-Bayesian Risk Bounds and Learning Algorithms for the Regression Approach to Structured Output Prediction	
	<i>S. Giguère, F. Laviolette, M. Marchand, and A. Rolland</i>	239
10.1	Introduction	240
10.2	From Structured Output Prediction to Vector-Valued Regression	241
10.3	A PAC-Bayesian Bound with Isotropic Gaussians	244

10.4	A Sample Compressed PAC-Bayesian Bound	247
10.5	Empirical Results	252
10.6	Conclusion	255
10.7	Appendix	255
10.8	References	264
11	Optimizing the Measure of Performance	
	<i>J. Keshet</i>	267
11.1	Introduction	267
11.2	Structured Perceptron	269
11.3	Large Margin Structured Predictors	270
11.4	Conditional Random Fields	273
11.5	Direct Loss Minimization	274
11.6	Structured Ramp Loss	275
11.7	Structured Probit Loss	276
11.8	Risk Minimization Under Gibbs Distribution	278
11.9	Conclusions	279
11.10	References	279
12	Structured Learning from Cheap Data	
	<i>X. Lou, M. Kloft, G. Rätsch and F. A. Hamprecht</i>	281
12.1	Introduction	282
12.2	Running Example: Structured Learning for Cell Tracking . .	284
12.3	Strategy I: Structured Learning from Partial Annotations . .	287
12.4	Strategy II: Structured Data Retrieval via Active Learning .	294
12.5	Strategy III: Structured Transfer Learning	299
12.6	Discussion and Conclusions	303
12.7	References	303
13	Dynamic Structured Model Selection	
	<i>D. Weiss and B. Taskar</i>	307
13.1	Introduction	307
13.2	Meta-Learning a Myopic Value-Based Selector	310
13.3	Applications to Sequential Prediction	312
13.4	Meta-Learning a Feature Extraction Policy	318
13.5	Applications to Sequential Prediction Revisited	325
13.6	Conclusion	329
13.7	References	331
14	Structured Prediction for Event Detection	
	<i>M. Hoai and F. de la Torre</i>	333
14.1	Introduction	333

14.2	Structured Prediction for Event Detection	336
14.3	Early Event Detection	339
14.4	Sequence Labeling	345
14.5	Experiments	349
14.6	Summary	357
14.7	References	358
15	Structured Prediction for Object Boundary Detection in Images	
	<i>S. Todorovic</i>	363
15.1	Introduction	363
15.2	Related Work	365
15.3	Edge Extraction and Properties	367
15.4	Sequential Labeling of Edges	369
15.5	HC-Search	372
15.6	Results	375
15.7	Conclusion	385
15.8	References	386
16	Genome Annotation with Structured Output Learning	
	<i>J. Behr, G. Schweikert and G. Rätsch</i>	389
16.1	Introduction: The Genome Annotation Problem	390
16.2	Inference	397
16.3	Learning	400
16.4	Experiments	406
16.5	Conclusions	410
16.6	References	412

1 Introduction to Structured Prediction

Sebastian Nowozin

*Microsoft Research
Cambridge, United Kingdom*

Sebastian.Nowozin@Microsoft.com

Peter V. Gehler

*Max Planck Insitute for Intelligent Systems
72076 Tübingen, Germany*

Peter.Gehler@tuebingen.mpg.de

Jeremy Jancsary

*Nuance Communications
Vienna, Austria*

Jeremy.Jancsary@Nuance.com

Christoph H. Lampert

*IST Austria
A-3400 Klosterneuburg, Austria*

chl@ist.ac.at

Structured prediction refers to machine learning models that predict multiple interrelated and dependent quantities. These models are commonly used in computer vision, speech recognition, natural language processing, and computational biology to accurately reflect prior knowledge, task-specific relations, and constraints. A wide variety of types of models is used, and they are expressive and powerful, but exact computation in these models is often intractable. This difficulty, paired with the practical significance, has resulted in a broad research effort in recent years to design structured prediction models and approximate inference and learning procedures that are computationally efficient. This chapter gives an introduction to structured prediction and summarizes the main approaches. It includes a discussion of the research trends in the field since 2007 and provides further references for the interested reader.

1.1 Structured Prediction

The general structured prediction problem is defined as follows. Given an observation $x \in \mathcal{X}$, make a prediction $y \in \mathcal{Y}(x)$ as

$$y = f(x). \tag{1.1}$$

The set $\mathcal{Y}(x)$ is typically finite but exponentially large, and its size may depend on the input x . A popular choice is to use an index set $I = \{1, 2, \dots, m\}$ and define both input x and prediction y as

$$x = (x_1, \dots, x_m), \quad \text{and} \quad y = (y_1, \dots, y_m).$$

For example, I can index all words in a sentence or all pixels in an image.

