

Statistical and Computational Tradeoffs in
Stochastic Composite Likelihood

by

J. Dillon and G. Lebanon
Purdue University

Technical Report #08-04

Department of Statistics
Purdue University

September 2008

Statistical and Computational Tradeoffs in Stochastic Composite Likelihood

Joshua Dillon

School of Electrical and Computer Engineering
Purdue University

Guy Lebanon

Department of Statistics
Purdue University

June 4, 2008

Abstract

Maximum likelihood estimators are often of limited practical use due to the intensive computation they require. We propose a family of alternative estimators that maximize a stochastic variation of the composite likelihood function. We prove the consistency of the estimators, provide formulas for their asymptotic variance and computational complexity, and discuss experimental results in the context of Boltzmann machines and conditional random fields. The theoretical and experimental studies demonstrate the effectiveness of the estimators in achieving a predefined balance between computational complexity and statistical accuracy.

1 Introduction

Maximum likelihood estimation is by far the most popular point estimation technique in machine learning and statistics. Assuming that the data consists of n , m -dimensional vectors

$$D = \{X^{(1)}, \dots, X^{(n)}\} \subset \mathbb{R}^m, \quad (1)$$

and is sampled iid from a parametric distribution p_{θ_0} with $\theta_0 \in \Theta \subset \mathbb{R}^r$, a maximum likelihood estimator (mle) $\hat{\theta}_n^{\text{mle}}$ is a maximizer of the loglikelihood function

$$\ell_n(\theta; D) = \sum_{i=1}^n \log p_{\theta}(X^{(i)}). \quad (2)$$

The use of the mle is motivated by its consistency, i.e. $\hat{\theta}_n^{\text{mle}} \rightarrow \theta_0$ as $n \rightarrow \infty$ with probability 1 [4]. The consistency property ensures that as the number n of data samples grows, the estimator will converge to the true parameter value θ_0 governing the data generation process.

An even stronger motivation for the use of the mle is that it has an asymptotically normal distribution with mean vector θ_0 and variance matrix $(nI(\theta_0))^{-1}$. More formally, we have the following convergence in distribution as $n \rightarrow \infty$ [4]

$$\sqrt{n}(\hat{\theta}_n^{\text{mle}} - \theta_0) \rightsquigarrow N(0, I^{-1}(\theta_0)), \quad (3)$$

where $I(\theta)$ is the $r \times r$ Fisher information matrix

$$I(\theta) = \mathbb{E}_{p_{\theta}} \{ \nabla \log p_{\theta}(X) (\nabla \log p_{\theta}(X))^{\top} \} \quad (4)$$

with ∇f represents the $r \times 1$ gradient vector of $f(\theta)$ with respect to θ . The convergence (3) is especially striking since according to the Cramer-Rao lower bound, the asymptotic variance $(nI(\theta_0))^{-1}$ of the mle is the smallest possible variance for any estimator. Since it achieves the lowest possible asymptotic variance, the mle (and other estimators which share this property) is said to be asymptotically efficient.

The consistency and asymptotic efficiency of the mle motivate its use in many circumstances. Unfortunately, in some situations the maximization or even evaluation of the loglikelihood (2) and its derivatives is impossible due to computational considerations. This has led to the proposal of alternative estimators under the premise that a loss of asymptotic efficiency is acceptable—in return for reduced computational complexity. Consistency however, is typically viewed as less negotiable and inconsistent estimators should be avoided if at all possible.

In this paper, we propose a family of estimators, for use in situations where the computation of the mle is intractable. In contrast to previously proposed approximate estimators, our estimators are statistically consistent and admit a precise quantification of both computational complexity and statistical accuracy through their asymptotic variance. Due to the continuous parameterization of the estimator family, we obtain an effective framework for achieving a predefined problem-specific balance between computational tractability and statistical accuracy. For the sake of concreteness, we focus on the case of estimating the parameters associated with Markov random fields. In this case, we provide a detailed discussion of the accuracy complexity tradeoff and experimental results for the Boltzmann machine and conditional random fields.

2 Related Work

There is a large body of work dedicated to tractable learning techniques. Two popular categories are Markov chain Monte Carlo (MCMC) and variational methods. MCMC is a general purpose technique for approximating expectations and can be used to approximate the normalization term and other intractable portions of the loglikelihood and its gradient [3]. Variational methods are techniques for conducting inference and learning based on tractable bounds. Despite the substantial work on MCMC and variational methods, there are few results that are general enough to be practical while preserving clear results concerning convergence and approximation rate.

