



OPEN Explainable AI via learning to optimize

Howard Heaton^{1✉} & Samy Wu Fung^{2✉}

Indecipherable black boxes are common in machine learning (ML), but applications increasingly require explainable artificial intelligence (XAI). The core of XAI is to establish transparent and interpretable data-driven algorithms. This work provides concrete tools for XAI in situations where prior knowledge must be encoded and untrustworthy inferences flagged. We use the “learn to optimize” (L2O) methodology wherein each inference solves a data-driven optimization problem. Our L2O models are straightforward to implement, directly encode prior knowledge, and yield theoretical guarantees (e.g. satisfaction of constraints). We also propose use of interpretable certificates to verify whether model inferences are trustworthy. Numerical examples are provided in the applications of dictionary-based signal recovery, CT imaging, and arbitrage trading of cryptoassets. Code and additional documentation can be found at <https://xai-l2o.research.typl.academy>.

A paradigm shift in machine learning is to construct explainable and transparent models, often called explainable AI (XAI)¹. This is crucial for sensitive applications like medical imaging and finance (e.g. see recent work on the role of explainability^{2–5}). Yet, many commonplace models (e.g. fully connected feed forward) offer limited interpretability. Prior XAI works give explanations via tools like sensitivity analysis⁵ and layer-wise propagation^{6,7}, but these neither quantify trustworthiness nor necessarily shed light on how to correct “bad” behaviours. Our work shows how learning to optimize (L2O) can be used to directly embed explainability into models.

The scope of this work is machine learning (ML) applications where domain experts can create approximate models by hand. In our setting, the inference $\mathcal{N}_\Theta(d)$ of a model \mathcal{N}_Θ with input d solves an optimization problem. That is, we use

$$\mathcal{N}_\Theta(d) \triangleq \arg \min_{x \in \mathcal{C}_\Theta(d)} f_\Theta(x; d), \quad (1)$$

where f_Θ is a function and $\mathcal{C}_\Theta(d) \subseteq \mathbb{R}^n$ is a constraint set (e.g. encoding prior information like physical quantities), and each (possibly) includes dependencies on weights Θ . Note the model \mathcal{N}_Θ is *implicit* since its output is defined by an optimality condition rather than an explicit computation. To clarify the scope of the word *explainable* in this work, we adopt the following conventions. We say a model is explainable provided a domain expert can identify the core design elements of a model and how they translate to expected inference properties. We say a *particular inference* is explainable provided its properties can be linked to the model’s design and intended use. Explainable models and inferences are achieved via L2O with our proposed certificates.

A standard practice in software engineering is to code post-conditions after function calls return. Post-conditions are criteria used to validate what the user expects from the code and ensure code is not executed under the wrong assumptions⁸. We propose use of these for ML model inferences (see Fig. 1 and Supplementary Fig. A1). These conditions enable use of certificates with labels—pass, warning or fail—to describe each model inference. We define an inference to be *trustworthy provided it satisfies all provided post-conditions*.

Two ideas, optimization and certificates, form a concrete notion of XAI. Prior and data-driven knowledge can be encoded via optimization, and this encoding can be verified via certificates. To illustrate, consider inquiring why a model generated a “bad” inference (e.g. an inference disagrees with observed measurements). The first diagnostic step is to check certificates. If no fails occurred, the model was not designed to handle the instance encountered. In this case, the model in (1) can be redesigned to encode prior knowledge of the situation. Alternatively, each failed certificate shows a type of error and often corresponds to portions of the model (see Figs. 1 and 2). The L2O model allows debugging of algorithmic implementations and assumptions to correct errors. In a sense, this setup enables one to manually backpropagate errors to fix models (similar to training).

Contributions. This work brings new explainability and guarantees to deep learning applications using prior knowledge. We propose novel implicit L2O models with intuitive design, memory efficient training, infer-

¹Typl Academy, Richland, USA. ²Department of Applied Mathematics and Statistics, Colorado School of Mines, Golden, USA. ✉email: research@typl.academy; swufung@mines.edu

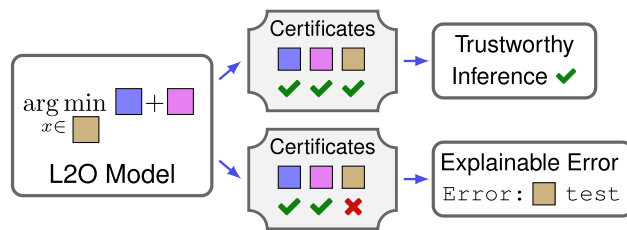


Figure 1. The L2O model is composed of parts (shown as colored blocks) based on prior knowledge or data. L2O inferences solve the optimization problem for given model inputs. Certificates label if each inference is consistent with training. If so, it is trustworthy; otherwise, the faulty model part errs.