Researchers working in structured prediction are concerned with the *representation* of the function f , procedures for *evaluating* $f(x)$ for a given input x , and *learning* f from a class of functions \mathcal{F} given annotated training data consisting of pairs (x, y) of data instances.

We describe these three aspects below, but first we would like to define how to measure the quality of a structured prediction model (1.1) by means of loss functions.

Loss Functions and Decision Rules. A generally accepted criterion for assessing the quality of our model is that of the *expected loss* of our model as a function of the true generating probability distribution $q(x, y)$ and a *loss function*¹ $\ell : \mathcal{Y}(x) \times \mathcal{Y}(x) \rightarrow \mathbb{R}$. The distribution $q(x, y)$ is the sampling distribution we encounter when our model (1.1) is used; for example, it could be the joint distribution of emails x and spam/no-spam decisions y that are sent to a particular email address. We do not know q , but a standard assumption is that we are able to obtain independent and identically distributed (iid) samples from it. The loss function $\ell(z, y)$ quantifies—on an arbitrary but fixed scale—the loss suffered if z happens to be the truth and we decide for y . The quality of a structured prediction model can now be quantified as the *risk*,

$$\mathcal{R}(f, q, \ell) = \mathbb{E}_{(x, y) \sim q} [\ell(y, f(x))]. \tag{1.2}$$

1. Alternatively, an equivalent definition can be made using *utility functions*; we want to maximize utility or minimize loss, but except for a change in sign, both definitions are identical. The loss function can be more generally defined as $\ell : \mathcal{Y} \times \mathcal{D} \rightarrow \mathbb{R}$, where \mathcal{D} is the *decision domain*, which can differ from \mathcal{Y} .

Because \mathcal{R} depends on the unknown distribution q , the expectation is approximated using a data set $D = \{(x^{(i)}, y^{(i)})\}_{i=1, \dots, N}$ sampled iid from q , yielding the *empirical risk*,

$$\mathcal{R}_{\text{emp}}(f, D, \ell) = \frac{1}{N} \sum_{i=1}^N \ell(y^{(i)}, f(x^{(i)})). \quad (1.3)$$

While there are different philosophies with respect to how to best build structured prediction models and which loss functions are relevant to an application, the criterion (1.3) is widely accepted.

The best possible risk, which is the lowest possible, is known as the *Bayes risk*. It is defined by making the optimal decisions with the knowledge of q , that is, $\mathcal{R}_{\text{Bayes}}(q, \ell) = \mathcal{R}(f_{\text{Bayes}}, q, \ell)$, where f_{Bayes} is the Bayes-optimal predictor,

$$f_{\text{Bayes}}(x) = \underset{y \in \mathcal{Y}(x)}{\operatorname{argmin}} \mathbb{E}_{z \sim q(z|x)} [\ell(z, y)]. \quad (1.4)$$

Representation. For representing the function f , different choices exist; one popular branch of the literature defines $f(x)$ as the maximizer of an auxiliary optimization problem,

$$f(x) = \underset{y \in \mathcal{Y}(x)}{\operatorname{argmax}} F(x, y, \theta), \quad (1.5)$$

where $\theta \in \Theta$ are model parameters. In many applications, solving (1.5) corresponds to solving a combinatorial optimization problem. The function $F(x, y, \theta)$ to be maximized is commonly parametrized as a linear form,

$$F(x, y, \theta) = \langle \phi(x, y), \theta \rangle, \quad (1.6)$$

where $\Theta = \mathbb{R}^d$ and $\phi(x, y)$ is a *joint feature map*, transforming x and y into a large but fixed size feature vector. The class of functions is now indexed by θ , and we have

$$\mathcal{F} = \{F(\cdot, \cdot, \theta) \mid \theta \in \mathbb{R}^d\}. \quad (1.7)$$

Another approach to construct structured prediction functions is by starting with a probabilistic model and applying Bayesian decision theory (Berger, 1985). For this we assume that we have a model for the conditional distribution $p(y|x; \theta)$ over $\mathcal{Y}(x)$. Together with a loss function ℓ , we can then use the *Bayes decision rule*,

$$f(x) = \underset{y \in \mathcal{Y}(x)}{\operatorname{argmin}} \mathbb{E}_{z \sim p(z|x)} [\ell(z, y)]. \quad (1.8)$$

This rule is identical to (1.4), except we replaced the unknown distribution q with our model p . Intuitively (1.8) selects our prediction so that we minimize our expected loss under every possibility z , weighted by our beliefs about the state of the world as encoded in $p(z|x)$. The similarity between (1.8) and (1.4) implies that if p equals the true distribution q , then our decisions made using the Bayes decision rule will be optimal, that is, they will achieve the *Bayes risk*.