Our work draws on Lindsay’s composite likelihood method [7] of parameter estimation which in turn generalized Besag’s pseudo likelihood [2]. A selection of more recent studies on pseudo and composite likelihood are [1, 6, 9, 8, 5]. Most of the recent studies in this area examine the behavior of the pseudo or composite likelihood in a particular modeling situation. We believe that the present paper is the first to systematically examine statistical and computational tradeoffs in a general quantitative framework. Possible exceptions are [11] which is a mostly experimental study in context of MRFs for texture generation and [10] which is focused on inference rather than parameter estimation.

3 Stochastic Composite Likelihood

In many cases, the absence of a closed form expression for the normalization term prevents the computation of the loglikelihood (2) and its derivatives thereby severely limiting the use of the mle. A popular example are Markov random fields, wherein the computation of the normalization term is often intractable (see Section 5 for more details). In this paper we propose alternative estimators based on the maximization of a stochastic variation of the composite likelihood function.

We start by defining Besag’s pseudo loglikelihood function [2] associated with the data D (1)

$$p\ell_n(\theta; D) = \sum_{i=1}^n \sum_{j=1}^m \log p_\theta(X_j^{(i)} | \{X_k^{(i)} : k \neq j\}). \quad (5)$$

The maximum pseudo likelihood estimator (mple) $\hat{\theta}_n^{\text{mple}}$ is consistent, but possesses considerably higher asymptotic variance than that of the mle $(nI(\theta_0))^{-1}$. Its main advantage is that it does not require the computation of the normalization term as it cancels out in the probability ratio defining conditional distributions, viz.

$$p_\theta(X_j | \{X_k : k \neq j\}) = p_\theta(X) / \sum_{X_j'} p_\theta(X_1, \dots, X_{j-1}, X_j', X_{j+1}, \dots, X_m). \quad (6)$$

The mle and mple represent two different ways of resolving the tradeoff between asymptotic variance and computational complexity. The mle has low asymptotic variance but high computational complexity while the mple has higher asymptotic variance but low computational complexity. It is desirable to obtain additional estimators realizing alternative resolutions of the accuracy complexity tradeoff. To this end we define the stochastic composite likelihood function whose maximization provides a family of consistent estimators with statistical accuracy and computational complexity spanning the entire accuracy-complexity spectrum.

Stochastic composite likelihood generalizes the likelihood and pseudo likelihood functions by constructing an objective function that is a stochastic sum of likelihood objects. We start by defining the notion of m -pairs and likelihood objects and then proceed to stochastic composite likelihood.

Definition 1. An m -pair (A, B) is a pair of sets $A, B \subset \{1, \dots, n\}$ satisfying $A \neq \emptyset = A \cap B$. The likelihood object associated with an m -pair (A, B) and X is $S_\theta(A, B) = \log p_\theta(X_A | X_B)$ where $X_S \stackrel{\text{def}}{=} \{X_j : j \in S\}$. We similarly define likelihood objects with respect to a dataset $D = \{X^{(1)}, \dots, X^{(n)}\}$ as

$$S_\theta(n, A, B) = \sum_{i=1}^n \log p_\theta(X_A^{(i)} | X_B^{(i)}).$$

The composite loglikelihood function, proposed by Lindsay [7], is a collection of likelihood objects defined by a finite sequence of m -pairs $(A_1, B_1), \dots, (A_k, B_k)$

$$cl_n(\theta; D) = \sum_{j=1}^k S_\theta(n, A_j, B_j) = \sum_{i=1}^n \sum_{j=1}^k \log p_\theta(X_{A_j}^{(i)} | X_{B_j}^{(i)}). \quad (7)$$

There exists a certain lack of flexibility associated with the composite likelihood framework. Since each likelihood object $S_\theta(n, A, B)$ is either selected or not, there is no allowance for some objects to be selected more frequently than others. Allowing stochastic, rather than deterministic, selection of likelihood objects provides a higher degree of flexibility and a richer parametric family of estimators. Furthermore, the discrete parameterization of (7) defined by the sequence $(A_1, B_1), \dots, (A_k, B_k)$ is less convenient for theoretical analysis than the continuous parameterization underlying the stochastic variation of composite likelihood defined below.

Definition 2. The stochastic composite loglikelihood (scl) associated with a finite sequence of m -pairs $(A_1, B_1), \dots, (A_k, B_k)$ is

$$scl_n(\theta; D) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k \beta_j Z_{ij} \log p_\theta(X_{A_j}^{(i)} | X_{B_j}^{(i)}). \quad (8)$$

where $\beta_j > 0$ and $Z_{ij} \sim \text{Ber}(\lambda_j)$ are independent binary Bernoulli rv with parameters $\lambda_j \in [0, 1]$.