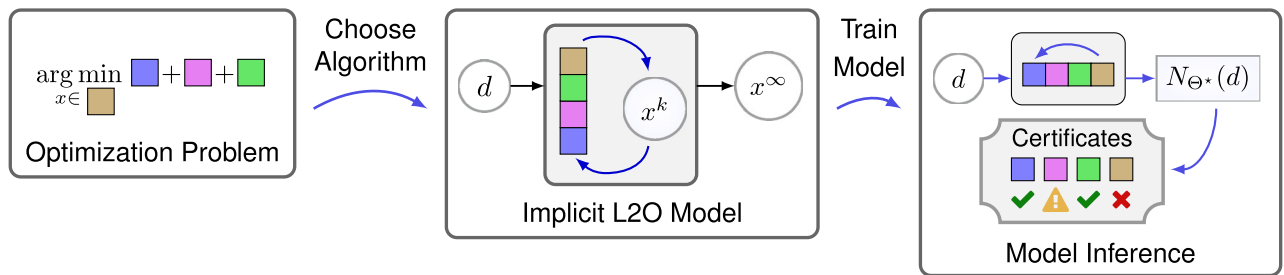


Figure 2. Left shows learning to optimize (L2O) model. Colored blocks denote prior knowledge and data-driven terms. Middle shows an iterative algorithm formed from the blocks (e.g. via proximal/gradient operators) to solve optimization problem. Right shows a trained model's inference $\mathcal{N}_{\Theta^*}(d)$ and its certificates. Certificates identify if properties of inferences are consistent with training data. Each label is associated with properties of specific blocks (indicated by labels next to blocks in right schematic). Labels take value pass \checkmark , warning \triangle , or fail \times , and values identify if inference features for model parts are trustworthy.

ences that satisfy optimality/constraint conditions, and certificates that either indicate trustworthiness or flag inconsistent inference features.

Related works. Closely related to our work is deep unrolling, a subset of L2O wherein models consist of a fixed number of iterations of a data-driven optimization algorithm. Deep unrolling has garnered great success and provides intuitive model design. We refer readers to recent surveys^{9–12} for further L2O background. Downsides of unrolling are growing memory requirements with unrolling depth and a lack of guarantees.

Implicit models circumvent these two shortcomings by defining models using an equation (e.g. as in (1)) rather than prescribe a fixed number of computations as in deep unrolling. This enables inferences to be computed by iterating until convergence, thereby enabling theoretical guarantees. Memory-efficient training techniques were also developed for this class of models, which have been applied successfully in games¹³, music source separation¹⁴, language modeling¹⁵, segmentation¹⁶, and inverse problems^{17,18}. The recent work¹⁸ most closely aligns with our L2O methodology.

Related XAI works use labels/cards. Model Cards¹⁹ document intended and appropriate uses of models. Care labels^{20,21} are similar, testing properties like expressivity, runtime, and memory usage. FactSheets²² are modeled after supplier declarations of conformity and aim to identify models' intended use, performance, safety, and security. These works provide statistics at the distribution level, complementing our work for trustworthiness of individual inferences.

Explainability via optimization

Model design. The design of L2O models is naturally decomposed into two steps: optimization formulation and algorithm choice. The first step is to identify a tentative objective to encode prior knowledge via regularization (e.g. sparsity) or constraints (e.g. unit simplex for classification). We may also add terms that are entirely data-driven. Informally, this step identifies a special case of (1) of the form

$$\mathcal{N}_{\Theta}(d) \triangleq \arg \min_x (\text{prior knowledge}) + (\text{data-driven terms}), \quad (2)$$

where the constraints are encoded in the objective using indicator functions, equaling 0 when constraint is satisfied and ∞ otherwise. The second design step is to choose an algorithm for solving the chosen optimization problem (e.g. proximal-gradient or ADMM²³). We use iterative algorithms, and the update formula for each iteration is given by a *model operator* $T_{\Theta}(x; d)$. Updates are typically composed in terms of gradient and proximal operations. Some parameters (e.g. step sizes) may be included in the weights Θ to be tuned during training. Given data d , computation of the inference $\mathcal{N}_{\Theta}(d)$ is completed by generating a sequence $\{x_d^k\}$ via the relation

L2O	Implicit	Flags	Obtainable model property
✓			Intuitive design
	✓		Memory efficient
✓	✓		Satisfy constraints + (above)
		✓	Trustworthy inferences
✓	✓	✓	Explainable errors + (above)

Table 1. Summary of design features and corresponding model properties. Design features yield additive properties, as indicated by “+ (above).” Proposed implicit L2O models with certificates have intuitive design, memory efficient training, inferences that satisfy optimality/constraint conditions, certificates of trustworthiness, and explainable errors.

$$x_d^{k+1} = T_{\Theta}(x_d^k; d), \quad \text{for all } k \in \mathbb{N}. \quad (3)$$

By design, $\{x_d^k\}$ converges to a solution of (1), and we set

$$\mathcal{N}_{\Theta}(d) = \lim_{k \rightarrow \infty} x_d^k. \quad (4)$$

In our context, each model inference $\mathcal{N}_{\Theta}(d)$ is defined to be an optimizer as in (1). Hence *properties of inferences can be explained via the optimization model* (1); note this is unlike blackbox models where one has no way of explaining why a particular inference is made. The iterative algorithm is applied successively until stopping criteria are met (i.e. in practice we choose an iterate K , possibly dependent on d , so that $\mathcal{N}_{\Theta}(d) \approx x_d^K$). Because $\{x_d^k\}$ converges, we may adjust stopping criteria to approximate the limit to arbitrary precision, which implies we may provide guarantees on model inferences (e.g. satisfying a linear system of equations to a desired precision^{13,17,18}). The properties of the implicit L2O model (1) are summarized by Table 1.

Example of model design. To make the model design procedure concrete, we illustrate this process on a classic problem: sparse recovery from linear measurements. These problems appear in many applications such as radar imaging²⁴ and speech recognition²⁵. Here the task is to estimate a signal x_d^* via access to linear measurements d satisfying $d = Ax_d^*$ for a known matrix A .

