# Parametric Majorization for Data-Driven Energy Minimization Methods

Jonas Geiping          Michael Moeller

Department of Electrical Engineering and Computer Science, University of Siegen

{jonas.geiping, michael.moeller}@uni-siegen.de

## Abstract

*Energy minimization methods are a classical tool in a multitude of computer vision applications. While they are interpretable and well-studied, their regularity assumptions are difficult to design by hand. Deep learning techniques on the other hand are purely data-driven, often provide excellent results, but are very difficult to constrain to predefined physical or safety-critical models. A possible combination between the two approaches is to design a parametric energy and train the free parameters in such a way that minimizers of the energy correspond to desired solution on a set of training examples. Unfortunately, such formulations typically lead to bi-level optimization problems, on which common optimization algorithms are difficult to scale to modern requirements in data processing and efficiency. In this work, we present a new strategy to optimize these bi-level problems. We investigate surrogate single-level problems that majorize the target problems and can be implemented with existing tools, leading to efficient algorithms without collapse of the energy function. This framework of strategies enables new avenues to the training of parameterized energy minimization models from large data.*

## 1. Introduction

Energy minimization methods, also referred to as variational methods, are a classical tool in computer vision [83, 18, 32, 14]. The idea is to define a data-dependent cost function $E$ that assigns a value to each candidate solution $x$. The desired optimal solution is then the target solution with the lowest energy value. This methodology has several advantages, for one, it is characterized by an *explicit model* - namely the energy function to be minimized - and an implicit inference method - how we compute the minimizer of this energy is a separate problem. This duality allows a fruitful analysis, leading to controllable methods with provable guarantees that are paramount in many critical applications [80, 78, 98]. Furthermore, explicit knowledge over the model structure allows for explainable and clear modifications when the method is applied in a related task [26].

Conversely, deep learning approaches [60], specifically deep feed-forward neural networks work by very different principles. The methodology of deep learning is characterized by *implicit* models and explicit inference. The solution to the problem at hand is given directly by the output of the learned feed-forward structure. This is advantageous in practice and crucial for the efficient training of neural networks, however the underlying model of the problem structure is now only implicitly contained in the responses of the network. Deep neural networks have fundamentally changed the state-of-the-art in various computer vision applications, due to these properties as the inference operations are learned directly from large amounts of training data. These approaches are able to learn expressive and convincing mechanisms, examples of which can be found not only in recognition tasks (e.g. [56]), but also in denoising [99], optical flow [70, 49] or segmentation tasks [64, 81, 21]. Yet, as the underlying model is only implicitly defined and 'hidden' in the network structure, it is difficult to modify it for applications in other domains or to guarantee specific outputs. Domain adaptation is still an active field of research and several examples, for instance in medical imaging [3, 38], have demonstrated the need for possibly model-based physically plausible output restrictions. This problem is most strikingly demonstrated by the phenomenon of adversarial examples [89] - the existence of input data, that, when fed through the network, leads to highly erroneous solutions. While one would expect that such behaviour is possibly unavoidable in recognition tasks [87, 71], it should not be a factor in low-level computer vision applications.

Reviewing these two methodologies, we would - of course - prefer to have the best of both worlds. We would like to use both the large amounts of data at our disposal and our far-reaching domain knowledge in many tasks to train explicit models with a significant number of free parameters, so that their optimal solutions are similar to directly trained feed-forward networks.

A promising candidate for such a combination of learning- and model based approaches are *parametrized energy minimization methods*. The idea of such methods is to

define an energy $E$ that depends on the candidate solutions $x$, the input data $y$ and parameters $\theta$,

$$
\begin{aligned}
E : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^s &\to \mathbb{R}, \\
(x, y, \theta) &\mapsto E(x, y, \theta),
\end{aligned} \tag{1}
$$

such that for a good choice of parameters $\theta$, the argument $x(\theta) = \arg\min_x E(x, y, \theta)$ that minimizes the energy over all $x$ is as close a possible to the desired true solution $x^*$.

To train such parametric energies, assume we are given $N$ training samples $\{(x_i^*, y_i)\}_{i=1}^N$ and a continuous *higher-level* loss function $l : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$, which measures the deviation of solutions of the model to the given training samples. Determining the optimal parameters $\theta$ then becomes a *bi-level optimization problem* combining both the higher-level loss function and the lower-level energy,

$$
\min_{\theta \in \mathbb{R}^s} \sum_{i=1}^N l(x_i^*, x_i(\theta)), \tag{2}
$$

$$
\text{subject to} \qquad x_i(\theta) = \arg\min_{x \in \mathbb{R}^n} E(x, y_i, \theta). \tag{3}
$$

Usual first-order learning methods are difficult to apply in this setting. For every gradient computation it is necessary to compute a derivative of the $\arg\min$ operation of the lower-level problem, which is even further complicated if we consider parametrized non-smooth energy models which are wide-spread in computer vision [32, 14].

Therefore, the goal of this paper is to analyze bi-level optimization problems and identify strategies that allow for efficient approximate solutions. We investigate single-level minimization problems with simple constraints without second-order differentiation, which are applicable even to non-smooth energies. Such forms allow scaling the previously limited training of energy minimization methods in computer vision to larger datasets and increase the effectiveness in applications where it is critical that the solution follows a specific model structure.

In the remainder of this paper we analyze the bi-level optimization problem to develop a rigorous understanding of sufficient conditions for a single-level surrogate strategy for continuous loss functions $l$ and convex, non-smooth lower-level energies $E$ to be successful. We introduce the concept of a *parametric majorization function*, show relations to structured support vector machines and provide several levels of parametric majorization functions with varying levels of exactness and computational effort. We extend our approximations to an iterative scheme, allowing for repeated evaluations of the approximation, before illustrating the proposed strategies in computer vision applications.

## 2. Related Work

The straightforward way of optimizing bi-level problems is to consider *direct descent methods* [55, 85, 30]. These methods directly differentiate the higher-level loss function with respect to the minimizing argument and descend in the direction of this gradient. An incomplete list of examples in image processing is [13, 26, 24, 25, 33, 34, 41, 45, 46]. This strategy requires both the higher- and lower-level problems to be smooth and the minimizing map to be invertible. This is usually facilitated by implicit differentiation, as discussed in [84, 57, 25, 26]. In more generality, the problem of directly minimizing $\theta$ without assuming that smoothness in $E$ leads to optimization problems with equilibrium constraints (MPECs), see [9] for a discussion in terms of machine learning or [36, 35, 37] and [30]. This approach also applies to the optimization layers of [2], which lend themselves well to a reformulation as a bi-level optimization problem.

*Unrolling* is a prominent strategy in applied bi-level optimization across fields, i.e. MRF literature [4, 69] in deep learning [100, 22, 19, 63] and in variational settings [73, 59, 58, 43, 44, 77]. The problem is transformed into a single level problem by choosing an optimization algorithm $\mathcal{A}$ that produces an approximate solution to the lower level problem after a fixed number of iterations. $x(\theta)$ is then replaced by $\mathcal{A}(y, \theta)$. Automatic differentiation [42] allows for an efficient evaluation of the gradient of the upper-level loss w.r.t to this reduced objective

$$
\min_{\theta} \sum_{i=1}^N l(x_i^*, \mathcal{A}(y_i, \theta)). \tag{4}
$$

In general these strategies are very successful in practice, *because* they combine the model and its optimization method into a single feed-forward process, where the model is again only implicitly present. Later works [27, 23, 43, 44] allow the lower-level parameters to change in between the fixed number of iterations, leading to structures that model differential equations and stray further from underlying modelling. As pointed out in [53], these strategies are more aptly considered as a set of nested quadratic lower-level problems.

Several techniques have been developed in the field of structured support vector machines (SSVMs) [92, 28, 1, 95] that are very relevant to the task of learning energy models, as SSVMs can be understood as bi-level problems with a lower-level energy that is linear in $\theta$ and often a non-continuous higher-level loss. Various strategies such as margin rescaling [92], slack rescaling [95, 97], softmax-margins [40] exist and have also been applied recently in the training of computer vision models in [54, 29], we will later return to their connection to the investigated strategies.

## 3. Bi-Level Learning

We now formalize our learning problem. We assume the lower-level energy $E$ from (1) to be convex (but not necessarily smooth) in its first variable $x \in \mathbb{R}^n$ and to depend

continuously on input data $y \in \mathbb{R}^m$ and parameters $\theta \in \mathbb{R}^s$. We assume its minimizer $x(\theta)$ to be unique. For our higher-level loss function (2) $l : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$, we assume that it fulfills $l(x, y) \geq 0, l(x, x) = 0$ for all $x, y$ and is differentiable in its second argument.

Note that this formulation of bi-level optimization problems directly generalizes classical supervised (deep) learning with a network $\mathcal{N}(\theta, y)$ via the quadratic energy $E(x, y_i, \theta) = \frac{1}{2}||x - \mathcal{N}(\theta, y_i)||^2$, for which $x_i(\theta) = \mathcal{N}(\theta, y_i)$.

*Preliminaries (Convex Analysis):* Let us summarize our notation and some fundamental results from convex analysis. We refer the reader to [6] for more details. We denote by $\partial E(x)$ the set of subgradients of a convex function $E$ at $x$. We define the Bregman distance between two vectors relative to a convex function $E$ by $D_E^p(x, y) = E(x) - E(y) - \langle p, x - y \rangle$ for a subgradient $p \in \partial E(y)$, intuitively the Bregman distance measures the difference of the energy at $x$ to its linear lower bound around $y$. $E^*(p) = \sup_x \langle p, x \rangle - E(x)$ is the convex conjugate of $E$. $x$ is a minimizer of the energy $E$ if and only if $0 \in \partial E(x)$ or equivalently by convex duality $x \in \partial E^*(0)$. $E$ is $m$-strongly convex if $D_E^p(x, y) \geq \frac{m}{2}||x - y||^2$ for all $x, y$. Conversely, if $E$ is $m$-strongly convex, then $E^*$ is $\frac{1}{m}$-strongly smooth, i.e. $D_{E^*}(p, q) \leq \frac{2}{m}||p - q||^2$. Furthermore $D_E^p(x, y) = D_{E^*}^x(p, q), q \in \partial E(x)$ holds for all Bregman distances [11]. We consider parametrized energies in several variables, yet we always assume (sub)-gradients, Bregman distances and convex conjugates to be with respect to the first argument $x$.

## 3.1. Majorization of Bi-level Problems

As previously discussed, directly solving the bi-level problem as posed in Eq. (2) and (3) is tricky. We need to implicitly differentiate the minimizing argument $x_i(\theta)$ for all $N$ samples just to apply a first-order method in $\theta$ - which is in stark contrast to our goal of finding efficient and scalable algorithms.

