

Supplementary Material: Revisiting Uncertainty in Graph Cut Solutions

Daniel Tarlow
 Dept. of Computer Science
 University of Toronto
 dtarlow@cs.toronto.edu

Ryan P. Adams
 School of Engineering and Applied Sciences
 Harvard University
 rpa@seas.harvard.edu

Analysis of Kohli-Torr as Approximate Gibbs Distribution Inference

Here we treat the min-marginal uncertainty of Kohli & Torr (KT) [1] as a procedure for computing approximate marginal probabilities under the Gibbs distribution that is standard in conditional random field (CRF) models. In the very simple case where there are no interactions between variables in the model (i.e., a “unary only” model), the approximate KT marginals match the true marginals. The main result here is that even after adding just submodular pairwise interactions to the model, the marginals estimated by KT can be off by a factor that is exponentially large in the number of variables in the network.

The crux of the problem (which will be made more precise below) is fairly simple: if we approximate the sum of probabilities over a subset of configurations by instead maximizing over the probabilities of states in the subset, how bad can the approximation be? The analysis simply considers how to parameterize models such that the summation required to compute marginals gives as different value as possible from the maximizations used by KT.

We are concerned with this question, because when learning with the KT procedure, approximate gradients are calculated using these approximate marginals. This analysis adds some theoretical explanation for the empirical behavior of the KT learning procedure in our experiments.

We will consider models over a set of binary variables $\mathbf{y} \in \mathcal{Y} = \{0, 1\}^D$. We assume we have an energy function $E(\mathbf{y})$ (for simplicity dropping the dependence on input \mathbf{x} and the explicit dependence on parameters \mathbf{w}), then under the Gibbs distribution, the probability of a configuration \mathbf{y} is

$$p'(\mathbf{y}) = \exp\{-E(\mathbf{y})\} \quad (1)$$

$$\mathcal{Z} = \sum_{\mathbf{y}' \in \mathcal{Y}} p'(\mathbf{y}') \quad (2)$$

$$p(\mathbf{y}) = \frac{1}{\mathcal{Z}} p'(\mathbf{y}). \quad (3)$$

The marginal probability for y_d under the Gibbs distribution

is

$$p(y_d) = \sum_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y_d} p(\hat{\mathbf{y}}) \quad (4)$$

$$= \frac{\sum_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y_d} \exp\{-E(\hat{\mathbf{y}})\}}{\sum_{\mathbf{y}' \in \{0,1\}} \sum_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y'_d} \exp\{-E(\hat{\mathbf{y}})\}}. \quad (5)$$

That is, we are summing over the probability of all configurations \mathbf{y}' that are consistent with the given setting of y_d .

The KT approximation is to define

$$\tilde{p}(y_d) = \frac{\exp\{-\min_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y_d} E(\hat{\mathbf{y}})\}}{\sum_{\mathbf{y}' \in \{0,1\}} \exp\{-\min_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y'_d} E(\hat{\mathbf{y}})\}} \quad (6)$$

$$= \frac{\max_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y_d} \exp\{-E(\hat{\mathbf{y}})\}}{\sum_{\mathbf{y}' \in \{0,1\}} \max_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y'_d} \exp\{-E(\hat{\mathbf{y}})\}}. \quad (7)$$

Comparing Eq. 5 to Eq. 7, we see that they are identical, except the maximizations in Eq. 7 are replaced by sums to get the expression Eq. 5 for the true marginals. For convenience, let \tilde{p}' be the unnormalized value of the maximization:

$$\tilde{p}'(y_d) = \max_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y_d} \exp\{-E(\hat{\mathbf{y}})\}. \quad (8)$$

and similarly

$$p'(y_d) = \sum_{\hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = y_d} \exp\{-E(\hat{\mathbf{y}})\}. \quad (9)$$

The question we ask in the following is, “how large can we make the approximation factor $f = \tilde{p}(y_d)/p(y_d)$, for some d and y_d ?” In other words, in the worst case, how bad can the approximation of KT be?

0.1. General Models

We begin with the simplest case, where there are no restrictions on the parameterization of the energy function. In this case, we have the freedom to assign an independent probability value to each of the 2^D joint configurations \mathbf{y} , so long as we restrict the probabilities to be non-negative, and the sum to be equal to 1.