Step 1: Choose model Since true signals are known to be sparse, we include ℓ_1 regularization. To comply with measurements, we add a fidelity term. Lastly, to capture hidden features of the data distribution, we also add a data-driven regularization. Putting these together gives the problem

$$\min_{x \in \mathbb{R}^n} \underbrace{\tau \|x\|_1}_{\text{sparsity}} + \underbrace{\|Ax - d\|_2^2}_{\text{fidelity}} + \underbrace{\|W_1 Ax\|_2^2 + \langle x, W_2 d \rangle}_{\text{data-driven regularizer}}, \quad (5)$$

where $\tau > 0$ and W_1 and W_2 are two tunable matrices. This model encodes a balance of three terms—sparsity, fidelity, data-driven regularization—each quantifiable via (5).

Step 2: Choose Algorithm The proximal-gradient scheme generates a sequence $\{z^k\}$ converging to a limit which solves (5). By simplifying and combining terms, the proximal-gradient method can be written via the iteration

$$z^{k+1} = \eta_{\tau\lambda}(z^k - \lambda W(Az^k - d)), \quad \text{for all } k \in \mathbb{N}, \quad (6)$$

where $\lambda > 0$ is a step-size, W is a matrix defined in terms of W_1, W_2 , and A^{\top} , and η_{θ} is the shrink operator given by

$$\eta_{\theta}(x) \triangleq \text{sign}(x) \max(|x| - \theta, 0). \quad (7)$$

From the update on the right hand side of (6), we see the step size λ can be “absorbed” into the tunable matrix W and the shrink function parameter can be set to $\theta > 0$. That is, this example model has weights $\Theta = (W, \theta, \tau)$ with model operator

$$T_{\Theta}(x; d) \triangleq \eta_{\theta}(x - W(Ax - d)), \quad (8)$$

which resembles the updates of previous L2O works^{26–28}. Inferences are computed via a sequence $\{x_d^k\}$ with updates

$$x_d^{k+1} = T_{\Theta}(x_d^k; d), \quad \text{for all } k \in \mathbb{N}. \quad (9)$$

The model inference is the limit x_d^{∞} of this sequence $\{x_d^k\}$.

Convergence. Evaluation of the model $\mathcal{N}_{\Theta}(d)$ is well-defined and tractable under a simple assumption. By a classic result²⁹, it suffices to ensure, for all d , $T_{\Theta}(\cdot; d)$ is *averaged*, i.e. there is $\alpha \in (0, 1)$ and Q such that $T_{\Theta}(x; d) = (1 - \alpha)x + \alpha Q(x; d)$, where Q is 1-Lipschitz. When this property holds, the sequence $\{x_d^k\}$ in (3) converges to a solution x_d^* . This may appear to be a strong assumption; however, common operations in convex optimization algorithms (e.g. proximals and gradient descent updates) are averaged. For entirely data-driven portions of T_{Θ} , several activation functions are 1-Lipschitz^{30,31} (e.g. ReLU and softmax), and libraries like

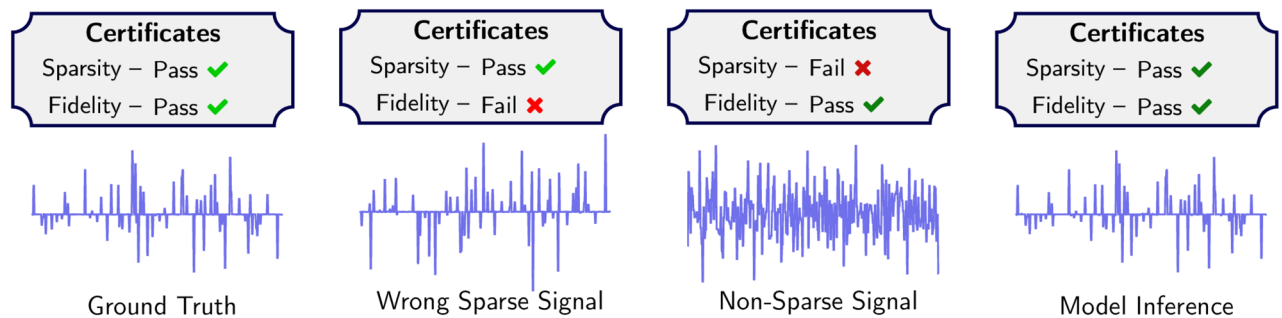


Figure 3. Example inferences for test data d . The sparsified version Kx of each inference x is shown (c.f. Fig. 5) along with certificates. Ground truth was taken from test dataset of implicit dictionary experiment. The second from left is sparse and inconsistent with measurement data. The second from right complies with measurements but is not sparse. The rightmost is generated using our proposed model (IDM), which approximates the ground truth well and is trustworthy.

Concept	Quantity	Formula
Sparsity	Nonzeros	$\ x\ _0$
Measurements	Relative error	$\ Ax - d\ /\ d\ $
Constraints	Distance to set \mathcal{C}	$d_{\mathcal{C}}(x)$
Smooth images	Total variation	$\ \nabla x\ _1$
Classifier Confidence	Probability short of one-hot label	$1 - \max_i x_i$
Convergence	Iterate residual	$\ x^k - x^{k-1}\ $
Regularization	Proximal residual	$\ x - \text{prox}_{f_{\Omega}}(x)\ $

Table 2. Certificate examples. Each certificate is tied to a high-level concept, and then quantified in a formula. For classifier confidence, we assume x is in the unit simplex. The proximal is a data-driven update for f_{Ω} with weights Ω .

PyTorch³² include functionality to force affine mappings to be 1-Lipschitz (e.g. spectral normalization). Furthermore, by making $T_{\Theta}(\cdot; d)$ a contraction, a unique fixed point is obtained. We emphasize, even without forcing T_{Θ} to be averaged, $\{x^k\}$ is often observed to converge in practice^{15,17,18} upon tuning the weights Θ .