Evaluation. Using either definition (1.5) or (1.8), in order to make predictions, we need to solve an optimization problem, an instance of an *inference problem*. Depending on the structure of F and $\mathcal{Y}(x)$, problem (1.5) may be intractable to solve exactly, and we need to develop *approximate inference methods*. When using (1.5), such methods are often called *energy minimization methods*, and a large part of the structured prediction literature is concerned with their properties. In case (1.8) is used, the tractability depends on the distribution p , the loss function ℓ , and the set $\mathcal{Y}(x)$. For example, if the so-called 0/1-loss $\ell_{0/1}(z, y) = 1_{\{y \neq z\}}$ is used, the problem reduces to the *maximum-a-posteriori* (MAP) decision rule,

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}(x)} p(y|x). \quad (1.9)$$

If the loss function decomposes additively over individual dimensions of its arguments, then we can solve (1.8) in two steps, where first a set of low-dimensional marginal distributions $p(y_i|x)$ is inferred, and then decisions are independently made by minimizing $\mathbb{E}_{z_i \sim p(z_i|x)} [\ell_i(z_i, y_i)]$ (Marroquin et al., 1987). Inferring the marginal distributions $p(y_i|x)$, also known as *marginal beliefs*, requires probabilistic inference methods for the model. In the last fifteen years, a large number of approximate inference methods have been developed to this end.

One important class of methods, the linear programming relaxations, apply to discrete graphical models (Wainwright and Jordan, 2008). For these models (1.5) can be reformulated as an integer linear program, which can be relaxed to a polynomial-time solvable linear program for which specialized message-passing algorithms have been developed. These algorithms are now popularly used and provide robust inference for otherwise challenging models, but until recently, understanding the structure and limitations of the linear programming relaxation approach has been an open question.

Learning. Structured prediction models can be learned in different ways from a given data set of iid samples. If the direct form of the predictor (1.5) is adopted, then the most popular choice is *regularized risk minimization*

(Vapnik and Chervonenkis, 1974), in which we minimize the regularized empirical risk,

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \Omega(f) + \frac{1}{N} \sum_{i=1}^N \ell(y^{(i)}, f(x^{(i)})). \quad (1.10)$$

Here $\Omega(f)$ is a regularizer that controls the capacity of the learned model \hat{f} . While (1.10) has served as motivation for a large number of general machine learning methods, the application to structured prediction problems was only enabled through the work of Tsochantaridis et al. (2004), who showed how (1.10) can be implemented in the structured case when the linear form (1.6) for the definition of f is used.

They propose the *structured support vector machine*, which learns the parameters θ of the predictor f by solving the problem

$$\operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{2} \|\theta\|^2 + \frac{\lambda}{N} \sum_{i=1}^N L(y^{(i)}, x^{(i)}, \theta), \quad (1.11)$$

where $\lambda > 0$ is a regularization parameter, and we define

$$L(y^{(i)}, x^{(i)}, \theta) = \max_{y \in \mathcal{Y}(x^{(i)})} \left[\ell(y^{(i)}, y) - F(x^{(i)}, y^{(i)}, \theta) + F(x^{(i)}, y, \theta) \right]. \quad (1.12)$$

It can be shown that $L(y^{(i)}, x^{(i)}, \theta) \geq \ell(y^{(i)}, f(x^{(i)}))$, that is, L is an upper bound of ℓ for any θ . Therefore, (1.11) is an upper bound of the empirical risk (1.3), and by minimization of the upper bound, we can find model parameters with low empirical risk. Given enough training data, the empirical risk will be close to the true risk (1.2). While it is not trivial to solve (1.11), it is a convex optimization problem, and the tractability of the formulation has enabled a large number of structured prediction applications.