Without loss of generality, suppose that the KT approximate marginal is $\tilde{p}(y_d = 1) = \lambda$, and thus $\tilde{p}(y_d = 0) = 1 - \lambda$, for $\lambda \in [0, 1]$. This implies that

$$\frac{\tilde{p}'(y_d = 1)}{\tilde{p}'(y_d = 0) + \tilde{p}'(y_d = 1)} = \lambda. \quad (10)$$

Again without loss of generality, we can further assume that $\tilde{p}'(y_d = 0) + \tilde{p}'(y_d = 1) = 1$, so

$$\tilde{p}'(y_d = 1) = \lambda, \text{ and} \quad (11)$$

$$\tilde{p}'(y_d = 0) = 1 - \lambda. \quad (12)$$

Let S_0 be the set of states where $y_d = 0$:

$$S_0 = \{\hat{\mathbf{y}} \mid \hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = 0\} \quad (13)$$

and define S_1 similarly:

$$S_1 = \{\hat{\mathbf{y}} \mid \hat{\mathbf{y}} \in \mathcal{Y}, \hat{y}_d = 1\} \quad (14)$$

So far, we know that

$$\max_{\hat{\mathbf{y}} \in S_1} p'(\hat{\mathbf{y}}) = \lambda \quad (15)$$

$$\max_{\hat{\mathbf{y}} \in S_0} p'(\hat{\mathbf{y}}) = 1 - \lambda. \quad (16)$$

The size of these sets are $|S_0| = |S_1| = 2^{D-1}$, so the largest that $p'(y_d = 1)$ could be is $\lambda \cdot 2^{D-1}$. The proof of this is straightforward: suppose that $p'(y_d = 1) > \lambda \cdot 2^{D-1}$. Then there must be some configuration $\mathbf{y} \in S_1$ such that $p'(\mathbf{y}) > \lambda$. But this contradicts Eq. 15. Finally, $p'(y_d = 1) = \lambda \cdot 2^{D-1}$ is achievable by assigning $p'(\mathbf{y}) = \lambda$ for all $\mathbf{y} \in S_1$.

The smallest that $p'(y_d = 0)$ could be is $1 - \lambda$, because probabilities are non-negative, and at least one configuration $\mathbf{y} \in S_0$ must have $p'(\mathbf{y}) = 1 - \lambda$ to be consistent with Eq. 16.

Since S_0 and S_1 are disjoint, the minimum achievable value of $p(y_d = 0)$ consistent with Eq. 15 and Eq. 16 comes from simultaneously maximizing $p'(y_d = 1)$ and minimizing $p'(y_d = 0)$. This yields

$$p(y_d = 0) = \frac{p'(y_d = 0)}{p'(y_d = 0) + p'(y_d = 1)} \quad (17)$$

$$= \frac{1 - \lambda}{1 - \lambda + \lambda \cdot 2^{D-1}}. \quad (18)$$

The approximation factor is then:

$$\frac{\tilde{p}(y_d = 0)}{p(y_d = 0)} = 1 - \lambda + \lambda \cdot 2^{D-1}, \quad (19)$$

which is exponentially large in D for constant values of λ (say $\lambda = .5$).

0.2. Pairwise Submodular Models

We now consider the case where we restrict the parameterization, so that we cannot choose arbitrary values for the probability assigned to a joint assignment \mathbf{y} . Instead, we require that the model be a pairwise graphical model with only submodular interactions. In this section, we show that there are such models where the same basic strategy as before is applicable, and thus the approximation factor is still exponentially bad.

Consider a pairwise model where there is a single root, y_r , which shares an edge with all other variables. There are no other edges, so this is also a tree-structured model. Let θ be the set of potentials, and define the unnormalized probability as a typical graphical model, where the only unary potential is on the root:

$$p'(\mathbf{y}) = \theta_r(y_r) \prod_d \theta_{rd}(y_r, y_d). \quad (20)$$

Let $\theta_r(1) = \lambda$, $\theta_r(0) = 1 - \lambda$ for $\lambda \in [0, 1]$, and parameterize pairwise interactions for all d as follows:

$$\theta_{rd}(y_r = 0, y_d = 0) = 1 \quad (21)$$

$$\theta_{rd}(y_r = 0, y_d = 1) = 0 \quad (22)$$

$$\theta_{rd}(y_r = 1, y_d = 0) = 1 \quad (23)$$

$$\theta_{rd}(y_r = 1, y_d = 1) = 1. \quad (24)$$

This is clearly a submodular interaction. Now, if y_r is constrained to be 1, there are 2^{D-1} joint configurations with support, each getting an unnormalized probability of λ . If y_r is constrained to 0, then there is only a single joint configuration with support (the all 0's assignment), and it gets unnormalized probability of $1 - \lambda$. This model is then equivalent to the model constructed in the previous section, and thus the approximation factor is also exponentially bad.