In other words, the scl is a stochastic version of (7) where for each sample $X^{(i)}, i = 1, \dots, n$, the likelihood objects $S(A_1, B_1), \dots, S(A_k, B_k)$ are selected independently with probabilities $\lambda_1, \dots, \lambda_k$. The positive weights β_j provide additional flexibility by emphasizing different components more than others.

In analogy to the mle and the mple, the maximum scl estimator (mscle) $\hat{\theta}_n^{\text{mscl}}$ estimates θ_0 by maximizing the scl function. In contrast to the loglikelihood and pseudo loglikelihood functions, the scl function and its maximizer are random variables that depend on the indicator variables Z_{ij} in addition to D . As such, its behavior should be summarized by examining its expectation or its behavior in the limit $n \rightarrow \infty$. Different selections of the continuous parameters $(\lambda, \beta) \in [0, 1]^k \times \mathbb{R}_+^k$ underlying the scl function result in different asymptotic variance and computational complexity. As a result the accuracy and complexity of $\hat{\theta}_n^{\text{mscl}}$ become continuous functions over the parametric space $[0, 1]^k \times \mathbb{R}_+^k$ which include as special cases the mle, mple, and maximum quasi likelihood [5] estimators. Different selections of $(\lambda, \beta) \in [0, 1]^k \times \mathbb{R}_+^k$ represent estimators $\hat{\theta}_n^{\text{mscl}}$ achieving different resolutions of the accuracy-complexity tradeoff—each appropriate in a different situation.

4 Statistical Properties of $\hat{\theta}_n^{\text{mscl}}$

The statistical properties of the mscl estimator depend on the selection probabilities and positive weights $(\lambda, \beta) \in [0, 1]^k \times \mathbb{R}_+^k$ while the computational properties depend only on λ . Under some mild conditions $\hat{\theta}_n^{\text{mscl}}$ may be shown to be a consistent estimator whose asymptotic distribution is Gaussian with a certain variance matrix that is larger or equal to the optimal variance expressed by the inverse Fisher information. For simplicity, we assume that the random vector X is discrete and $p_\theta(x)$ is a probability mass function, rather than a density.

Definition 3. A sequence of m -pairs $(A_1, B_1), \dots, (A_k, B_k)$ ensures identifiability of p_θ if the map $\{p_\theta(X_{A_j}|X_{B_j}) : j = 1, \dots, k\} \mapsto p_\theta(X)$ is injective. In other words, there exists only a single collection of conditionals $\{p_\theta(X_{A_j}|X_{B_j}) : j = 1, \dots, k\}$ that does not contradict the joint $p_\theta(X)$.

Proposition 1. Let $(A_1, B_1), \dots, (A_k, B_k)$ be a sequence of m -pairs that ensures identifiability of $p_\theta, \theta \in \Theta$ and $\alpha_1, \dots, \alpha_k$ positive constants. Then

$$\sum_{j=1}^k \alpha_k D(p_\theta(X_{A_j}|X_{B_j}) || p_{\theta'}(X_{A_j}|X_{B_j})) \geq 0$$

where equality holds iff $\theta = \theta'$.

Proof. The inequality follows from applying Jensen's inequality for each conditional KL divergence

$$-D(p_\theta(X_{A_j}|X_{B_j}) || p_{\theta'}(X_{A_j}|X_{B_j})) = E_{p_\theta} \log \frac{p_{\theta'}(X_{A_j}|X_{B_j})}{p_\theta(X_{A_j}|X_{B_j})} \leq \log E_{p_\theta} \frac{p_{\theta'}(X_{A_j}|X_{B_j})}{p_\theta(X_{A_j}|X_{B_j})} = \log 1 = 0.$$

For equality to hold we need each term to be 0 which follows only if $p_\theta(X_{A_j}|X_{B_j}) \equiv p_{\theta'}(X_{A_j}|X_{B_j})$ for all j which, assuming identifiability, holds iff $\theta = \theta'$. \square

Proposition 2. Let $\lambda \in [0, 1]^k$ and $(A_1, B_1), \dots, (A_k, B_k)$ be a sequence of m -pairs for which $\{(A_j, B_j) : \forall j \text{ such that } \lambda_j > 0\}$ ensures identifiability. We also assume that $\Theta \subset \mathbb{R}^r$ is an open set and p_θ is continuous and smooth in θ . Then there exists a strongly consistent sequence of scl maximizers, i.e. $\hat{\theta}_n^{\text{msl}} \rightarrow \theta_0$ as $n \rightarrow \infty$ with probability 1.