Let us instead look at the problem from a very different angle and entertain the idea that the loss function $l$ is actually of secondary importance to us. We really only want to find parameters $\theta$ so that our training samples are well reconstructed, $x_i^* \approx x_i(\theta)$. If we go so far as to assume that the loss value of our optimal parameters $\theta^*$ is zero, meaning that minimizers of our energy are perfectly able to reconstruct our training samples, then the bi-level problem is reduced to a single-level problem, inserting $x_i^* = x_i(\theta^*)$:

$$\min_\theta \quad \text{s.t.} \ 0 \in \partial E(x_i^*, y_i, \theta), \tag{5}$$

which we could solve via

$$\min_\theta \sum_{i=1}^N ||q_i||^2 \quad \text{s.t.} \ q_i \in \partial E(x_i^*, y_i, \theta) \tag{6}$$

This train of thought is closely interconnected to the notion of separability in Support Vector Machine methods [96], where it is assumed that given training samples are linearly separable, which is equivalent to assuming that the classification loss is zero on the training set.

However minimizing Eq. (6) is often not a good choice. A simple example is $E(x, y, \theta) = (\theta x - y)^2$, i.e. we simply try to learn a positive scaling factor $\theta$ between $x$ and $y$. Problem (5) can then be written as $\min_\theta \sum_i (\theta^2 x_i^* - \theta y_i)^2$ and is trivially minimized by $\theta = 0$. Such a solution makes $E$ independent of $x$ such that every $x$ becomes a minimizer. This phenomenon is referred to as *collapse* of the energy function [62, 61] in machine learning literature, and clearly cannot be a good strategy to learn a scaling factor.

Interestingly, the scaling problem can be reformulated into a reasonable (non-collapsing) problem, if we require (6) to *majorize* the bilevel problem: If we consider the higher-level loss function $l(x_i^*, x_i(\theta)) = (x_i^* - x_i(\theta))^2$, then our surrogate problem $\sum_i (\theta^2 x_i^* - \theta y_i)^2$ is clearly not a majorizer for arbitrary $\theta$. However, if we consider a reformulation of the energy to $E(x) = (x - \frac{1}{\theta}y)^2$, then this reformulation leads to a *majorizing* surrogate $\sum_i (x_i^* - \frac{1}{\theta}y_i)^2$. Minimizing $\theta$ now leads to learning the desired scaling factor.

Our toy example motivates us to formalize the concept of majorizing surrogates:

**Definition 1** (Parametrized Majorizer)**.** Given a bi-level optimization problem in the higher level loss $l(x, y)$ and lower-level energy $E(x, y, \theta)$, we call the function $S(x, y, \theta) : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^s \to \mathbb{R}$ a parametrized majorizer, if

$$\forall \theta \in \mathbb{R}^s : \qquad l(x, x(\theta)) \leq S(x, y, \theta)$$
$$\forall \theta \in \mathbb{R}^s \quad \text{s.t.} \quad l(x, x(\theta)) = 0 \implies S(x, y, \theta) = 0$$

hold for any $x, y \in \mathbb{R}^n \times \mathbb{R}^m$.

This definition allows us to formalize our objective further. We investigate replacing the bi-level optimization problem (2), (3) by the minimization of a suitable parametrized majorizer, i.e.

$$\min_{\theta \in \mathbb{R}^s} \sum_{i=1}^N S(x_i^*, y_i, \theta). \tag{7}$$

An immediate conclusion of Definition 1 is that the function $S$ now certifies our progress as $S(x, y, \theta) = 0$ implies $l(x, x(\theta)) = 0$. Moreover, our goal is to choose majorizers $S$ in such a way that they yield **single-level** problems (7), meaning it is not necessary to differentiate an $\arg\min$ operation to minimize them or to solve an equally difficult reformulation, making them significantly easier to solve.

## 3.2. Single-Level Majorizers

One possible way to find a majorizer that satisfies the previously postulated properties is by considering the majorizer naturally induced through the Bregman distance of the lower level energy. We assume the following condition

$$l(x, z) \leq D_{E_\theta}(x, z) \quad \forall x, z \in \mathbb{R}^n, \theta \in \mathbb{R}^s, \qquad (8)$$

and propose the surrogate problem

$$\min_\theta \ \sum_{i=1}^N D_{E_\theta}\left(x_i^*, x_i(\theta)\right). \qquad (9)$$

Condition (8) is an assumption on both the loss function and the energy. It thus delineates the class of bi-level problems that can be attacked with this majorization strategy. However this condition is quite general. For a large class of loss functions, we only need the energy to contain a term that also induces the loss function, a property also known as (relative) strong convexity [94, 65]:

**Proposition 1.** *If the loss function $l(x, y)$ is a Bregman distance induced by a strictly convex function $w : \mathbb{R}^n \to \mathbb{R}$, i.e. $l(x, y) = D_w(x, y)$, then assumption (8) is fulfilled if the energy $E$ is $w$-strongly convex, i.e. if $E(x) - w(x)$ is still a convex function.*

*Proof:* We write $E$ as $E(x) = \hat{E}(x) + w(x)$ and apply the additive separability of Bregman distances to find $D_E(x, y) = D_{\hat{E}}(x, y) + D_w(x, y)$, which is greater than or equal to $D_w(x, y)$, as $D_{\hat{E}}(x, y)$ is non-negative due to the convexity of $\hat{E}$. For the usual euclidean loss, this property reduces to strong convexity:

**Example 1.** *If the loss function is given by a squared Euclidean loss, $l(x, y) = \frac{1}{2}||x - y||^2$ and the energy is $m$-strongly convex, then assumption (8) is fulfilled for the energy $\frac{1}{m}E$.*

The question remains whether the proposed surrogate problem (9) is efficiently solvable. We especially wanted to circumvent the differentiation of $x(\theta)$. However $D_E\left(x_i^*, x_i(\theta)\right)$ is much easier to solve, in comparison to the original bi-level problem, as we can see in both its primal and its dual formulation. First, from a primal viewpoint, we have

$$\begin{aligned} & D_E\left(x_i^*, x_i(\theta)\right) \\ =& E(x_i^*, y_i, \theta) - E(x_i(\theta), y_i, \theta) - \langle p_i, x_i^* - x_i(\theta)\rangle, \end{aligned}$$

for some subgradient $p_i \in \partial E(x_i(\theta))$ which we have not specified yet. But, as $0 \in \partial E(x_i(\theta))$ as $x_i(\theta)$ is by definition a solution to the lower-level problem, we may take $p = 0$ and simplify to

$$E(x_i^*, y_i, \theta) - E(x_i(\theta), y_i, \theta).$$

Now $x_i(\theta)$ is contained solely in $E$ and we can write

**Bregman Surrogate:**

$$D_{E_\theta}^0\left(x_i^*, x_i(\theta)\right) = \max_{x \in \mathbb{R}^n} E(x_i^*, y_i, \theta) - E(x, y_i, \theta). \qquad (10)$$

This surrogate function is already much simpler than the original bi-level problem. We can minimize (10) either by alternating minimization in $\theta$ and maximization in $x$ or by jointly optimizing both variables. However, the problem is still set up as a saddle-point problem which is not ideal for optimization.

*Remark.* Interestingly, this discriminative formulation is not wholly unfamiliar. We can understand this as an appropriate generalization of generalized perceptron training [62, 61, 90] as discussed as far back as [82]. See the appendix for further details. In vein of this comparison, conditions 1 and 2 from e.g. [62], i.e. conditions on the existence of a margin between the optimal solution and other candidate solutions central to (S)SVM methods [96, 91, 93] are reflected in Proposition 1 in the convex continuous setting. Due to continuity of the energy and loss function we cannot obey a fixed margin, yet we impose that the energy grows at least as fast as the loss function, when we move away from the optimal solution.

We can resolve the saddle-point question by analyzing the surrogate (9) from a dual standpoint, as by Bregman duality [10]

$$D_{E_\theta}^0\left(x_i^*, x_i(\theta)\right) = D_{E_\theta^*}^{x_i^*}(0, q_i) \qquad (11)$$

for $q_i \in \partial E(x_i^*, y, \theta)$. Contrasting this formulation with our initial goal of penalizing the subgradient as in Eq. (6), we see that the Bregman distance induced by $E^*$ is the natural 'distance' by which to penalize the subgradient in the sense that penalizing the subgradient at $x_i^*$ with this generalized distance recovers a majorizing surrogate.

We can further simplify the dual formulation by applying Fenchel's theorem:

$$D_{E_\theta^*}^{x_i^*}(0, q_i) = E(x_i^*, y_i, \theta) + E^*(0, y_i, \theta). \qquad (12)$$

Computing $E^*(0)$ is exactly as difficult as minimizing $E$ (as $E^*(0) = \min_x E(x)$), so we need to rewrite this surrogate in a tractable manner. To do so, we assume that $E$ can be additively decomposed into two parts,

$$E(x, y, \theta) = E_1(x, y, \theta) + E_2(x, y, \theta), \qquad (13)$$

where both $E_1$ and $E_2$ are convex in their first argument and their convex conjugates are simple to compute. Exploiting that $E^*(0) = \min_z E_1^*(-z) + E_2^*(z)$ yields

$$D_{E_\theta^*}^{x_i^*}(0, q_i) = \min_{z \in \mathbb{R}^n} E(x_i^*, y_i, \theta) + E_1^*(-z, y, \theta) + E_2^*(z, y, \theta). \qquad (14)$$

In comparison to the primal formulation in Eq (10), we have now reformulated the problem from a saddle point problem (minimizing in $\theta$ and maximizing in $x$) to a pure minimization problem, which is easier to handle. This is a generalization of the dual formulation discussed in the linear context of SSVMs for example in [91, 93].

However for both variants we still need to handle an auxiliary variable. We can trade some of this computational effort for a weaker majorizer by making specific choices for $z$ in Eq. (14). To illuminate these choices we introduce the function $W_E(p, x) = E^*(p) + E(x) - \langle p, x \rangle$ [76, 12], which allows us to write

$$D_{E_\theta^*}^{x_i^*}(0, q_i) = \min_{z \in \mathbb{R}^n} W_{E_1, \theta}(-z, x_i^*) + W_{E_2, \theta}(z, x_i^*). \quad (15)$$

Note that $W_E(p, x) = 0$ if $p \in \partial E(x)$. As such choosing either $-z \in \partial E_1(x_i^*)$ or $z \in \partial E_2(x_i^*)$ allows us to simplify the problem further. This is especially attractive if $E$ is differentiable, as then both surrogates can be computed without auxiliary variables. We will denote these as *partial* surrogates, owing to the fact that we minimize only one term in (15)

**Partial Surrogate:**

$$\min_{z \in \partial E_2(x_i^*, y_i, \theta)} W_{E_1, \theta}(-z, x_i^*). \quad (16)$$

Effectively, this reduces the requirements of (14), as only the convex conjugate of $E_1$ needs to be computed. By symmetry, the other partial surrogate follows analogously.