Trustworthiness certificates. Explainable models justify whether each inference is trustworthy. We propose providing justification in the form of certificates, which verify various properties of the inference are consistent with those of the model inferences on training data and/or prior knowledge. Each certificate is a tuple of the form (name, label) with a property name and a corresponding label which has one of three values: pass, warning, or fail (see Fig. 3). Each certificate label is generated by two steps. The first is to apply a function that maps inferences (or intermediate states) to a *nonnegative scalar value* α quantifying a property of interest. The second step is to map this scalar to a label. Labels are generated via the flow:

$$\text{Inference} \rightarrow \text{Property Value} \rightarrow \text{Certificate Label.} \quad (10)$$

Property value functions. Several quantities may be used to generate certificates. In the model design example above, a sparsity property can be quantified by counting the number of nonzero entries in a signal, and a fidelity property can use the relative error $\|Ax - d\|/\|d\|$ (see Fig. 3). To be most effective, property values are chosen to coincide with the optimization problem used to design the L2O model, i.e. to quantify structure of prior and data-driven knowledge. This enables each certificate to clearly validate a portion of the model (see Fig. 2). Since various concepts are useful for different types of modeling, we provide a brief (and non-comprehensive) list of concepts and possible corresponding property values in Table 2.

One property concept deserves particular attention: data-driven regularization. This regularization is important for discriminating between inference features that are qualitatively intuitive but difficult to quantify by hand. Rather than approximate a function, implicit L2O models directly approximate gradients/proximals. These provide a way to measure regularization indirectly via gradient norms/residual norms of proximals. Moreover, these norms (e.g. see last row of Table 2) are easy to compute and equal zero only at local minima of regularizers. To our knowledge, this is the first work to *quantify* trustworthiness using the quality of inferences with respect to data-driven regularization.

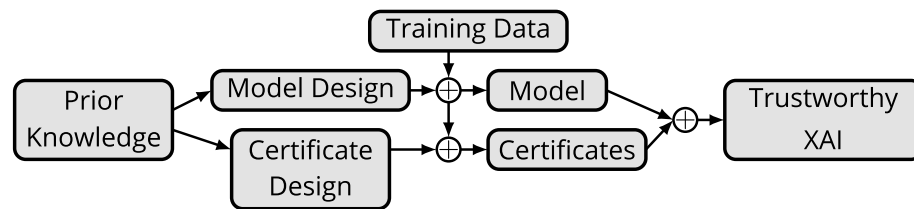


Figure 4. This diagram illustrates relationships between certificates, models, training data, and prior knowledge. Prior knowledge is embedded directly into model design via the L2O methodology. This also gives rise to quantities to measure for certificate design. The designed model is tuned using training data to obtain the “optimal” L2O model (shown by arrows touching top middle + sign). The certificates are tuned to match the test samples and/or model inferences on training data (shown by arrows with bottom middle + sign). Together the model and certificates yield inferences with certificates of trustworthiness.

Certificate labels. Typical certificate labels should follow a trend where inferences often obtain a pass label to indicate trustworthiness while warnings occur occasionally and failures are obtained in extreme situations. Let the samples of model inference property values $\alpha \in [0, \infty)$ come from distribution $\mathbb{P}_{\mathcal{A}}$. We pick property value functions for which small α values are desirable and the distribution tail consists of larger α . Intuitively, smaller property values of α resemble property values of inferences from training and/or test data. Thus, labels are assigned according to the probability of observing a value less than or equal to α , i.e. we evaluate the cumulative distribution function (CDF) defined for probability measure $\mathbb{P}_{\mathcal{A}}$ by

$$\text{CDF}(\alpha) = \int_0^{\alpha} d\mathbb{P}_{\mathcal{A}}, \quad (11)$$

Labels are chosen according to the task at hand. Let p_p , p_w , and $p_f = 1 - p_p - p_w$ be the probabilities for pass, warning, and fail labels, respectively. Labels are made for α via

$$\text{Label}(\alpha) = \begin{cases} \text{pass} & \text{if } \text{CDF}(\alpha) < p_p \\ \text{warning} & \text{if } \text{CDF}(\alpha) \in [p_p, 1 - p_f) \\ \text{fail} & \text{otherwise.} \end{cases} \quad (12)$$

The remaining task is to estimate the CDF value for a given α . Recall we assume access is given to property values $\{\alpha_i\}_{i=1}^N$ from ground truths or inferences on training data, where N is the number of data points. To this end, given an α value, we estimate its CDF value via the empirical CDF:

$$\text{CDF}(\alpha) \approx \frac{|\{\alpha_i : \alpha_i \leq \alpha, 1 \leq i \leq N\}|}{N} \quad (13a)$$

$$= \frac{\# \text{ of } \alpha_i\text{'s } \leq \alpha}{N}, \quad (13b)$$

where $|\cdot|$ denotes set cardinality. Figure 4 shows how these certificates can be combined with the L2O methodology.

Certificate implementation. As noted in the introduction, trustworthiness certificates are evidence an inference satisfies post-conditions (i.e. passes various tests). Thus, they are to be used in code in the same manner as standard software engineering practice. Consider the snippet of code in Supplementary Fig. A1. As usual, an inference is generated by calling the model. However, alongside the inference SPSVERBc1, certificates SPSVERBc2 are returned that label whether the inference SPSVERBc1 passes tests that identify consistency with training data and prior knowledge.