Proof. The scl function, modified slightly by multiplication and addition with constants in θ is

$$scl'(\theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k \beta_j \left(Z_{ij} \log p_\theta(X_{A_j}^{(i)}|X_{B_j}^{(i)}) - \lambda_j \log p_{\theta_0}(X_{A_j}^{(i)}|X_{B_j}^{(i)}) \right).$$

By the strong law of large numbers, the above expression converges to its expectation

$$\mu(\theta) = - \sum_{j=1}^k \beta_j \lambda_j D(p_\theta(X_{A_j}|X_{B_j}) || p_{\theta_0}(X_{A_j}|X_{B_j})).$$

Due to Proposition 1 the above function is non-positive and is zero iff $\theta = \theta_0$. As a result, if we restrict ourselves to the compact set $S = \{\theta : c_1 \leq \|\theta - \theta_0\| \leq c_2\}$ the function $\mu(\theta)$ would attain its maximum δ on S which would be strictly negative. This means that there exists N such that for all $n > N$ the scl maximizers on S would achieve strictly negative values of $scl'(\theta)$ with probability 1. However, since $scl'(\theta)$ can be made to achieve arbitrarily close to zero values under $\theta = \theta_0$ we have that $\hat{\theta}_n^{\text{msl}} \notin S$ for $n > N$. Since c_1, c_2 were chosen arbitrarily $\hat{\theta}_n^{\text{msl}} \rightarrow \theta_0$ with probability 1. \square

For example, m -pair sequences containing the pseudo likelihood sequence $A_i = \{i\}, B_i = \{1, \dots, m\} \setminus A_i, i = 1, \dots, k$ as a subsequence ensure identifiability and consequently the consistency of the mscl estimator.

Proposition 3. Assuming the assumptions of Proposition 2 as well as convexity of $\Theta \subset \mathbb{R}^r$ we have

$$\sqrt{n}(\hat{\theta}_n^{\text{msl}} - \theta_0) \rightsquigarrow N(0, \Upsilon \Sigma \Upsilon) \tag{9}$$

where $\Upsilon^{-1} = \sum_{j=1}^k \beta_j \lambda_j \text{Var}_{\theta_0}(V_j), V_j = \nabla S_{\theta_0}(A_j, B_j), \Sigma = \text{Var}_{\theta_0}(\sum_{j=1}^k \beta_j \lambda_j V_j)$.

The notation $\text{Var}_{\theta_0}(Y)$ represents the covariance matrix of the random vector Y under p_{θ_0} while the notations $\xrightarrow{p}, \rightsquigarrow$ in the proof below denote convergences in probability and in distribution [4].

Proof. By the mean value theorem and convexity of Θ there exists $\eta \in (0, 1)$ for which $\theta' = \theta_0 + \eta(\hat{\theta}_n^{\text{msl}} - \theta_0)$ and

$$\nabla \text{scl}_n(\hat{\theta}_n^{\text{msl}}) = \nabla \text{scl}_n(\theta_0) + \nabla^2 \text{scl}_n(\theta')(\hat{\theta}_n^{\text{msl}} - \theta_0)$$

where $\nabla f(\theta)$ and $\nabla^2 f(\theta)$ are the $r \times 1$ gradient vector and $r \times r$ matrix of second order derivatives of $f(\theta)$. Since $\hat{\theta}_n$ maximizes the scl, $\nabla \text{scl}_n(\hat{\theta}_n^{\text{msl}}) = 0$ and

$$\sqrt{n}(\hat{\theta}_n^{\text{msl}} - \theta_0) = -\sqrt{n}(\nabla^2 \text{scl}_n(\theta'))^{-1} \nabla \text{scl}_n(\theta_0). \quad (10)$$

By Proposition 2 we have $\hat{\theta}_n^{\text{msl}} \xrightarrow{p} \theta_0$ which implies that $\theta' \xrightarrow{p} \theta_0$ as well. Furthermore, by the law of large numbers and the fact that if $W_n \xrightarrow{p} W$ then $g(W_n) \xrightarrow{p} g(W)$ for continuous g ,

$$\begin{aligned} (\nabla^2 \text{scl}_n(\theta'))^{-1} &\xrightarrow{p} (\nabla^2 \text{scl}_n(\theta_0))^{-1} \xrightarrow{p} \left(\sum_{j=1}^k \beta_j \lambda_j \mathbb{E}_{\theta_0} \nabla^2 S_{\theta_0}(A_j, B_j) \right)^{-1} \\ &= - \left(\sum_{j=1}^k \beta_j \lambda_j \text{Var}_{\theta_0}(\nabla S_{\theta_0}(A_j, B_j)) \right)^{-1}. \end{aligned} \quad (11)$$