We can finally also return to the previously discussed gradient penalty (6). If our energy $E$ is $m(\theta, y)$-strongly convex, then its convex conjugate is strongly smooth and we can bound the dual formulation (11) via

**Gradient Penalty**

$$\frac{1}{m(\theta, y_i)} ||q_i||^2 \quad \text{s.t. } q_i \in \partial E(x_i^*, y_i, \theta). \quad (17)$$

While this formulation allows us to minimize an upper bound on the bi-level problem without either auxiliary variables or knowledge about $E_1^*$ or $E_2^*$, it also is the crudest over-approximation among the considered surrogates as the following proposition illustrates.

**Proposition 2** (Ordering of parametric majorizers)**.** *Assuming the condition $l(x, z) \leq D_{E_\theta}(x, z)$ from Eq. (8), we find that the presented parametric majorizers can be ordered in the following way:*

$$l(x_i^*, x(\theta)) \leq D_{E_\theta}^0(x_i^*, x_i(\theta)) = D_{E_\theta^*}^{x_i^*}(0, q_i)$$
$$\leq \min_{z \in \partial E_2(x_i^*)} W_{E_1}(-z, x_i^*)$$
$$\leq \frac{1}{m(\theta, y)} ||q_i||^2 \quad \text{s.t. } q_i \in \partial E(x_i^*, y, \theta).$$
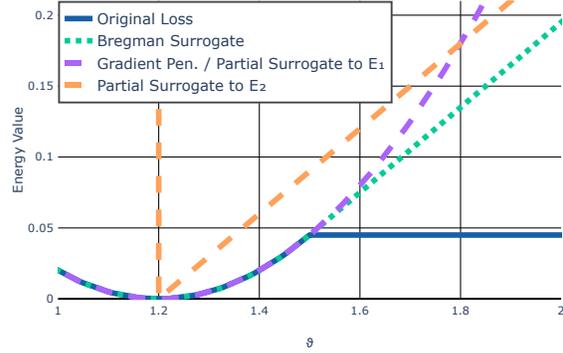


Figure 1. Visualization of surrogate functions for the bi-level problem given in Eq. (18). The blue line marks the original bi-level problem, the green dots marks the Bregman distance surrogate discussed in Eq. (10). The orange curve marks the partial surrogate obtained from (15) by inserting $z = \nabla E_1(x^*)$, whereas the purple line marks the other partial surrogate (16) which is equivalent to the gradient penalty (17) here.

*The Bregman surrogate (10) majorizes the original loss function and is in turn majorized by the partial surrogate (16) which is majorized by the gradient penalty (17) under the assumption of strong convexity.*

*Proof.* See appendix. □

As a clarifying example, we can simplify these majorizers in the differentiable setting:

**Example 2** (Differentiable Energy)**.** *Let $E$ be differentiable and $m(\theta, y)$-strongly convex, then the majorizers in Prop. 2 are given by*

$$l(x_i^*, x(\theta)) \leq D_{E_\theta}(x_i^*, x_i(\theta)) = D_{E_\theta^*}(0, \nabla E(x_i^*, y_i, \theta))$$
$$\leq W_{E_1}(-\nabla E_2(x_i^*), x_i^*)$$
$$\leq \frac{1}{m(\theta, y)} ||\nabla E(x_i^*, y_i, \theta)||^2.$$

## 3.3. Intermission: One-Dimensional Example

Let us illustrate our discussion with a toy example. We consider the non-smooth bi-level problem of learning the optimal sparsity parameter $\theta$ in the bi-level problem:

$$\min_{\theta \in \mathbb{R}} \frac{1}{2} |x^* - x(\theta)|^2, \quad (18)$$

$$\text{subject to} \quad x(\theta) = \arg\min_x \frac{1}{2}|x - y|^2 + \theta|x|. \quad (19)$$

As the lower-level energy is 1-strongly convex and the upper level loss is quadratic $l(x, y) \leq D_{E_\theta}(x, y)$ holds. Detailed derivations of all three surrogate functions of this example can be found in the appendix. Figure 1 visualizes

these surrogates, plotting their energy values relative to $\theta$. Due to the low dimensionality of the problem, all surrogate functions coincide with the original loss function at the optimal value of $\theta$. It is further interesting to note that the Bregman surrogate is exactly identical with the original loss function in the vicinity of the optimal value, due to the low dimensionality of the example.

### 3.4. Iterative Majorizers

We used subsection 3.2 to construct a series of upper bounds to facilitate a trade-off between efficiency and exactness. However what happens if we are not satisfied with the exactness of the Bregman surrogate (9)? This setting can happen especially if $x^*$ and $x(\theta)$ are significantly incompatible and subsequently $l(x^*, x(\theta))$ is large, even for optimal $\theta$. For example if we try to optimize only a few hyper-parameters we might not at all expect $x(\theta)$ to be close to $x^*$. This discussion can again be linked to the notion of 'separability' in SVM approaches [96]: The quality of the majorizing strategy is directly related to the level of 'separability' of the bi-level problem.

However, we can use the previously introduced majorizers iteratively. To do so we need to develop a majorizer that depends on a given estimate $\bar{x}$.

**Proposition 3.** *Under the standing assumption that $l(x, y) \leq D_{E_\theta}(x, y)$ (8) and if the loss function is induced by a strictly convex function $w : \mathbb{R}^n \to \mathbb{R}$, i.e. $l(x, y) = D_w(y, x)$, we have the following inequality:*

$$l(x, y) \leq l(x, z) + \langle \nabla_z l(x, z), y - z \rangle + D_E(z, y). \quad (20)$$

*Proof.* It holds that $l(x, y) = D_w(y, x)$ which is equivalent to $D_w(y, z) + D_w(z, x) - \langle \nabla w(x) - \nabla w(z), z - y \rangle$ by the Bregman 3-Point inequality [20, 94]. Using the standing assumption and that $\nabla w(x) - \nabla w(z) = \nabla_x D_w(x, z)$ we find the proposed inequality. $\square$

Assume we are given an estimated solution $\bar{x}_i$, then we can use this estimate to rewrite our bound to

$$\begin{aligned} l(x_i^*, x_i(\theta)) \leq & l(x_i^*, \bar{x}_i) + \langle \nabla l(x_i^*, \bar{x}_i), x_i(\theta) - \bar{x}_i \rangle \\ & + D_E(\bar{x}_i, x_i(\theta)). \end{aligned} \quad (21)$$

This is a linearized variant of the parametric majorization bound and as such a nonconvex composite majorizer in the sense of [39], as such a key property of majorization-minimization techniques remains in the parametrized setting, choosing $\bar{x}_i = x_i(\theta^k)$:

**Proposition 4** (Descent Lemma)**.** *The iterative procedure given by repeatedly minimizing the right-hand side of Eq. (21) in $\theta$ and setting $\bar{x}_i = x_i(\theta^k)$ is guaranteed to be stable, i.e. not to increase the bi-level loss:*

$$\sum_{i=1}^N l\left(x_i^*, x_i(\theta^{k+1})\right) \leq \sum_{i=1}^N l\left(x_i^*, x_i(\theta^k)\right) \quad (23)$$

*Proof.* See appendix. $\square$

However this algorithm cannot be applied directly, as we would still need to differentiate $x_i(\theta)$ appearing in the linearized part. Nevertheless, we can use both Fenchel's inequality $\langle p, x \rangle \leq E(x) + E^*(p)$ and the previously established $D_{E_\theta}(x, x(\theta)) = E(x, y, \theta) - E(x(\theta), y, \theta)$ to find an over-approximation to the iterative majorizer of Prop. 4:

$$\begin{aligned} & l(x_i^*, x_i(\theta)) \\ \leq & \ l(x_i^*, \bar{x}_i) - \langle \nabla l(x_i^*, \bar{x}_i), \bar{x}_i \rangle \\ & + E^* \left(\nabla l(x_i^*, \bar{x}_i), y_i, \theta\right) + E(x_i(\theta), y_i, \theta) \\ & + E(\bar{x}_i, y_i, \theta) - E(x_i(\theta), y_i, \theta) \\ = & \ l(x_i^*, \bar{x}_i) - \langle \nabla l(x_i^*, \bar{x}_i), \bar{x}_i \rangle \\ & + E(\bar{x}_i, y_i, \theta) + E^* \left(\nabla l(x_i^*, \bar{x}_i), y_i, \theta\right) \end{aligned}$$

This estimate reveals that we can approximate the iterative majorizer much like the previously discussed surrogates:

---

**Iterative Surrogate**

$$E(\bar{x}_i, y, \theta) + E^* \left(\nabla l(x_i^*, \bar{x}_i), y_i, \theta\right) + C, \quad (22)$$

---

as the constant $C = l(x_i^*, \bar{x}_i) - \langle \nabla l(x_i^*, \bar{x}_i), \bar{x}_i \rangle$ does not depend on $\theta$. We essentially return to Eq. (12) and only the input to $E$ and $E^*$ changes with respect to $\bar{x}_i$. This strategy recovers the previous majorizer as a special case:

**Corollary 1.** *If we linearize around $\bar{x}_i = x_i^*$, then we recover the Bregman surrogate of (9).*

*Proof.* If $\bar{x}_i = x_i^*$, then $l(x_i^*, \bar{x}_i) = 0$ and $\nabla l(x_i^*, \bar{x}_i) = 0$ by the properties of the differentiable loss function. As such the constant term $C$ is zero and $E^* \left(\nabla l(x_i^*, \bar{x}_i), y_i, \theta\right) = E^*(0, y_i, \theta)$ so that we recover (12) which is equivalent to the Bregman surrogate (9). $\square$

We can use this surrogate to form an efficient approximation to a classical majorization-minimization strategy as in [88, 67, 66, 48]. Notably the 'tightness' of the majorization is violated by the over-approximation, i.e. inserting $\theta^k$ into the majorizer does not recover $l(x_i^*, x_i(\theta^k))$. We iterate

$$\begin{aligned} \theta^{k+1} = \arg\min_\theta \sum_{i=1}^N & E^* \left(\nabla l(x_i^*, x_i(\theta^k)), y_i, \theta\right) \\ & + E\left(x(\theta^k), y_i, \theta\right) \end{aligned} \quad (23)$$

As the application of this iterative scheme reduces to a simple change from Eq (12) to Eq. (22), we can easily apply it in practice to further increase the fidelity of the surrogate by solving a sequence of fast surrogate optimizations. We initialize the scheme with $\bar{x}_i = x_i^*$ as suggested from Corollary 1 and either stop iterating or reduce the step size of the surrogate solver if the higher-level objective is increased after an iteration.