Experiments

Each numerical experiment shows an application of novel implicit L2O models, which were designed directly from prior knowledge. Associated certificates of trustworthiness are used to emphasize the explainability of each model and illustrate use-cases of certificates. Experiments were coded using Python with the PyTorch library³², the Adam optimizer³³, and, for ease of re-use, were run via Google Colab. We emphasize these experiments are for illustration of intuitive and novel model design and trustworthiness and are not benchmarked against state-of-the-art models. The datasets generated and/or analysed during the current study are available in the following repository: github.com/typal-research/xai-l2o. All methods were performed in accordance with the relevant guidelines and regulations.

Algorithms. To illustrate evaluation of L2O model used herein, we begin with an example L2O model and algorithm. Specifically, models used for the first two experiments take the form

$$\min_{x \in \mathbb{R}^n} f(Kx) + h(x) \quad \text{s.t.} \quad \|Mx - d\| \leq \delta, \quad (14)$$

where K and M are linear operators, $\delta \geq 0$ is a noise tolerance, and f and g are proximable functions. Introducing auxiliary variables w and p and dual variable $v = (v_1, v_2)$, linearized ADMM³⁴ (L-ADMM) can be used to iteratively update the tuple (p, w, v, x) of variables via

$$p^{k+1} = \text{prox}_{\lambda f} \left(p^k + \lambda(v_1^k + \alpha(Kx^k - p^k)) \right) \quad (15a)$$

$$w^{k+1} = \text{proj}_{B(d, \delta)} \left(w^k + \lambda(v_2^k + \alpha(Mx^k - w^k)) \right) \quad (15b)$$

$$v_1^{k+1} = v_1^k + \alpha(Kx^k - p^{k+1}) \quad (15c)$$

$$v_2^{k+1} = v_2^k + \alpha(Mx^k - w^{k+1}) \quad (15d)$$

$$r^k = K^\top (2v_1^{k+1} - v_1^k) + M^\top (2v_2^{k+1} - v_2^k) \quad (15e)$$

$$x^{k+1} = \text{prox}_{\beta h} (x^k - \beta r^k), \quad (15f)$$

where $\text{proj}_{B(d, \delta)}$ is the Euclidean projection onto the Euclidean ball of radius δ centered at d , prox_f is the proximal operator for a function f , and the scalars $\alpha, \beta, \lambda > 0$ are appropriate step sizes. Further details, definitions, and explanations are available in the appendices. We note the updates are ordered so that x^{k+1} is the final step to make it easy to backprop through the final x^k update.

Implicit model training. Standard backpropagation cannot be used for implicit models as it requires memory capacities beyond existing computing devices. Indeed, storing gradient data for each iteration in the forward propagation (see (3)) scales the memory during training linearly with respect to the number of iterations. Since the limit x^∞ solves a fixed point equation, implicit models can be trained by differentiating implicitly through the fixed point to obtain a gradient. This implicit differentiation requires further computations and coding. Instead of using gradients, we utilize Jacobian-Free Backpropagation (JFB)³⁵ to train models. JFB further simplifies training by only backpropagating through the final iteration, which was proven to yield preconditioned gradients. JFB trains using fixed memory (with respect to the K steps used to estimate $\mathcal{N}_\Theta(d)$) and avoids numerical issues arising from computing exact gradients³⁶, making JFB and its variations^{37,38} apt for training implicit models.

Implicit dictionary learning. *Setup.* In practice, high dimensional signals often approximately admit low dimensional representations^{39–44}. For illustration, we consider a linear inverse problem where true data admit sparse representations. Here each signal $x_d^* \in \mathbb{R}^{250}$ admits a representation $s_d^* \in \mathbb{R}^{50}$ via a transformation M (i.e. $x_d^* = Ms_d^*$). A matrix $A \in \mathbb{R}^{100 \times 250}$ is applied to each signal x_d^* to provide linear measurements $d = Ax_d^*$. Our task is to recover x_d^* given knowledge of A and d without the matrix M . Since the linear system is quite under-determined, schemes solely minimizing measurement error (e.g. least squares approaches) fail to recover true signals; additional knowledge is essential.

Model design. All convex regularization approaches are known lead to biased estimators whose expectation does not equal the true signal⁴⁶. However, the seminal work⁴⁷ of Candes and Tao shows ℓ_1 minimization (rather than additive regularization) enables exact recovery under suitable assumptions. Thus, we minimize a sparsified signal subject to linear constraints via the implicit dictionary model (IDM)

$$\mathcal{N}_\Theta(d) \triangleq \arg \min_{x \in \mathbb{R}^{250}} \|Kx\|_1 \quad \text{s.t.} \quad Ax = d. \quad (16)$$

The square matrix K is used to leverage the fact x has a low-dimensional representation by transforming x into a sparse vector. Linearized ADMM³⁴ (L-ADMM) is used to create a sequence $\{x_d^k\}$ as in (3). The model \mathcal{N}_Θ has weights $\Theta = K$. If it exists, the matrix K^{-1} is known as a dictionary and $K\mathcal{N}_\Theta(d)$ is the corresponding sparse code; hence the name IDM for (16). To this end, we emphasize K is learned during training and is *different* from M , but these matrices are related since we aim for the product $Kx_d^* = KMs_d^*$ to be sparse. Note we use L-ADMM to *provably* solve (16), and \mathcal{N}_Θ is easy to train. More details are in Appendix C.