For the remaining term in (10) we have

$$\sqrt{n} \nabla \text{scl}_n(\theta_0) = \sum_{j=1}^k \beta_j \sqrt{n} \frac{1}{n} \sum_{i=1}^n W_{ij}$$

where the random vectors $W_{ij} = Z_{ij} \nabla \log p_{\theta}(X_{A_j}^{(i)} | X_{B_j}^{(i)})$ have expectation 0 and variance matrix $\text{Var}_{\theta_0}(W_{ij}) = \lambda_j \text{Var}_{\theta_0}(\nabla S_{\theta_0}(A_j, B_j))$. By the central limit theorem

$$\sqrt{n} \frac{1}{n} \sum_{i=1}^n W_{ij} \rightsquigarrow N(0, \lambda_j \text{Var}_{\theta_0}(\nabla S_{\theta_0}(A_j, B_j))).$$

The sum $\sqrt{n} \nabla \text{scl}_n(\theta_0) = \sum_{j=1}^k \beta_j \sqrt{n} \frac{1}{n} \sum_{i=1}^n W_{ij}$ is asymptotically Gaussian as well with mean zero since it converges to a sum of Gaussian distributions with mean zero. Since in the general case the random variables $\sqrt{n} \frac{1}{n} \sum_{i=1}^n W_{ij}$, $j = 1, \dots, k$ are correlated, the asymptotic variance matrix of $\sqrt{n} \nabla \text{scl}_n(\theta_0)$ needs to account for cross covariance terms leading to

$$\sqrt{n} \nabla \text{scl}_n(\theta_0) \rightsquigarrow N \left(0, \text{Var}_{\theta_0} \left(\sum_{j=1}^k \beta_j \lambda_j \nabla S_{\theta_0}(A_j, B_j) \right) \right). \quad (12)$$

We finish the proof by combining (10), (11) and (12) using Slutsky's theorem. \square

5 Stochastic Composite Likelihood for Markov Random Fields

Markov random fields (MRF) are some of the more popular statistical models for complex high dimensional data. Approaches based on pseudo likelihood and composite likelihood are naturally well-suited in this case due to the cancellation of the normalization term in the probability ratios defining conditional distributions. More specifically, a MRF with respect to a graph $G = (V, E)$, $V = \{1, \dots, m\}$ with a clique set \mathcal{C} is given by the following exponential family model

$$P_{\theta}(x) = \exp \left(\sum_{\mathcal{C} \in \mathcal{C}} \theta_{\mathcal{C}} f_{\mathcal{C}}(x_{\mathcal{C}}) - \log Z(\theta) \right), \quad Z(\theta) = \sum_x \exp \left(\sum_{\mathcal{C} \in \mathcal{C}} \theta_{\mathcal{C}} f_{\mathcal{C}}(x_{\mathcal{C}}) \right). \quad (13)$$

The primary bottlenecks in obtaining the maximum likelihood are the computations $\log Z(\theta)$ and $\nabla \log Z(\theta)$. Their computational complexity is exponential in the graph's treewidth and for many cyclic graphs, such as the Ising model or the Boltzmann machine, it is exponential in $|V| = m$.

In contrast, the conditional distributions that form the composite likelihood of (13) are given by

$$P_{\theta}(x_A|x_B) = \frac{\sum_{x'_{(A \cup B)^c}} \exp\left(\sum_{C \in \mathcal{C}} \theta_C f_C((x_A, x_B, x'_{(A \cup B)^c})_C)\right)}{\sum_{x'_{(A \cup B)^c}} \sum_{x'_A} \exp\left(\sum_{C \in \mathcal{C}} \theta_C f_C((x'_A, x_B, x'_{(A \cup B)^c})_C)\right)} \frac{Z(\theta)}{Z'(\theta)}. \quad (14)$$

The computation of (14) depends on the size of the sets A and $(A \cup B)^c$ and their intersections with the cliques in \mathcal{C} . In general, selecting small $|A|$ leads to efficient computation of the composite likelihood and its gradient. For example, in the case of $|A_j| = l, |B_j| = m - l$ with $l \ll m$ we have that $k \leq m!/(l!(m-l)!)$ and the complexity of computing the $cl(\theta)$ function and its gradient may be shown to require time that is at most exponential in l and polynomial in m . Computing the $scl(\theta)$ function and its gradient depends on the Bernoulli parameters $\lambda \in [0, 1]^k$ and the sequence of m -pairs $(A_1, B_1), \dots, (A_k, B_k)$. Selecting a sequence of m pairs that includes all $A_i = \{i\}, B_i = \{1, \dots, m\} \setminus A_i$ pairs ensures consistency. Adding pairs (A_j, B_j) with larger sets $|A_j|$ enables obtaining a specific a complexity number within a wide spectrum of available complexities by choosing appropriate mixing parameters λ . We omit the details due to lack of space.