## 4. Examples

This section will feature several experiments[1] in which we will illustrate the application of the investigated methods. We will show two concepts of new applications that are possible in parametrized variational settings, 4.1 and 4.2. We then show an application to image denoising in 4.3.

### 4.1. Computed Tomography

Making only specific parts of a variational model learnable is especially interesting for computed tomography (CT). An image $x$ is to be reconstructed from data $y = Ax + n$ that is formed by applying the radon transform to the image $x$ and adding noise $n$. While first fully-learning based solutions to this problem exist (e.g. [50, 51]), suitable networks are difficult to find not only due to the ill-posedness of the underlying problem, but also due to the well-justified concerns about fully learning-based approaches in medical imaging [3]. To benefit from the explicit control of the data fidelity of the reconstruction, we consider to introduce a learnable linear correction term into an otherwise classical reconstruction technique via

$$x_i(\theta) = \arg\min_x \frac{1}{2}\|Ax - y_i\|_2^2 + \beta R(x) + \langle x, \mathcal{N}(\theta, y_i)\rangle,$$

for a suitable network $\mathcal{N}$ (we chose 8 blocks of $3 \times 3$ convolutions with 32 filters, ReLU activations, and batch-normalization, and a final $5 \times 5$ convolution), and $R$ denoting the Huber loss of the discrete gradient of $x$.

As both convex conjugates are difficult to evaluate in closed-form, we choose the gradient penalty (17), which is a parametric majorizer for euclidean loss if $A$ has full rank (and practically even works beyond this setting, as it majorizes $\|A(x-y)\|^2$ even for rank-deficient $A$). According to (17) we consider

$$\min_{\theta \in \mathbb{R}^s} \sum_{i=1}^n \|A^*Ax_i^* - A^*y_i + \beta\nabla R(x_i^*) + \mathcal{N}(\theta, y_i)\|_2^2,$$

train on simulated noisy data and test our model on the widely-used Shepp-Logan phantom. Figure 2 illustrates the resulting reconstruction, as well as the best reconstruction using the variational approach without the additional linear correction term after a grid-search for the optimal $\beta$. As we can see, the surrogate trained the linear correction term well enough to improve the PSNR of the reconstruction by almost 2dB. Moreover, the influence of the linear correction term can still be visualized and the data fidelity can easily be controlled via a suitable weighting. We visualize the correction map in the appendix.



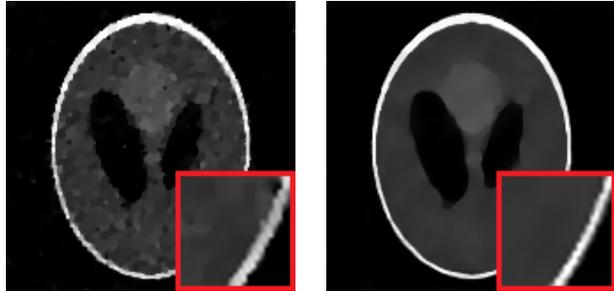Huber-TV, PSNR 23.9          Learned cor., PSNR 25.8

Figure 2. Learning a linear correction term for a Huber-regularized CT reconstruction problem using the gradient penalty (6).

### 4.2. Variational Segmentation

For a very different (and non-smooth) example, consider the task of learning a variational segmentation model [18, 15, 32, 72]. We are interested in learning a model whose minimizer coincides with a (semantic) segmentation of the input data. The lower-level problem is given by

$$x(\theta) = \arg\min_x -\langle \mathcal{N}(\theta, y), x\rangle + \|Dx\|_1 + h(x), \quad (24)$$

where $h(x) = \sum_{j=1}^n x_i \log(x_i) + I_\Delta(x)$ is the entropy function on the unit simplex $\Delta$ [7]. $\mathcal{N}(\theta, y)$ is some parametrized function that computes the potential of the segmentation model, this can be a deep neural network, as we only require convexity in $x$ and not in $\theta$. $D$ is a finite-differences operator, so that the overall total variation (TV) term $\|Dx\|_1$ measures the perimeter of a segmentation $x$ if $x \in \{0,1\}^n$. The entropy function crucially not only leads to a strictly convex model but also represents the structure of a usual learned segmentation method. Without the perimeter term, a solution to the lower-level problem would be given by

$$x(\theta) = \nabla h^*(\mathcal{N}(\theta, y)). \quad (25)$$

Due to [79, P.148], $\nabla h^*$ is exactly the $\mathrm{softmax}$ function, so that Eq. (25) is equivalent to applying a parametrized function $\mathcal{N}$ and then applying the $\mathrm{softmax}$ function to arrive at the final output, a usual image recognition pipeline during training. As a higher-level loss, we choose $\log$ loss

$$\sum_{i=1}^N -\langle x_i^*, \log(x_i(\theta))\rangle = \sum_{i=1}^N D_h(x_i^*, x_i(\theta)) \quad (26)$$

so that the bi-level problem without the perimeter term is equivalent to minimizing the cross-entropy loss of $\mathcal{N}(\theta, y)$. With the inclusion of the perimeter term, however, we cannot find a closed-form solution for $x(\theta)$ need to consider bi-level optimization. But, as the log-loss (26) can be written as a Bregman distance relative to $h$, our primary assumption $l(x, z) \leq D_{E_\theta}(x, z)$ (8) is fulfilled and we can consider the
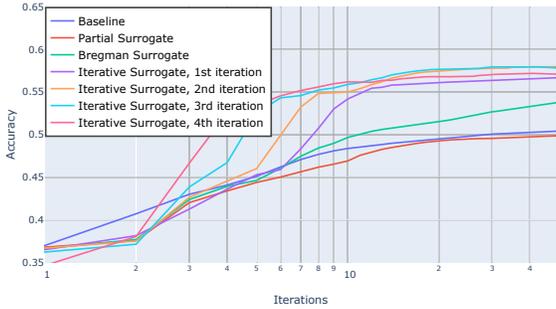
---

Figure 3. Training accuracy for the variational segmentation model discussed in Section 4.2 for a linear model $\mathcal{N}(\theta, y_i)$. Directly training a cross-entropy loss without the perimeter term, training the Bregman surrogate Eq (28), the Partial surrogate Eq (16) and four iterations of the iterative scheme are compared. We find that the end-to-end training with the perimeter term increases the segmentation accuracy. We also see that a small number of iterations in the iterative scheme is sufficient for a practical CV task.

Bregman surrogate problem in the dual setting of Eq. (14):

$$\min_{\theta} \sum_{i=1}^{N} \min_{z_i} W_h(\mathcal{N}(\theta, y_i) - z_i, x_i^*) + W_{TV}(z_i, x_i^*), \quad (27)$$

which we can rewrite to

$$\min_{\theta} \sum_{i=1}^{N} \min_{||p_i|| \leq 1} h^* \left( \mathcal{N}(\theta, y_i) - D^T p_i \right) \\ - \langle \mathcal{N}(\theta, y_i), x_i^* \rangle + ||Dx_i^*||_1. \quad (28)$$

We note that this is essentially a cross-entropy loss with an additional additive term $p_i$, that is able to balance out incoherent output of $\mathcal{N}(\theta, y_i)$ that would lead to erroneous segmentations with a higher perimeter. Furthermore, the training process is still convex w.r.t to $\mathcal{N}(\theta, y_i)$, in contrast to unrolling schemes. The iterative model (23) has a very similar structure, including the gradient of the loss into (28).

To validate this setup, we choose $\mathcal{N}$ to be given by a simple convolutional linear model. We draw a small subset of the `cityscapes` dataset and compare the cross entropy model of Eq (25) with the total variation bi-level model of Eq. (28) and its partial and iterative applications. Figure 3 visualizes the training accuracy over training iterations. We find that the proposed approach is able to improve the segmentation accuracy of the linear model significantly. We refer to the appendix for further details.

### 4.3. Analysis Operator Models

Finally, we illustrate the behaviour of our approach on a practically relevant model, learning a set of optimal convolutional filters for denoising [83, 26]. We consider the

| Model | PSNR | T | PSNR(Iter.) | TT |
|---|---|---|---|---|
| Total Variation | 27.41 | - | - | - |
| 3 3x3 Filters | 26.66 | 00:34 | 27.66 | 02:21 |
| 48 7x7 Filters | 27.41 | 02:45 | 28.03 | 03:11 |
| 96 9x9 Filters | 27.46 | 01:43 | 28.03 | 02:22 |

Table 1. Training time (T) in minutes for each surrogate computation and PSNR on the test dataset for various gray-scale filters for the energy model in Eq. (30) with and without the iterative process of Eq (22) and total time (TT) for the iterative process are compared to total variation with optimal regularization parameter. Note that training time varies mostly due to differing iteration counts. The results of the convex model of [26] are reproduced.

parametric energy model

$$x(\theta) = \arg\min_x \frac{1}{2}||x - y_i||^2 + ||D(\theta)x||_1, \quad (29)$$

with $D(\theta)$ denoting the convolution operator to be learned, which is prototypical for many other image processing tasks. We consider square loss $l(x, y) = \frac{1}{2}||x - y||^2$ as a higher loss function and apply our approach. A Bregman surrogate for this model has the form

$$\min_{\theta} \sum_{i=1}^{N} \min_{||p_i|| \leq 1} ||D(\theta)x_i^*||_1 + \frac{1}{2}||D^T(\theta)p_i - y_i||^2. \quad (30)$$

Model (29) was previously considered in [26, 24], where it was solved via implicit differentiation. We repeat the setup of [26] and train a denoising model on the `BSDS` dataset [68]. Refer to the appendix for the experimental setup and optimization strategy.

Table 1 shows both PSNR values achieved when training $D(\theta)$ as convolutional filters as well as training time. In comparison to [26], we find strikingly, that we can train a convex model with similar performance to the convex model in [26], while being an order of magnitude faster than the original approach. Furthermore in [26], the necessary training time jumps from 24 hours for 48 7x7 filters to 20 days for 96 9x9 filters - in our experiment the training time is almost unaffected by the number of parameters, and in this example actually smaller as the larger model converges faster. Also this analysis validates that the iterative process is crucial to reaching competitive PSNR values.

## 5. Conclusions

We investigated approximate training strategies for data-driven energy minimization methods by introducing *parametric majorizers*. We systematically studied such strategies in the framework of convex analysis, and proposed the Bregman distance induced by the lower level energy as well as over-approximations thereof as suitable majorizers. We discussed an iterative scheme that shows promise for applications in computer vision, particularly due to its scalability as shown by its application to image denoising.