Bootstrap Analysis Results

See Fig. 1 and Fig. 2.

References

- [1] P. Kohli and P. H. S. Torr. Measuring uncertainty in graph cut solutions. *Computer Vision and Image Understanding*, 112(1):30–38, 2008. 1

		Log Lik	Accuracy	AUC
Aero	KT-final	$-.35 \pm .005$ ($-.29 \pm .004$)	$87.3 \pm .2$ ($89.7 \pm .1$)	$.85 \pm .018$ ($.87 \pm .020$)
	KT- f_{best}	$-.32 \pm .007$ ($-.28 \pm .005$)	$88.1 \pm .3$ ($89.9 \pm .2$)	$.85 \pm .018$ ($.87 \pm .019$)
	Ours	$-.26 \pm .007$ ($-.24 \pm .005$)	$88.9 \pm .3$ ($90.4 \pm .2$)	$.90 \pm .014$ ($.90 \pm .014$)
Car	KT-final	$-.48 \pm .008$ ($-.65 \pm .008$)	$84.1 \pm .3$ ($78.7 \pm .3$)	$.69 \pm .036$ ($.65 \pm .022$)
	KT- f_{best}	$-.39 \pm .007$ ($-.51 \pm .006$)	$86.2 \pm .3$ ($80.0 \pm .3$)	$.66 \pm .029$ ($.62 \pm .023$)
	Ours	$-.35 \pm .009$ ($-.51 \pm .006$)	$86.1 \pm .5$ ($81.4 \pm .3$)	$.76 \pm .035$ ($.65 \pm .032$)
Cow	KT-final	$-.54 \pm .009$ ($-.64 \pm .019$)	$79.7 \pm .3$ ($75.9 \pm .8$)	$.76 \pm .040$ ($.77 \pm .038$)
	KT- f_{best}	$-.47 \pm .004$ ($-.52 \pm .008$)	$80.9 \pm .3$ ($76.7 \pm .5$)	$.66 \pm .040$ ($.65 \pm .038$)
	Ours	$-.38 \pm .005$ ($-.41 \pm .006$)	$82.5 \pm .4$ ($79.9 \pm .4$)	$.84 \pm .023$ ($.82 \pm .034$)
Dog	KT-final	$-.52 \pm .009$ ($-.45 \pm .004$)	$81.8 \pm .4$ ($84.7 \pm .1$)	$.64 \pm .029$ ($.66 \pm .027$)
	KT- f_{best}	$-.43 \pm .004$ ($-.38 \pm .003$)	$84.1 \pm .2$ ($86.7 \pm .1$)	$.62 \pm .023$ ($.66 \pm .026$)
	Ours	$-.38 \pm .004$ ($-.34 \pm .004$)	$84.0 \pm .2$ ($86.8 \pm .2$)	$.76 \pm .025$ ($.79 \pm .028$)

Figure 1. Results for binary models. Format is “Train (Test)”. \pm indicates the standard deviation of the result under bootstrap resamplings.

		n/U Before	n/U After	Change
Aero	KT-final	63.4 ± 2.2	63.9 ± 2.1	$.5 \pm .3$
	KT- f_{best}	60.7 ± 2.1	61.8 ± 1.9	$1.1 \pm .4$
	Ours	64.1 ± 2.2	66.6 ± 2.0	$2.5 \pm .7$
Car	KT-final	43.5 ± 1.4	44.2 ± 1.2	$.7 \pm .3$
	KT- f_{best}	40.7 ± 1.0	$43.3 \pm .9$	$2.6 \pm .5$
	Ours	41.8 ± 1.0	45.9 ± 1.0	$4.1 \pm .8$
Cow	KT-final	51.2 ± 3.7	51.7 ± 3.6	$.5 \pm .5$
	KT- f_{best}	40.7 ± 1.8	47.1 ± 2.7	6.4 ± 2.0
	Ours	46.5 ± 1.9	52.8 ± 2.8	6.3 ± 1.8
Dog	KT-final	52.1 ± 2.6	52.0 ± 2.5	$-.1 \pm .5$
	KT- f_{best}	44.5 ± 1.1	47.2 ± 1.2	$2.7 \pm .6$
	Ours	48.9 ± 1.4	55.3 ± 1.9	6.4 ± 1.3

Figure 2. Test results for maximizing surrogate expected $\frac{n}{U}$ score. “Before” corresponds to predicting the mode of Q; “After” is the prediction from our expected score maximization routine. \pm indicates the standard deviation of the result under bootstrap resamplings.