6 Experiments

We demonstrate the asymptotic properties of $\hat{\theta}_n^{msl}$ for the Boltzmann machine and explore the complexity-accuracy tradeoff associated with several stochastic versions of $scl(\theta)$ for CRFs.

6.1 Boltzmann Machines

We illustrate the improvement in asymptotic variance of the mscl associated with adding higher order likelihood components with increasing probabilities in context of the Boltzmann machine $p_{\theta}(x) = \exp(\sum_{i < j} \theta_{ij} x_i x_j - \log \psi(\theta)), x \in \{0, 1\}^m$. To be able to accurately compute the asymptotic variance we use $m = 5$ with θ being a $\binom{5}{2}$ dimensional vector with half the components $+1$ and half -1 . Since the asymptotic variance of $\hat{\theta}_n^{msl}$ is a matrix we summarize its size using either its trace or determinant with the former having the statistical interpretation of sum of marginal variances.

We plot in Figure 1 the asymptotic variance, relative to the minimal variance of the mle, for the cases of full likelihood (FL), pseudo likelihood ($|A_j| = 1$) PL_1 , stochastic combination of pseudo likelihood and 2nd order pseudo likelihood ($|A_j| = 2$) components $\alpha PL_2 + (1 - \alpha) PL_1$, stochastic combination of 2nd order pseudo likelihood and 3rd order pseudo likelihood ($|A_j| = 3$) components $\alpha PL_3 + (1 - \alpha) PL_2$, and stochastic combination of 3rd order pseudo likelihood and 4th order pseudo likelihood ($|A_j| = 4$) components $\alpha PL_4 + (1 - \alpha) PL_3$. The graph demonstrates the computation-accuracy tradeoff as follows: (a) pseudo likelihood is the fastest but also the least accurate, (b) full likelihood is the slowest but the most accurate, (c) adding higher order components reduces the asymptotic variance but also requires more computation, (d) the variance reduces with the increase in the selection probability α of the higher order component, and (e) adding 4th order components brings the variance very close the lower limit and with each successive improvement becoming smaller and smaller according to a law of diminishing returns.

6.2 Conditional Random Fields

To demonstrate the complexity-accuracy tradeoff in a more realistic scenario we experimented with regularized maximum scl estimators for conditional random fields (CRF). We trained and tested the CRF models on local sentiment prediction data obtained in previous work [censored]. The data consisted of 249 movie review documents having an average of 30.5 sentences each with an average of 12.3 words from a 12633 word vocabulary. Each sentence was manually labeled as one of five sentimental designations: very negative, negative, objective, positive, or very positive.

Figure 2 contains the contour plots of train and test loglikelihood as a function of the scl parameters: weight β and selection probability λ . The likelihood components were mixtures of full and pseudo ($|A_j| = 1$) likelihood (rows 1,3) and pseudo and 2nd order pseudo ($|A_j| = 2$) likelihood (rows 2,4). Results were averaged over 100 cross validation iterations with 50% train-test split. We used BFGS quasi-Newton method for maximizing the regularized scl functions. Figure 2 demonstrates how the train loglikelihood increases

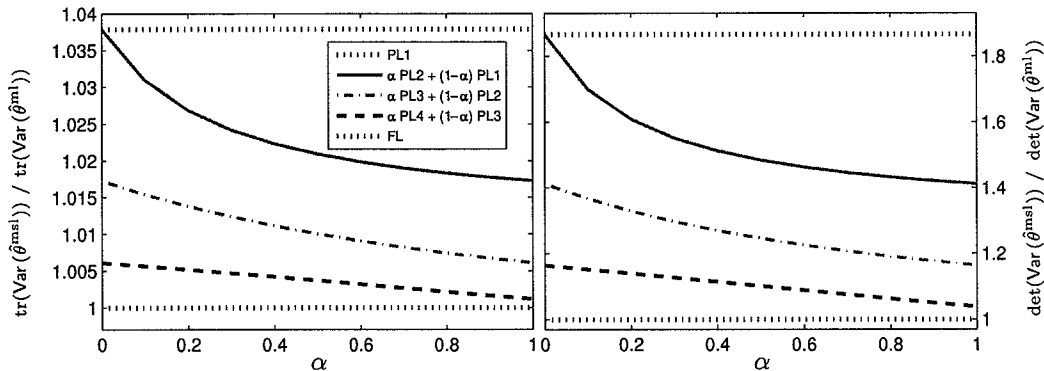


Figure 1: Asymptotic variance matrix, as measured by trace (left) and determinant (right), as a function of the selection probabilities for different stochastic versions of the scl function.

with increasing the weight and selection probability of full likelihood in rows 1,3 and of 2nd order pseudo likelihood in rows 2,4. This increase in train loglikelihood is also correlated with an increase in computational complexity as higher order likelihood components require more computation.