# A. Appendix

## A.1. Convex Analysis in Section 3

### A.1.1  Details for Derivation of Eqs. (11), (12)

Eq. (11) above describes the application of Bregman duality:

$$D^0_{E_\theta}(x_i^*, x_i(\theta)) = D^{x_i^*}_{E_\theta^*}(0, q_i) \quad q_i \in \partial E(x_i^*, y_i, \theta), \quad (11)$$

which is a common application of the following identity [11, 10]:

**Lemma 1** (Bregman Identity). *Consider a convex lsc. function $E : \mathbb{R}^n \to \mathbb{R}$ with a subgradient $p \in \partial E(y)$. Then, the following identity holds:*

$$D^p_E(x, y) = D^x_{E^*}(p, q), \quad q \in \partial E(x)$$

*Proof.* This property follows from equality (Fenchel's identity) in the Fenchel-Young inequality $E(x) + E^*(p) = \langle p, x \rangle \iff p \in \partial E(x)$. To see this we write

$$D^p_E(x, y) = E(x) - \langle p, x \rangle - E(y) + \langle p, y \rangle$$

and apply Fenchel's identity for $p, y$ to find

$$D^p_E(x, y) = E(x) - \langle p, x \rangle + E^*(p)$$

We then introduce any $q \in \partial E(x)$ by writing $\langle p, x \rangle = \langle p - q + q, x \rangle$ and apply Fenchel's identity again:

$$D^p_E(x, y) = E^*(p) - E^*(q) - \langle x, p - q \rangle = D^x_{E^*}(p, q)$$

$\square$

The step from Eq. (11) to Eq.(12) is simply the first step of this derivation:

$$
\begin{aligned}
D_{E_\theta}(x_i^*, x_i(\theta)) &= E(x_i^*, y_i, \theta) - \langle 0, x_i^* \rangle + E^*(0, y_i, \theta) \\
= D^{x_i^*}_{E_\theta^*}(0, q_i) &= E(x_i^*, y_i, \theta) + E^*(0, y_i, \theta) \quad (12)
\end{aligned}
$$

as $p_i = 0$ is a subgradient of $E$ at $x_i(\theta)$ and $q_i$ at $x_i^*$.

### A.1.2  Details for Derivation of Eq. (14) to (15)

A crucial subtlety of Lemma 1 is that this identity holds for any $q \in \partial E(x)$ and the choice of subgradients is irrelevant, the Bregman distance is equal for all choices. This motivates the introduction of the $W$-function $W_E(p, x) = E^*(p) + E(x) - \langle p, x \rangle$. This function is convex in either $p$ or $x$ and always non-negative. It can be understood as measuring the deviation of $p$ from subgradients of $x$ as a direct implementation of the Fenchel-Young inequality. As such it is 0 exactly if $p \in \partial E(x)$. Previous usage of this function can be found for example in [12, 76]. For Legendre functions [5], i.e. functions where both $E$ and $E^*$ are

(essentially) smooth, the connection to Bregman distances is immediate:

$$W_E(p, x) = D^p_E(x, \nabla E^*(p)),$$

for non-smooth functions this is also a part of the proof of Lemma 1, replacing $\nabla E^*(p)$ by $y \in \partial E^*(p)$. As such, we can write Eq. (12) as

$$D^{x_i^*}_{E^*}(0, q_i) = W_{E_\theta}(0, x_i^*). \quad (12)$$

The introduction of this function then allows us to show that

$$W_E(0, x_i^*) = \min_z W_{E_1, \theta}(-z, x_i^*) + W_{E_2, \theta}(z, x_i^*) \quad (15)$$

under the assumption in Eq.(13), that $E$ can be written as $E_1 + E_2$, with both functions convex. We recognize this as the clear extension of the infimal convolution property $E^*(0) = \min_z E_1^*(-z) + E_2^*(z)$ (which itself can be understood as Fenchel's duality theorem applied to $E_1, E_2$) to these functions, in the smooth setting this could be written via

$$
\begin{aligned}
D^{x_i^*}_{E^*}(0, \nabla E(x_i^*)) = \min_z \quad & D_{E_1^*}(-z, \nabla E_1(x_i^*)) \\
& + D_{E_2^*}(z, \nabla E_2^*(x_i^*)).
\end{aligned}
$$

We arrive at Eq. (15) from Eq. (14) by rewriting $E$ in Eq.(14):

$$
\begin{aligned}
\min_z \quad & E_1(x_i^*, y_i, \theta) + E_2(x_i^*, y_i, \theta) \\
& + E_1^*(-z, y_i, \theta) + E_2^*(z, y_i, \theta) \\
= \min_z \quad & E_1(x_i^*, y_i, \theta) + E_2(x_i^*, y_i, \theta) + \langle z, x_i^* \rangle \\
& + E_1^*(-z, y_i, \theta) + E_2^*(z, y_i, \theta) - \langle z, x_i^* \rangle \\
= \min_z \quad & W_{E_1, \theta}(-z, x_i^*) + W_{E_2, \theta}(z, x_i^*). \quad (15)
\end{aligned}
$$
(14)

### A.1.3  Proof of Proposition 2

**Proposition 2** (Ordering of parametric majorizers). *Assuming the condition $l(x, z) \leq D_{E_\theta}(x, z)$ from Eq. (8), we find that the presented parametric majorizers can be ordered in the following way:*

$$
\begin{aligned}
l(x_i^*, x(\theta)) &\leq D^0_{E_\theta}(x_i^*, x_i(\theta)) = D^{x_i^*}_{E_\theta^*}(0, q_i) \\
&\leq \min_{z \in \partial E_2(x_i^*)} W_{E_1}(-z, x_i^*) \\
&\leq \frac{1}{m(\theta, y)} \|q_i\|^2 \quad \text{s.t. } q_i \in \partial E(x_i^*, y, \theta).
\end{aligned}
$$

*The Bregman surrogate (10) majorizes the original loss function and is in turn majorized by the partial surrogate (16) which is majorized by the gradient penalty (17) under the assumption of $m(\theta, y)$ - strong convexity of $E_1$.*
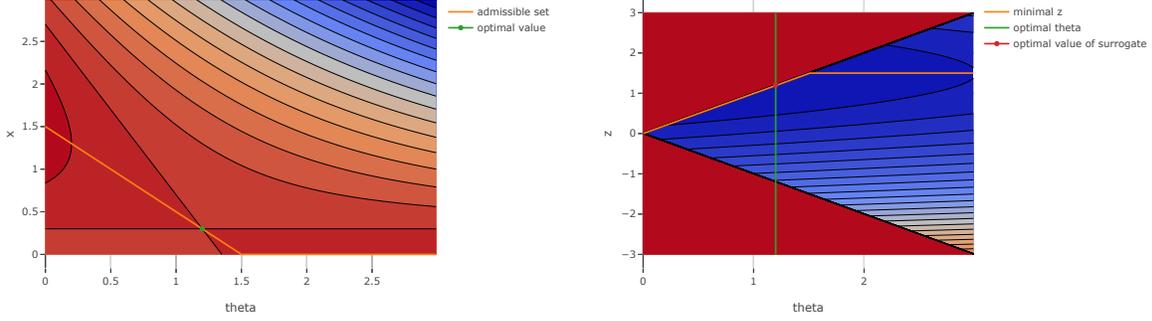
9

Figure 4. Visualization of the Bregman surrogate problem in primal formulation (left) and dual formulation (right). The problem in visualized over all $(x, \theta)$, respectively $(z, \theta)$. The admissible $x(\theta)$ are marked in orange in the left contour plot and the optimal $z(\theta)$ one the right. The optimal value in $\theta$ is marked in green in both plots.

*Proof.* The first inequality follows directly by the assumption $l(x, z) \leq D_{E_\theta}(x, z)$. The second inequality is the application of Bregman Duality discussed in Lemma 1. From Eq.(15) we now see that $D_{E_\theta}^{x_i^*}(0, q_i)$, $q_i \in \partial E(x_i^*, y_i, \theta)$ can be written as a minimum over $z$. Clearly choosing a non-optimal $z$ yields an upper bound to this minimal value. Without loss of generality, we choose $z \in \partial E_2(x_i^*)$ so that $W_{E_2, \theta}(z, x_i^*)$ is equal to zero.

Now we assume that $E$ is $m(\theta, y)$-strongly convex. We subsume this strong convexity term in $E_1$ again without loss of generality so that $E_1$ is strongly convex. By convex duality [6], this implies that $E_1^*$ is $m(\theta, y)$ strongly smooth, i.e. $D_{E_1^*}^x(p, q) \leq \frac{1}{2m(\theta, y)}||p - q||^2$. Following Eq.(12), we write

$$W_{E_1^*}(-z, x_i^*) = D_{E_1^*}^{x_i^*}(-z, r) \quad z \in \partial E_2(x_i^*, y_i, \theta),$$
$$r \in \partial E_1(x_i^*, y_i, \theta)$$
$$\leq \frac{1}{2m(\theta, y)}|| -z - r||^2$$
$$= \frac{1}{2m(\theta, y)}||q_i||^2 \qquad q_i \in \partial E(x_i^*, y_i, \theta),$$

under mild assumptions on the additivity of subgradients of $E_1$ and $E_2$. $\qquad \square$

### A.1.4 Derivation of the surrogate functions for the example in subsection 3.3

Section 3.3 discusses the non-smooth bi-level problem given in Eqs. (18) and (19):

$$\min_{\theta \in \mathbb{R}} \frac{1}{2}|x^* - x(\theta)|^2, \tag{18}$$

$$\text{subject to} \quad x(\theta) = \arg\min_x \frac{1}{2}|x - y|^2 + \theta|x|. \tag{19}$$

for both $x^*, y \in \mathbb{R}$. In this setting, the 'primal' formulation of the Bregman surrogate is given by

$$\min_\theta \max_x \frac{1}{2}|x^* - y|^2 - \frac{1}{2}|x - y|^2 + \theta(|x^*| - |x|) \tag{10 ex.}$$

whereas the 'dual' formulation is given by

$$\min_\theta \min_{|z| \leq \theta} \frac{1}{2}|x^* - y|^2 + \theta|x^*| + \frac{1}{2}|z - y|^2. \tag{12 ex.}$$

Note that this problem is convex in $z, \theta$ as the epigraph constraint $|z| \leq \theta$ is convex. Both (equivalent!) variants are visualized in Figure 4. We see that the saddle-point of the primal formulation and the minimizer of the dual formulation correctly coincide with the optimal $\theta$.