Discussion. IDM combines intuition from dictionary learning with a reconstruction algorithm. Two properties are used to identify trustworthy inferences: sparsity and measurement compliance (i.e. fidelity). Sparsity and fidelity are quantified via the ℓ_1 norm of the sparsified inference (i.e. $K\mathcal{N}_\Theta(d)$) and relative measurement error. Figure 5 shows the training the model yields a sparsifying transformation K . Figure 3 shows the proposed certificates identify “bad” inferences that might, at first glance, appear to be “good” due to their compatibility with

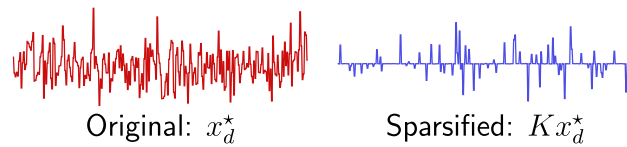


Figure 5. Training IDM yields sparse representation of inferences. Diagram shows a sample true data x (left) from test dataset and its sparsified representation Kx (right).

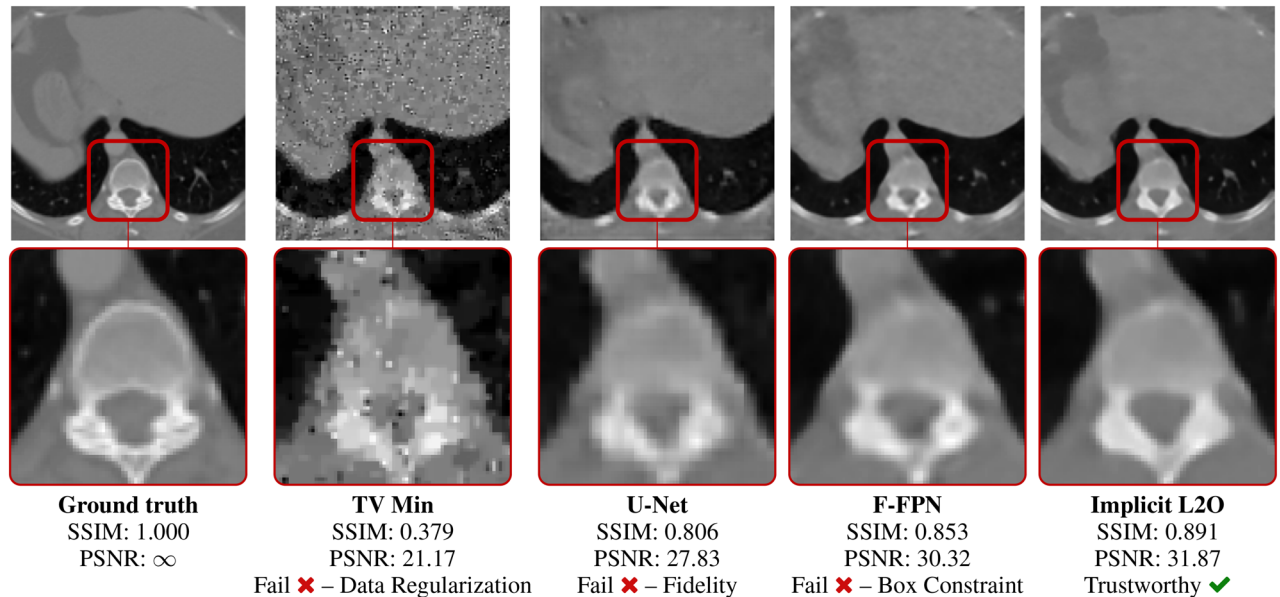


Figure 6. Reconstructions on test data computed via U-Net⁴⁵, TV minimization, F-FPNs¹⁷, and Implicit L2O (left to right). Bottom row shows expansion of region indicated by red box. Pixel values outside $[0, 1]$ are flagged. Fidelity is flagged when images do not comply with measurements, and regularization is flagged when texture features of images are sufficiently inconsistent with true data (e.g. grainy images). Labels are provided beneath each image (*n.b.* fail is assigned to images that are worse than 95% of L2O inferences on training data). Shown comparison methods fail while the Implicit L2O image passes all tests.

constraints. Lastly, observe the utility of learning K , rather than approximating M , is K makes it is easy to check if an inference admits a sparse representation. Using M to check for sparsity is nontrivial.

CT image reconstruction. *Setup.* Comparisons are provided for low-dose CT examples derived from the Low-Dose Parallel Beam dataset (LoDoPab) dataset⁴⁸, which has publically available phantoms derived from actual human chest CT scans. CT measurements are simulated with a parallel beam geometry and a sparse-angle setup of only 30 angles and 183 projection beams, giving 5490 equations and 16,384 unknowns. We add 1.5% Gaussian noise to *each individual beam measurement*. Images have resolution 128×128 . To make errors easier to contrast between methods, the linear systems here are under-determined and have more noise than those in some similar works. Image quality is determined using the Peak Signal-To-Noise Ratio (PSNR) and structural similarity index measure (SSIM). The training loss was mean squared error. Training/test datasets have 20,000/2000 samples.

Model design. The model for the CT experiment extends the IDM. In practice, it has been helpful to utilize a sparsifying transform^{49,50}. We accomplish this via a linear operator K , which is applied and then this product is fed into a data-driven regularizer f_{Ω} with parameters Ω . We additionally ensure compliance with measurements from the Radon transform matrix A , up to a tolerance δ . In our setting, all pixel values are also known to be in the interval $[0, 1]$. Combining our prior knowledge yields the implicit L2O model

$$\mathcal{N}_{\Theta}(d) \triangleq \arg \min_{x \in [0,1]^n} f_{\Omega}(Kx) \quad \text{s.t.} \quad \|Ax - d\| \leq \delta. \quad (17)$$

Here \mathcal{N}_{Θ} has weights $\Theta = (\Omega, K, \alpha, \beta, \lambda)$ with α, β and λ step-sizes in L-ADMM. More details are in Appendix D.