When the probabilistic perspective is adopted, learning of the model is performed using the model likelihood, using either maximum likelihood estimation (MLE) or Bayesian inference (Koller and Friedman, 2009). The model can either be generative $p(x, y|\theta)$ or discriminative $p(y|x, \theta)$ as in conditional random fields (Lafferty et al., 2001). The generative model provides an explicit model for the inputs x , whereas the discriminative model always conditions on an observed x . For the following example, let us use a discriminative model. We specify a *prior distribution* $p(\theta)$ for the model parameters and then solve

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} p(\theta) \prod_{i=1}^N p(y^{(i)}|x^{(i)}, \theta), \quad (1.13)$$

for the maximum likelihood estimate $\hat{\theta}$ or use Bayes rule to define a posterior belief over θ given the data set D as

$$p(\theta|D) = \frac{p(\theta) p(D|\theta)}{p(D)} \quad (1.14)$$

$$= \frac{p(\theta) \prod_{i=1}^N p(y^{(i)}|x^{(i)}, \theta)}{\int_{\Theta} p(\theta) \prod_{i=1}^N p(y^{(i)}|x^{(i)}, \theta) d\theta} \quad (1.15)$$

$$\propto p(\theta) \prod_{i=1}^N p(y^{(i)}|x^{(i)}, \theta). \quad (1.16)$$

At test time, given a new observation x , we proceed as follows. In case the MLE is used, the predictive distribution is derived from the point estimate $\hat{\theta}$ simply as $p(y|x, \hat{\theta})$. In case the posterior $p(\theta|D)$ is used, the predictive distribution marginalizes over all parameter uncertainty as

$$p(y|x) = \int_{\Theta} p(y|x, \theta) p(\theta|D) d\theta. \quad (1.17)$$

Learning structured prediction models is challenging because it is usually intractable to perform exact computation of the required quantities, with few exceptions—for example, in so-called linear chain models. In both (1.13) and (1.16), we use $p(y^{(i)}|x^{(i)}, \theta)$, but this important distribution typically cannot be exactly computed. Likelihood-based learning of structured models has therefore required approximations; a large variety of approximate inference and estimation methods have been proposed. A rough grouping of these methods is into stochastic and deterministic approximations.

Stochastic approximations perform Monte Carlo simulations to approximate expectations and integrals in the learning objective. Examples are MCMC-MLE (Descombes et al., 1999) and stochastic maximum likelihood approaches such as contrastive divergence (Hinton, 2002; Carreira-Perpiñán and Hinton, 2005).

Deterministic approximations typically optimize an auxiliary objective function; the class of *variational approximations* such as *mean field* methods (Saul and Jordan, 1995; Xing et al., 2003), *loopy belief propagation* (Yedidia et al., 2004), or more generally methods derived from the minimization of statistical divergence measures (Minka, 2005) are commonly used for otherwise intractable models. Another class of deterministic approximations instead modify the likelihood function itself to obtain tractable estimators; these include the *pseudolikelihood* (Besag, 1972, 1977), more general *composite likelihoods* (Lindsay, 1988; Varin et al., 2010), and *score matching* (Hyvärinen, 2005).

All these approximations have different trade-offs with respect to computational effort, the robustness and accuracy of the inference results, as well as the theory that is known about them. It is fair to say that while great progress has been made, for most models, there is not yet a clear favorite among the above approximate methods.

1.2 Recent Developments

We now briefly summarize the most significant developments in the field of structured prediction since around 2007, when the previous volume in this series was published (Bakir et al., 2007).

Joint Optimization over Parameters and Inference. Meshi et al. (2010) introduced a clever method to approximately solve (1.10) for discrete graphical models. In their method they rewrite ℓ as a maximization problem over vectors μ in the *local polytope* so that (1.10) is of a min-max structure with multiple inner maximization problems. The inner problems corresponding to ℓ are then dualized using convex duality to obtain a joint min-min problem over parameters θ as well as dual message vectors. The advantage of this is that one can now interleave message passing updates and parameter updates, whereas previously every parameter update required a message passing scheme to run to convergence. The result is an efficient learning method. A similar but more general method has been proposed by Hazan and Urtasun (2010). In this volume, Chapter 9 by Justin Domke continues this line of work to learn more expressive non-linear model potentials in which the parameter update step is replaced by a non-linear logistic regression subproblem.

Integrated Estimation and Inference. A recent take on how to deal with the intractability of structured prediction models is due to Domke (2011), Ross et al. (2011b), and Stoyanov et al. (2011). The idea is to take a model and an iterative approximate inference procedure, such as loopy belief propagation, and to view them as one computational unit that iteratively transforms some initial state into an inference result. As such, when using a fixed number of inference iterations, it is just a non-linear differentiable mapping from parameters and observations to inference results. This mapping is parametrized by the original model parameters, and as long as we can compute the gradient with respect to these parameters, a gradient-based optimizer can be used to minimize the empirical risk (Domke, 2013; Ross et al., 2011b). The combination of model and approximate