Moving forward, we set $E_1(x, y) = \frac{1}{2}|x - y|^2$ and $E_2(x, \theta) = \theta|x|$ to compute the two partial surrogates. Firstly $W_{E_1, \theta}(-z, x^*)$, $z \in \partial E_2(x^*)$ leads to

$$\min_\theta \frac{1}{2}|x^* - y + q|^2, \quad q \in \partial|x^*|, \tag{16 ex.1}$$

where we take $q = \text{sign}(x^*)$ as $x^* \neq 0$ in our example. As $E_1$ is a quadratic function, this is also equivalent to the gradient penalty in Eq. (17). The second partial surrogate, $W_{E_2, \theta}(z, x^*)$, $z \in \partial E_1(x^*)$ can be written as

$$\min_\theta \theta|x^*| + I_{|\cdot| \leq \theta}(x^* - y) - \langle x^*, x^* - y \rangle \tag{16 ex.2}$$
$$= \min_{|x^* - y| \leq \theta} \theta|x^*| + C.$$

Figure 4 here and Figure 1 in the main paper both arise from the data point $x^* = 0.3, y = 1.5$.

To give some more details on the fact that the Bregman surrogate is exactly identical with the original loss function in the vicinity of the optimal value, note that this is caused by the special structure of the Bregman distance of

10

the absolute value, $D_{|\cdot|}(x,y)$ as $D_{E_\theta}(x,y)$ decomposes into $\frac{1}{2}|x-y|^2 + \theta D_{|\cdot|}(x,y)$. This function is equal to the higher-level loss function as soon as the signs of $x^*$ and $x(\theta)$ coincide and as such the majorizer is exact, even if it is much easier to compute.

### A.1.5  Proof of Proposition 4

Section 3.4 describes an iterative procedure for repeated application of the majorization strategies discussed in section 3.2. This scheme was based on the result of Proposition 3:

$$l(x,y) \leq l(x,z) + \langle \nabla_z l(x,z), y-z \rangle + D_E(z,y), \quad (20)$$

inserting $x = x_i^*, y = x_i(\theta), z = x_i(\theta^k)$ leads to

$$l(x_i^*, x_i(\theta)) \leq l(x_i^*, x_i(\theta^k)) + D_{E_\theta}(x_i(\theta^k), x_i(\theta)) \\ + \langle \nabla l(x_i^*, x_i(\theta^k)), x_i(\theta) - x_i(\theta^k) \rangle. \quad (20b)$$

Eq.(20), respectively (20b), lead to a monotone descent of the higher-level loss, as shown in Proposition 4:

**Proposition 4** (Descent Lemma). *The iterative procedure given by*

$$\theta^{k+1} = \arg\min_\theta \sum_{i=1}^N l(x_i^*, x_i(\theta^k)) \\ + \langle \nabla l(x_i^*, x_i(\theta^k)), x_i(\theta) - x_i(\theta^k) \rangle \\ + D_{E_\theta}^0(x_i(\theta^k), x_i(\theta))$$

*is guaranteed to be stable, i.e. not to increase the bi-level loss:*

$$\sum_{i=1}^N l\left(x_i^*, x_i(\theta^{k+1})\right) \leq \sum_{i=1}^N l\left(x_i^*, x_i(\theta^k)\right) \quad (23)$$

*Proof of Proposition 4.* $\theta^{k+1}$ is a minimizer of the iterative scheme. Therefore, evaluating the iteration at $\theta^{k+1}$ leads to a lower value than evaluating at $\theta^k$:

$$\sum_{i=1}^N l(x_i^*, x_i(\theta^k)) + \langle \nabla l(x_i^*, x_i(\theta^k)), x_i(\theta^{k+1}) - x_i(\theta^k) \rangle \\ + D_{E_{\theta^{k+1}}}^0(x_i(\theta^k), x_i(\theta^{k+1})) \\ \leq \sum_{i=1}^N l(x_i^*, x_i(\theta^k)) + \langle \nabla l(x_i^*, x_i(\theta^k)), x_i(\theta^k) - x_i(\theta^k) \rangle \\ + D_{E_{\theta^k}}^0(x_i(\theta^k), x_i(\theta^k)) \\ = \sum_{i=1}^N l(x_i^*, x_i(\theta^k))$$

Now the left-hand-side is also equivalent to Eq. (20b) evaluated at $\theta^{k+1}$. Applying the inequality in (20b) for all

$i = 1, \ldots, N$ we find

$$\sum_{i=1}^N l(x_i^*, x_i(\theta^{k+1})) \leq \sum_{i=1}^N l(x_i^*, x_i(\theta^k)).$$

$\square$

*Remark.* The iterative scheme given in Eq.(22), i.e.

$$\theta^{k+1} = \arg\min_\theta \sum_{i=1}^N E^* \left(\nabla l(x_i^*, x_i(\theta^k)), y_i, \theta\right) \\ + E\left(x(\theta^k), y_i, \theta\right). \quad (22)$$

is an over-approximation of the iterative scheme discussed in Proposition 4. As such we expect the results of Proposition 4 to hold only approximately as stated in the main paper.

### A.2. Experimental Setup

This section will add additional details to the experiments presented in the paper[2].

### A.2.1  CT - Additional Details

The implementation of the CT example in section 4.1 is straightforward. We generate pairs $(y_i^*, x_i^*)$ of noisy sinograms and ground truth images and optimize

$$\min_{\theta \in \mathbb{R}^p} \sum_{i=1}^n \|A^*Ax_i^* - A^*y_i + \beta \nabla R(x_i^*) + \mathcal{N}(\theta, y_i)\|_2^2.$$

We test our model on the widely-used Shepp-Logan phantom, comparing the learned model with a pure Huber-TV solution, for which we found the optimal parameter $\beta$ by grid search. This setup was implemented in Matlab. To visualize the linear correction term, we repeat an extended version of Figure 2 in Figure 5.

### A.2.2  Segmentation - Additional Details

The segmentation experiment shown in Figure 3 of the main paper shows the results of training the variational model in Eq.(25), which corresponds to an augmented cross-entropy term, as discussed in section 4.2.

The partial surrogate implemented in Figure 3 is a direct application of Eq.(16) to the segmentation setting, giving

$$\min_\theta \sum_{i=1}^N \min_{p_i \in \partial \|Dx_i^*\|} D_h\left(x_i^*, \nabla h^*\left(\mathcal{N}(\theta, y_i) - D^T p_i\right)\right),$$

---

[2]Refer also to the implementations hosted at https://github.com/JonasGeiping/ParametricMajorization

11

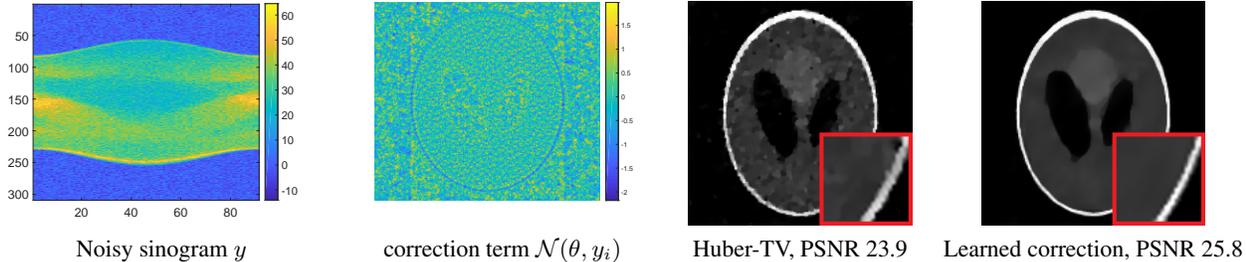| Noisy sinogram $y$ | correction term $\mathcal{N}(\theta, y_i)$ | Huber-TV, PSNR 23.9 | Learned correction, PSNR 25.8 |

Figure 5. Illustrate our results for learning a linear correction term for a Huber-regularized CT reconstruction problem. In reference to Figure 2 in the main paper we also visualize input data and the learned linear correction map. The predicted linear correction term can be visualized and inspected, and its influence can easily be quantified or explicitly scaled via a parameter.

where the computation of the auxiliary variable $p_i$ is simplified. Note further that the gradient penalty cannot be applied in this setting, as the segmentation energy $E$ is not strongly convex. Similarly, the iterative approach can be computed to be

$$
\min_\theta \sum_{i=1}^{N} \min_{||p_i|| \leq 1} h^* \left( \frac{x_i^*}{x_i(\theta^k)} + \mathcal{N}(\theta, y_i) - D^T p_i \right) \\
- \langle \mathcal{N}(\theta, y_i), x_i(\theta^k) \rangle
$$

which is still convex in $\mathcal{N}(\theta, y)$, but the input arguments now take previous solutions into account.

To emphasize the convexity of the setup, we choose $\mathcal{N}(\theta, y_i)$ as a linear convolutional network of $3x3x3$ filters for each target class. We accordingly optimize the resulting convex minimization problems by an optimal convex optimization method, namely FISTA [8]. To solve the inference problem in Eq. (25) we apply usual strategies and optimize via a primal-dual algorithm [16] - to increase the speed we adapt a recent variant [17] and consider the Bregman-Proximal operator in the primal sub-problem for which we use the entropy function $h$ described in the paper, paralleling [7, 74].

We draw four images and their corresponding segmentations from the `cityscapes` data set [31] and implement the proposed procedures in PyTorch [75]. For Figure 3 we drew the first four images, which we resized to 128x256 pixels. To visualize the improvement over the iterations, we initialize the subsequent iterations of the iterative scheme again with the initial value of $\theta$, so that the training accuracy curves in Figure 3 are comparable. This is of course not strictly necessary and $\theta$ could be initialized with the current estimate in every iteration. We also point out that we visualize the actual training accuracy in Figure 3, meaning the percentage of successfully segmented pixels after *hard argmax* of the results of the algorithms.

### A.2.3 Analysis Operators - Additional Details

For this experiment we considered the task of learning an 'analysis operator' $D(\theta)$, i.e. a set of convolutional filters

$\theta^k$ so that $D(\theta) = \sum_{k=1}^{K} \theta_k * x$ for a set of $K$ filters. Due to anisotropy, we can write the resulting minimization problem as

$$
x(\theta) = \arg\min_x \frac{1}{2} ||x - y||^2 + \sum_{k=1}^{K} ||\theta_k * x||_1.
$$

We repeat the experimental setup of [26] and train this model on image pairs $x^*, y$ of noise-free and noisy image patches, to learn filters that result in a convex denoising model [25, 26]. To do so we draw a batch of 200 $64x64$ image patches from the training set of the Berkeley Segmentation data set [68], convert the images to gray-scale and add Gaussian noise. To compare with [26] and [99] we do not clip the noisy images and use Matlab's `rgb2gray` routine to generate this data. Further, as in [26], we do not optimize directly for the convolutional filters, but instead decompose each filter into a DCT-II basis, where we learn the weight of each basis function, excluding the constant basis function [47]. Before training we initialize these weights by orthogonal initialization [86] with a factor of $0.01$, respectively $0.001$ for the larger 9x9 filters.