It is interesting to contrast the test loglikelihood behavior in the case of mild ($\sigma = 10$) and stronger ($\sigma = 1$) L_2 regularization. In the case of weaker or no regularization, the test loglikelihood shows different behavior than the train loglikelihood. Adding a lower order component such as pseudo likelihood acts as a regularizer that prevents overfitting. Thus, in cases that are prone to overfitting reducing higher order likelihood components improves both performance as well as complexity. This represents a win-win situation in contrast to the classical view where the mle has the lowest variance and adding lower order components reduces complexity but increases the variance.

Figure 3 displays the complexity and train (left) and test (right) negative loglikelihood of different scl estimators as points in a two dimensional space. The shaded area near the origin is unachievable as no scl estimator can achieve high accuracy and low computation at the same time. The optimal location in this 2D plane is the curved boundary of the achievable region with the exact position on that boundary depending on the required solution of the computation-accuracy tradeoff.

7 Discussion

The proposed estimator family facilitates computationally efficient estimation in complex graphical models. In particular, different parameterization of the stochastic likelihood enables the resolution of the complexity-accuracy tradeoff in a domain and problem specific manner. The framework is generally suited for Markov random fields, including conditional graphical models and is theoretically motivated. When the model is prone to overfit, stochastically mixing lower order components with higher order ones acts as a regularizer and results in a win-win situation of improving test-set accuracy and reducing computational complexity at the same time.