Discussion. Comparisons of our method (Implicit L2O) with U-Net⁴⁵, F-FPNs¹⁷, and total variation (TV) Minimization are given in Fig. 6 and Table 3. Table 3 shows the average PSNR and SSIM reconstructions. Our model

Method	Avg. PSNR	Avg. SSIM	Box constraint fail	Fidelity fail	Data Reg. fail	# Params
U-Net	27.32 dB	0.761	5.75 %	96.95%	3.20%	533,593
TV Min	28.52 dB	0.765	0.00 %	0.00%	25.40%	4
F-FPN [†]	30.46 dB	0.832	47.15%	0.40%	5.05%	96,307
Implicit L2O	31.73 dB	0.858	0.00%	0.00%	5.70%	59,697

Table 3. Average PSNR/SSIM for CT reconstructions on the 2000 image LoDoPab testing dataset. [†] Reported from original work¹⁷. U-Net was trained with filtered backprojection as in prior work⁴⁵. Three properties are used to check trustworthiness: box constraints, compliance with measurements (i.e. fidelity), and data-driven regularization (via the proximal residual in Table 2). Failed sample percentages are numerically estimated via (). Sample property values “fail” if they perform worse than 95% of the inferences on the training data, i.e. , its CDF value exceeds 0.95. Implicit L2O yields the most passes on test data.

obtains the highest average PSNR and SSIM values on the test data while using 11% and 62% as many weights as U-Net and F-FPN, indicating greater efficiency of the implicit L2O framework. Moreover, the L2O model is designed with three features: compliance with measurements (i.e. fidelity), valid pixel values, and data-driven regularization. Table 3 also shows the percentage of “fail” labels for these property values. Here, an inference fails if its property value is larger than 95% of the property values from the training/true data, i.e. we choose $p_p = 0.95$, $p_w = 0$, and $p_f = 0.05$ in (12). For the fidelity, our model never fails (due to incorporating the constraint into the network design). Our network fails 5.7% of the time for the data-driven regularization property. Overall, the L2O model generates the most trustworthy inferences. This is intuitive as this model outperforms the others and was specifically designed to embed all of our knowledge, unlike the others. To provide better intuition of the certificates, we also show the certificate labels for an image from the test dataset in Fig. 6. The only image to pass all provided tests is the proposed implicit L2O model. This knowledge can help identify trustworthy inferences. Interestingly, the data-driven regularization enabled certificates to detect and flag “bad” TV Minimization features (e.g. visible staircasing effects^{51,52}), which shows novelty of certificates as these features are intuitive, yet prior methods to quantify this were, to our knowledge, unknown.

Optimal cryptoasset trading. *Setup.* Ethereum is a blockchain technology anyone can use to deploy permanent and immutable decentralized applications. This technology enables creation of decentralized finance (DeFi) primitives, which can give censorship-resistant participation in digital markets and expand the use of stable assets^{53,54} and exchanges^{55–57} beyond the realm of traditional finance. Popularity of cryptoasset trading (e.g. GRT and Ether) is exploding with the DeFi movement^{58,59}.

Decentralized exchanges (DEXs) are a popular entity for exchanging cryptoassets (subject to a small transaction fee), where trades are conducted without the need for a trusted intermediary to facilitate the exchange. Popular examples of DEXs are constant function market makers (CFMMs)⁶⁰, which use mathematical formulas to govern trades. To ensure CFMMs maintain sufficient net assets, trades within CFMMs maintain constant total reserves (as defined by a function ϕ). A transaction in a CFMM tendering x assets in return for y assets with reserves assets r is accepted provided

$$\phi(r + \gamma x - y) \geq \phi(r), \quad (18)$$

with $\gamma \in (0, 1]$ a trade fee parameter. Here $r, x, y \in \mathbb{R}^n$ with each vector nonnegative and i -th entry giving an amount for the i -th cryptoasset type (e.g. Ether, GRT). Typical choices⁶¹ of ϕ are weighted sums and products, i.e.

$$\phi(r) = \sum_{i=1}^n w_i r_i \quad \text{and} \quad \phi(r) = \prod_{i=1}^n r_i^{w_i}, \quad (19)$$

where $w \in \mathbb{R}^n$ has positive entries. Figure 7 shows an example of a CFMM network.

This experiment aims to maximize arbitrage. Arbitrage is the simultaneous purchase and sale of equivalent assets in multiple markets to exploit price discrepancies between the markets. This can be a lucrative endeavor with cryptoassets⁶². For a given snapshot in time, our arbitrage goal is to identify a collection of trades that maximize the cryptoassets obtainable by trading between different exchanges, i.e. solve the (informal) optimization problem

$$\max_{\text{trade}} \text{Assets}(\text{trade}) \quad \text{s.t.} \quad \text{trade} \in \{\text{valid trades}\}. \quad (20)$$

The set of valid trades is all trades satisfying the transaction rules for CFMMs given by (18) with nonnegative values for tokens tendered and received (i.e. $x, y \geq 0$). Prior works^{61,63} deal with an idealistic noiseless setting while recognizing executing trades is not without risk (e.g. noisy information, front running⁶⁴, and trade delays). To show implications of trade risk, we incorporate noise in our trade simulations by adding noise $\varepsilon \in \mathbb{R}^n$ to CFMM asset observations, which yields noisy observed data $d = (1 + \varepsilon) \odot r$. Also, we consider trades with CFMMs where several assets can be traded simultaneously rather than restricting to pairwise swaps.

Model design. The aim is to create a model that infers a trade (x, y) maximizing utility. For a nonnegative vector $p \in \mathbb{R}^n$ of reference price valuations, this utility U is the net change in asset values provided by the trade, i.e.