To solve the training problem we minimize Eq. (33) in the paper jointly in $\theta, \{p_i\}_{i=1}^N$. We do this efficiently by taking steps toward the optimal weights with the 'Adam' optimization procedure [52] with a step size $\tau = 0.1$ (although gradient descent with momentum or FISTA [8] are also valid options). We use a standard accelerated primal-dual algorithm [16] to solve the convex inference problem. For the iterative procedure we repeat this process, computing $x(\theta^k)$ after every minimization of Eq.(33), inserting it as a factor into $E^*$ and repeating the optimization. If the iterative procedure increases the loss value, we reduce the step size $\tau$ of the majorizing problem and repeat the step. If reducing the step size does not successfully improve the result for several iterations, we terminate the algorithm.

We implement this setup in PyTorch [75] and refer to our reference implementation for further details.

For total variation denoising, which corresponds to choosing $D(\theta)$ as the gradient operator with appropriate scaling, $\alpha \nabla$, we use grid search to find the optimal scaling

parameter $\alpha$.

We report execution times for a single minimization of Eq.(33) for different filter sizes in Table 1 in the paper as well as total time for an iterative procedure. These timings are reported for a single *GeForce RTX 2080Ti* graphics card.

# References

[1] Yasemin Altun, Ioannis Tsochantaridis, and Thomas Hofmann. Hidden Markov Support Vector Machines. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, ICML'03, pages 3–10. AAAI Press, 2003. 2

[2] Brandon Amos and J. Zico Kolter. OptNet: Differentiable Optimization as a Layer in Neural Networks. In *International Conference on Machine Learning*, pages 136–145, July 2017. 2

[3] Vegard Antun, Francesco Renna, Clarice Poon, Ben Adcock, and Anders C. Hansen. On instabilities of deep learning in image reconstruction - Does AI come at a cost? *arXiv:1902.05300 [cs]*, Feb. 2019. 1, 7

[4] Adrian Barbu. Learning real-time MRF inference for image denoising. *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1574–1581, 2009. 2

[5] Heinz H Bauschke and Jonathan (Jon Borwein. Legendre Functions and the Method of Random Bregman Projections. *Journal of Convex Analysis*, 4(1):27–67, May 1997. 9

[6] Heinz H. Bauschke and Patrick L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. CMS Books in Mathematics. Springer New York, New York, NY, 2011. 3, 10

[7] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, May 2003. 7, 12

[8] Amir Beck and Marc Teboulle. A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, Jan. 2009. 12

[9] Kristin P. Bennett, Gautam Kunapuli, Jing Hu, and Jong-Shi Pang. Bilevel Optimization and Machine Learning. In *IEEE World Congress on Computational Intelligence, WCCI 2008*, Lecture Notes in Computer Science, pages 25–47. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008. 2

[10] Martin Benning and Martin Burger. Modern regularization methods for inverse problems. *Acta Numerica*, 27:1–111, May 2018. 4, 9

[11] Martin Burger. Bregman Distances in Inverse Problems and Partial Differential Equations. In *Advances in Mathematical Modeling, Optimization and Optimal Control*, Springer Optimization and Its Applications, pages 3–33. Springer International Publishing, Cham, 2016. 3, 9

[12] Dan Butnariu and Gabor Kassay. A Proximal-Projection Method for Finding Zeros of Set-Valued Operators. *SIAM J. Control Optim.*, 47(4):2096–2136, Jan. 2008. 5, 9

[13] Luca Calatroni, Chung Cao, Juan Carlos De Los Reyes, Carola-Bibiane Schönlieb, and Tuomo Valkonen. Bilevel approaches for learning of variational imaging models. *Variational Methods: In Imaging and Geometric Control*, 18:252, 2017. 2

[14] Antonin Chambolle, Vicent Caselles, Daniel Cremers, Matteo Novaga, and Thomas Pock. An introduction to total variation for image analysis. *Theoretical foundations and numerical methods for sparse recovery*, 9(263-340):227, 2010. 1, 2

[15] Antonin Chambolle, Daniel Cremers, and Thomas Pock. A Convex Approach to Minimal Partitions. *SIAM Journal on Imaging Sciences*, 5(4):1113–1158, Oct. 2012. 7

[16] Antonin Chambolle and Thomas Pock. A First-Order Primal-Dual Algorithm for Convex Problems with Applications to Imaging. *J Math Imaging Vis*, 40(1):120–145, May 2011. 12

[17] Antonin Chambolle and Thomas Pock. On the ergodic convergence rates of a first-order primal–dual algorithm. *Mathematical Programming*, 159(1-2):253–287, Sept. 2016. 12

[18] Tony F. Chan and Luminita A. Vese. Active contours without edges. *IEEE Transactions on image processing*, 10(2):266–277, 2001. 1, 7

[19] Siddhartha Chandra and Iasonas Kokkinos. Fast, Exact and Multi-scale Inference for Semantic Image Segmentation with Deep Gaussian CRFs. In *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science, pages 402–418. Springer International Publishing, 2016. 2

[20] Gong Chen and Marc Teboulle. Convergence Analysis of a Proximal-Like Minimization Algorithm Using Bregman Functions. *SIAM J. Optim.*, 3(3):538–543, Aug. 1993. 6

[21] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs. In *International Conference on Learning Representations (ICLR)*, 2015. 1

[22] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. *arXiv:1606.00915 [cs]*, June 2016. 2

[23] Yunjin Chen and Thomas Pock. Trainable Nonlinear Reaction Diffusion: A Flexible Framework for Fast and Effective Image Restoration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1256–1272, June 2017. 2

[24] Yunjin Chen, Thomas Pock, and Horst Bischof. Learning l1-based analysis and synthesis sparsity priors using bilevel optimization. In *Neural Information Processing Systems Conference (NIPS) 2012*, 2012. 2, 8

[25] Yunjin Chen, Rene Ranftl, and Thomas Pock. A bi-level view of inpainting - based image compression. In *Computer Vision Winter Workshop.* ., 2014. 2, 12

[26] Yunjin Chen, René Ranftl, and Thomas Pock. Insights Into Analysis Operator Learning: From Patch-Based Sparse Models to Higher Order MRFs. *IEEE Transactions on Image Processing*, 23(3):1060–1072, Mar. 2014. 1, 2, 8, 12

[27] Yunjin Chen, Wei Yu, and Thomas Pock. On Learning Optimized Reaction Diffusion Processes for Effective Image Restoration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5261–5269, 2015. 2

[28] Michael Collins. Discriminative Training Methods for Hidden Markov Models: Theory and Experiments with Perceptron Algorithms. In *Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - Volume 10*, EMNLP '02, pages 1–8, Stroudsburg, PA, USA, 2002. Association for Computational Linguistics. 2

[29] Aleksander Colovic, Patrick Knöbelreiter, Alexander Shekhovtsov, and Thomas Pock. End-to-End Training of Hybrid CNN-CRF Models for Semantic Segmentation using Structured Learning. In *Computer Vision Winter Workshop*, Feb. 2017. 2

[30] Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of Operations Research*, 153(1):235–256, June 2007. 2

[31] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes Dataset for Semantic Urban Scene Understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3223, 2016. 12

[32] Daniel Cremers, Thomas Pock, Kalin Kolev, and A. Chambolle. Convex Relaxation Techniques for Segmentation, Stereo and Multiview Reconstruction. In *Markov Random Fields for Vision and Image Processing*. MIT Press, Boston, 2011. 1, 2, 7

[33] Juan Carlos De Los Reyes, Carola-Bibiane Schönlieb, and Tuomo Valkonen. The structure of optimal parameters for image restoration problems. *Journal of Mathematical Analysis and Applications*, 434(1):464–500, Feb. 2016. 2

[34] Juan Carlos De Los Reyes, Carola-Bibiane Schönlieb, and Tuomo Valkonen. Bilevel Parameter Learning for Higher-Order Total Variation Regularisation Models. *J Math Imaging Vis*, 57(1):1–25, Jan. 2017. 2

[35] Stephan Dempe. *Foundations of Bilevel Programming*. Nonconvex Optimization and Its Applications. Springer US, 2002. 2

[36] Stephan Dempe and Joydeep Dutta. Is bilevel programming a special case of a mathematical program with complementarity constraints? *Math. Program.*, 131(1):37–48, Feb. 2012. 2

[37] Stephan Dempe, Vyacheslav Kalashnikov, Gerardo A. Pérez-Valdés, and Nataliya Kalashnykova. *Bilevel Programming Problems: Theory, Algorithms and Applications to Energy Networks*. Energy Systems. Springer-Verlag, Berlin Heidelberg, 2015. 2

[38] Samuel G. Finlayson, Hyung Won Chung, Isaac S. Kohane, and Andrew L. Beam. Adversarial Attacks Against Medical Deep Learning Systems. *arXiv:1804.05296 [cs, stat]*, Apr. 2018. 1

[39] Jonas Geiping and Michael Moeller. Composite Optimization by Nonconvex Majorization-Minimization. *SIAM J. Imaging Sci.*, pages 2494–2528, Jan. 2018. 6

[40] Kevin Gimpel and Noah A. Smith. Softmax-Margin CRFs: Training Log-Linear Models with Cost Functions. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 733–736, Los Angeles, California, June 2010. Association for Computational Linguistics. 2

[41] Stephen Gould, Basura Fernando, Anoop Cherian, Peter Anderson, Rodrigo Santa Cruz, and Edison Guo. On Differentiating Parameterized Argmin and Argmax Problems with Application to Bi-level Optimization. *arXiv:1607.05447 [cs, math]*, July 2016. 2

[42] Andreas Griewank. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2000. 2

[43] Kerstin Hammernik, Teresa Klatzer, Erich Kobler, Michael P. Recht, Daniel K. Sodickson, Thomas Pock, and Florian Knoll. Learning a variational network for reconstruction of accelerated MRI data. *Magn Reson Med*, 79(6):3055–3071, June 2018. 2

[44] Kerstin Hammernik, Tobias Würfl, Thomas Pock, and Andreas Maier. A Deep Learning Architecture for Limited-Angle Computed Tomography Reconstruction. In *Bildverarbeitung für die Medizin 2017*, Informatik aktuell, pages 92–97. Springer Berlin Heidelberg, 2017. 2