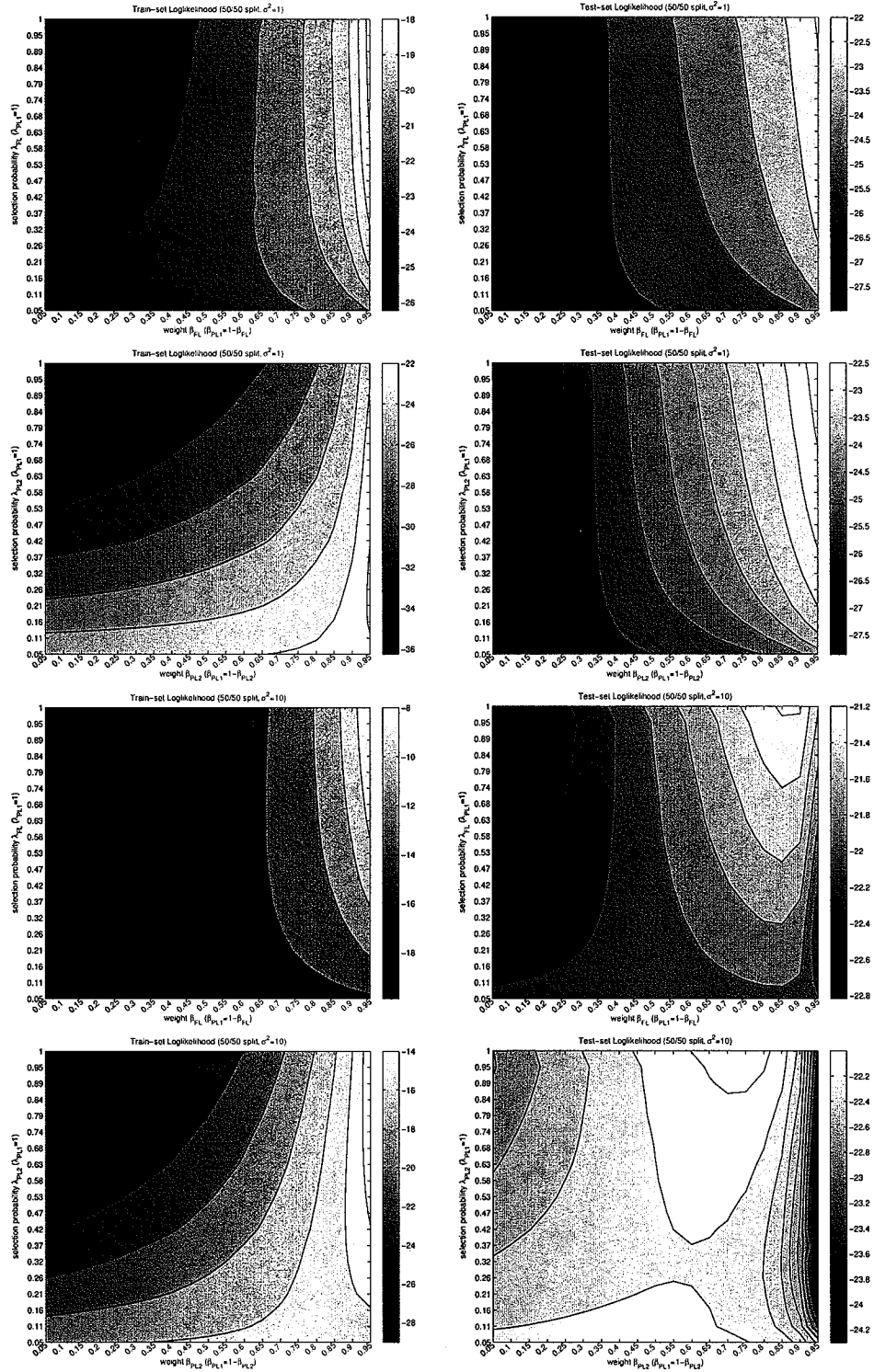


Figure 2: Train (left) and test (right) loglikelihood contours for maximum sel estimators for the CRF model. L_2 regularization parameters are $\sigma^2 = 1$ (rows 1,2) and $\sigma^2 = 10$ (rows 3,4). Rows 1,3 represent stochastic mixtures of full (FL) and pseudo (PL₁) likelihood components and rows 2,4 represent stochastic mixtures of pseudo (PL₁) and 2nd order pseudo (PL₂) likelihood components.

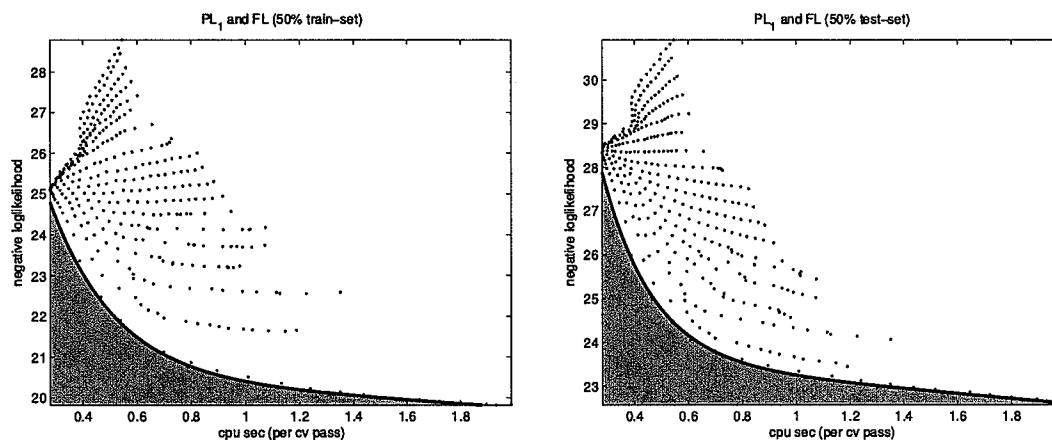


Figure 3: Scatter plot representing complexity and negative loglikelihood (left: train, right: test) of scl functions for CRFs with regularization parameter $\sigma^2 = 0.5$. The points represent different stochastic combinations of full and pseudo likelihood components. Unachievable region is shaded.

References

- [1] B. Arnold and D. Strauss. Pseudolikelihood estimation: some examples. *Sankhya B*, 53:233–243, 1991.
- [2] J. Besag. Spatial interaction and the statistical analysis of lattice systems (with discussion). *J Roy Statist Soc B*, 36(2):192–236, 1974.
- [3] R. Casella and C. Robert. *Monte Carlo Statistical Methods*. Springer Verlag, second edition, 2004.
- [4] T. S. Ferguson. *A Course in Large Sample Theory*. Chapman & Hall, 1996.
- [5] N. Hjort and C. Varin. ML, PL, and QL in markov chain models. *Scand J Stat*, 35(1):6482, 2008.
- [6] G. Liang and B. Yu. Maximum pseudo likelihood estimation in network tomography. *IEEE T Signal Proces*, 51(8):2043–2053, 2003.
- [7] B. G. Lindsay. Composite likelihood methods. *Contemporary Mathematics*, 80:221–239, 1988.
- [8] C. Sutton and A. McCallum. Piecewise pseudolikelihood for efficient training of conditional random fields. In *Proc. of the International Conference on Machine Learning*, 2007.
- [9] C. Varin and P. Vidoni. A note on composite likelihood inference and model selection. *Biometrika*, 92:519–528, 2005.
- [10] E. P. Xing, M. I. Jordan, and S. Russell. A generalized mean field algorithm for variational inference in exponential families. In *Uncertainty in Artificial Intelligence*, 2003.
- [11] S. C. Zhu and X. Liu. Learning in Gibbsian fields: How accurate and how fast can it be? *IEEE T Pattern Anal*, 24(7):1001–1006, 2002.