[45] Michael Hintermüller and Carlos N. Rautenberg. Optimal Selection of the Regularization Function in a Weighted Total Variation Model. Part I: Modelling and Theory. *J Math Imaging Vis*, 59(3):498–514, Nov. 2017. 2

[46] Michael Hintermüller, Carlos N. Rautenberg, Tao Wu, and Andreas Langer. Optimal Selection of the Regularization Function in a Weighted Total Variation Model. Part II: Algorithm, Its Analysis and Numerical Tests. *J Math Imaging Vis*, 59(3):515–533, Nov. 2017. 2

[47] Jinggang Huang and D. Mumford. Statistics of natural images and models. In *Proceedings. 1999 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (Cat. No PR00149)*, volume 1, pages 541–547 Vol. 1, June 1999. 12

[48] David R Hunter and Kenneth Lange. A Tutorial on MM Algorithms. *The American Statistician*, 58(1):30–37, Feb. 2004. 6

[49] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1647–1655, July 2017. 1

[50] Kyong Hwan Jin, Michael T. McCann, Emmanuel Froustey, and Michael Unser. Deep Convolutional Neural Network for Inverse Problems in Imaging. *IEEE Transactions on Image Processing*, 26(9):4509–4522, Sept. 2017. 7

[51] Eunhee Kang, Junhong Min, and Jong Chul Ye. A deep convolutional neural network using directional wavelets for low-dose X-ray CT reconstruction. *Medical Physics*, 44(10):e360–e375, Oct. 2017. 7

[52] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International Conference on*

*Learning Representations (ICLR)*, San Diego, May 2015. 12

[53] Teresa Klatzer, Kerstin Hammernik, Patrick Knobelreiter, and Thomas Pock. Learning joint demosaicing and denoising based on sequential energy minimization. In *2016 IEEE International Conference on Computational Photography (ICCP)*, pages 1–11, May 2016. 2

[54] Patrick Knobelreiter, Christian Reinbacher, Alexander Shekhovtsov, and Thomas Pock. End-to-End Training of Hybrid CNN-CRF Models for Stereo. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1456–1465, Honolulu, HI, July 2017. IEEE. 2

[55] Charles. D. Kolstad and Leon S. Lasdon. Derivative evaluation and computational experience with large bilevel mathematical programs. *J Optim Theory Appl*, 65(3):485–499, June 1990. 2

[56] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012. 1

[57] Karl Kunisch and Thomas Pock. A Bilevel Optimization Approach for Parameter Learning in Variational Models. *SIAM Journal on Imaging Sciences*, 6(2):938–983, Jan. 2013. 2

[58] Måns Larsson, Anurag Arnab, Fredrik Kahl, Shuai Zheng, and Philip Torr. A Projected Gradient Descent Method for CRF Inference Allowing End-to-End Training of Arbitrary Pairwise Potentials. In *Energy Minimization Methods in Computer Vision and Pattern Recognition*, Lecture Notes in Computer Science, pages 564–579. Springer International Publishing, 2018. 2

[59] Måns Larsson, Anurag Arnab, Shuai Zheng, Philip Torr, and Fredrik Kahl. Revisiting Deep Structured Models for Pixel-Level Labeling with Gradient-Based Inference. *SIAM J. Imaging Sci.*, pages 2610–2628, Jan. 2018. 2

[60] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, May 2015. 1

[61] Yann LeCun, Sumit Chopra, Raia Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006. 3, 4

[62] Yann LeCun and Fu Jie Huang. Loss functions for discriminative training of energy-based models. In *AISTATS 2005 - Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics*, pages 206–213, 2005. 3, 4

[63] Guosheng Lin, Chunhua Shen, Anton van den Hengel, and Ian Reid. Efficient Piecewise Training of Deep Structured Models for Semantic Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3194–3203, 2016. 2

[64] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully Convolutional Networks for Semantic Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 1

[65] Haihao Lu, Robert M. Freund, and Yurii Nesterov. Relatively Smooth Convex Optimization by First-Order Meth-

ods, and Applications. *SIAM J. Optim.*, pages 333–354, Jan. 2018. 4

[66] Julien Mairal. Optimization with First-order Surrogate Functions. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ICML'13, pages III–783–III–791, Atlanta, GA, USA, 2013. JMLR.org. 6

[67] Julien Mairal. Incremental Majorization-Minimization Optimization with Application to Large-Scale Machine Learning. *SIAM Journal on Optimization*, 25(2):829–855, Jan. 2015. 6

[68] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A Database of Human Segmented Natural Images and its Application to Evaluating Segmentation Algorithms and Measuring Ecological Statistics. In *Proceedings of 8th International Conference on Computer Vision*, volume 2, pages 416–423, July 2001. 8, 12

[69] André F. T. Martins, Noah A. Smith, and Eric P. Xing. Polyhedral outer approximations with application to natural language parsing. In *Proceedings of the 26th Annual International Conference on Machine Learning - ICML '09*, pages 1–8, Montreal, Quebec, Canada, 2009. ACM Press. 2

[70] Nikolaus Mayer, Eddy Ilg, Philip Häusser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4040–4048, June 2016. 1

[71] Jan Hendrik Metzen, Mummadi Chaithanya Kumar, Thomas Brox, and Volker Fischer. Universal Adversarial Perturbations Against Semantic Image Segmentation. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2774–2783, Oct. 2017. 1

[72] Claudia Nieuwenhuis, Eno Töppe, and Daniel Cremers. A Survey and Comparison of Discrete and Continuous Multi-label Optimization Approaches for the Potts Model. *Int J Comput Vis*, 104(3):223–240, Sept. 2013. 7

[73] Peter Ochs, René Ranftl, Thomas Brox, and Thomas Pock. Bilevel Optimization with Nonsmooth Lower Level Problems. In *Scale Space and Variational Methods in Computer Vision*, Lecture Notes in Computer Science, pages 654–665. Springer International Publishing, 2015. 2

[74] Peter Ochs, René Ranftl, Thomas Brox, and Thomas Pock. Techniques for Gradient-Based Bilevel Optimization with Non-smooth Lower Level Problems. *J Math Imaging Vis*, 56(2):175–194, Oct. 2016. 12

[75] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS 2017 Autodiff Workshop*, Long Beach, CA, 2017. 12

[76] Simeon Reich and Shoham Sabach. Existence and Approximation of Fixed Points of Bregman Firmly Nonexpansive Mappings in Reflexive Banach Spaces. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, Springer Optimization and Its Applications, pages 301–316. Springer, New York, NY, 2011. 5, 9

[77] Gernot Riegler, Matthias Rüther, and Horst Bischof. Atgv-net: Accurate depth super-resolution. In *European Conference on Computer Vision*, pages 268–284. Springer, 2016. 2

[78] Ludwig Ritschl, Frank Bergner, Christof Fleischmann, and Marc Kachelrie\s s. Improved total variation-based CT image reconstruction applied to clinical data. *Phys. Med. Biol.*, 56(6):1545–1561, Feb. 2011. 1

[79] R. Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, N.J, 1970. 7

[80] Torsten Rohlfing, Calvin R. Maurer, David A. Bluemke, and M. A. Jacobs. Volume-preserving nonrigid registration of MR breast images using free-form deformation with an incompressibility constraint. *IEEE Transactions on Medical Imaging*, 22(6):730–741, June 2003. 1

[81] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Lecture Notes in Computer Science, pages 234–241. Springer International Publishing, 2015. 1

[82] Frank Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, 1958. 4

[83] Leonid I. Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, Nov. 1992. 1, 8

[84] Kegan G. G. Samuel and Marshall F. Tappen. Learning optimized MAP estimates in continuously-valued MRF models. In *IEEE Conference on Computer Vision and Pattern Recognition*, Miami, FL, USA, June 2009. IEEE. 2

[85] Gilles Savard and Jacques Gauvin. The steepest descent direction for the nonlinear bilevel programming problem. *Operations Research Letters*, 15(5):265–272, June 1994. 2

[86] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv:1312.6120 [cond-mat, q-bio, stat]*, Dec. 2013. 12

[87] Ali Shafahi, W. Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. Are adversarial examples inevitable? In *ICLR 2019*, New Orleans, Sept. 2018. 1

[88] Ying Sun, Prabhu Babu, and Daniel P. Palomar. Majorization-Minimization Algorithms in Signal Processing, Communications, and Machine Learning. *IEEE Transactions on Signal Processing*, 65(3):794–816, Feb. 2017. 6

[89] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *arXiv:1312.6199 [Cs]*, Dec. 2013. 1

[90] Marshall F. Tappen, Ce Liu, Edward H. Adelson, and William T. Freeman. Learning Gaussian Conditional Random Fields for Low-Level Vision. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, June 2007. 4

[91] Ben Taskar, Vassil Chatalbashev, Daphne Koller, and Carlos Guestrin. Learning Structured Prediction Models: A Large Margin Approach. In *Proceedings of the 22nd International Conference on Machine Learning*, ICML '05, pages 896–903, New York, NY, USA, 2005. ACM. 4, 5

[92] Ben Taskar, Carlos Guestrin, and Daphne Koller. Max-Margin Markov Networks. In *Advances in Neural Information Processing Systems 16*, pages 25–32. MIT Press, 2004. 2

[93] Ben Taskar, Simon Lacoste-Julien, and Michael I. Jordan. Structured Prediction, Dual Extragradient and Bregman Projections. *Journal of Machine Learning Research*, 7(Jul):1627–1653, 2006. 4, 5

[94] Marc Teboulle. A simplified view of first order methods for optimization. *Math. Program.*, pages 1–30, May 2018. 4, 6

[95] Ioannis Tsochantaridis, Thorsten Joachims, Thomas Hofmann, and Yasemin Altun. Large Margin Methods for Structured and Interdependent Output Variables. *Journal of Machine Learning Research*, 6(Sep):1453–1484, 2005. 2

[96] Vladimir Vapnik. *Statistical Learning Theory. 1998*, volume 3. Wiley, New York, 1998. 3, 4, 6

[97] Sam Wiseman and Alexander M. Rush. Sequence-to-Sequence Learning as Beam-Search Optimization. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1296–1306, Austin, Texas, Nov. 2016. Association for Computational Linguistics. 2

[98] Wei Zhan, Jiachen Li, Yeping Hu, and Masayoshi Tomizuka. Safe and feasible motion generation for autonomous driving via constrained policy net. In *IECON 2017 - 43rd Annual Conference of the IEEE Industrial Electronics Society*, pages 4588–4593, Oct. 2017. 1

[99] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a Gaussian Denoiser: Residual Learning of Deep CNN for Image Denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, July 2017. 1, 12

[100] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip H. S. Torr. Conditional Random Fields as Recurrent Neural Networks. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1529–1537, Dec. 2015. 2