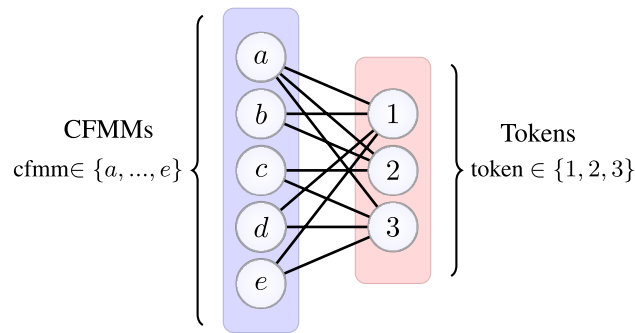


Figure 7. Network with 5 CFMMs and 3 tokens; structure replicates an experiment in recent work⁶³. Black lines show available tokens for trade in each CFMM.

$$U(x, y) \triangleq \underbrace{\sum_{j=1}^m \langle A^j p, A^j (y^j - x^j) \rangle}_{\text{net asset value change}}, \quad (21)$$

where A^j is a matrix mapping global coordinates of asset vector to the coordinates of the j -th CFMM (see Appendix E for details). For noisy data d , trade predictions can include a “cost of risk.” This is quantified by regularizing the trade utility, i.e. introducing a penalty term. For matrices W^j , we model risk by a simple quadratic penalty via

$$U_{\Theta}(x, y) \triangleq U(x, y) - \underbrace{\frac{1}{2} \cdot \sum_{j=1}^m \|A^j W^j (x - y)\|^2}_{\text{risk model}}. \quad (22)$$

The implicit L2O model infers optimal trades via U_{Θ} , i.e.

$$\mathcal{N}_{\Theta}(d) \triangleq (x_d, y_d) = \arg \max_{(x, y) \in \mathcal{C}_{\Theta}(d)} U_{\Theta}(x, y), \quad (23)$$

where $\mathcal{C}_{\Theta}(d)$ encodes constraints for valid transactions. The essence of \mathcal{N}_{Θ} is to output solutions to (20) that account for transaction risks. A formulation of Davis-Yin operator splitting⁶⁵ is used for model evaluation. Further details of the optimization scheme are in Appendix E.

Discussion. The L2O model contains three core features: profit, risk, and trade constraints. The model is designed to output trades that satisfy provided constraints, but note these are *noisy* and thus cannot be used to a priori determine whether a trade will be executed. For this reason, fail flags identify conditions to warn a trader when a trade should be aborted (due to an “invalid trade”). This can avoid wasting transaction fees (i.e. gas costs). Figure 8 shows an example of two trades, where we note the analytic method proposes a large trade that is *not* executed since it violates the trade constraints (due to noisy observations). The L2O method proposes a small trade that yielded arbitrage profits (i.e. $U > 0$) and has pass certificates. Comparisons are provided in Table 4 between the analytic and L2O models. Although the analytic method has “ideal” structure, it performs much worse than the L2O scheme. In particular, *no trades* are executable by the analytic scheme since the present noise always makes the proposed transactions fail to satisfy the actual CFMM constraints. Consistent with this, every proposed trade by the analytic trade is flagged as risky in Table 4. The noise is on the order of 0.2% Gaussian noise of the asset totals.

Conclusions

Explainable ML models can be concretely developed by fusing certificates with the L2O methodology. The implicit L2O methodology enables prior and data-driven knowledge to be directly embedded into models, thereby providing clear and intuitive design. This approach is theoretically sound and compatible with state-of-the-art ML tools. The L2O model also enables construction of our certificate framework with easy-to-read labels, certifying if each inference is trustworthy. In particular, our certificates provide a principled scheme for the detection of inferences with “bad” features via data-driven regularization. Thanks to this optimization-based model design (where inferences can be defined by fixed point conditions), failed certificates can be used to discard untrustworthy inferences and may help debugging the architecture. This reveals the interwoven nature of pairing implicit L2O with certificates. Our experiments illustrate these ideas in three different settings, presenting novel model designs and interpretable results. Future work will study extensions to physics-based applications where PDE-based physics can be integrated into the model^{66–68}.

Method	Predicted utility	Executed utility	Trade execution	Risk fail	Profitable fail	# Params
Analytic	11.446	0.00	0.00%	100.00%	0%	0
Implicit L2O	0.665	0.6785	88.20%	3.6%	11.80%	126

Table 4. Averaged results on test data for trades in CFMM network. The analytic method always predicts a profitable trade, but fails to satisfy the constraints (due to noise). This failure is predicted by the certificates “risk” certificate and reflected by the 0% trade execution. Alternatively, the L2O scheme makes conservative predictions regarding constraints, which limits profitability. However, using these certificates, executed L2O trades are always profitable and satisfy constraints.

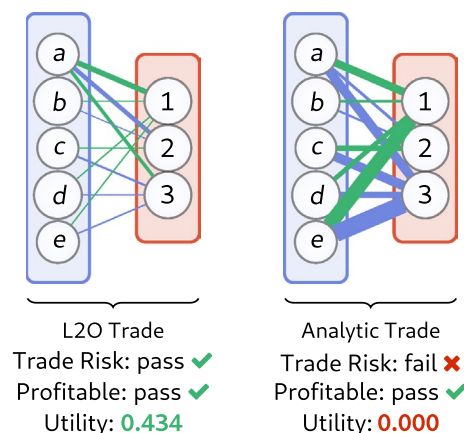


Figure 8. Example of proposed L2O (left) and analytic (right) trades with noisy data d . Blue and green lines show proposed cryptoassets x and y to tender and receive, respectively (widths show magnitude). The analytic trade is unable to account for trade risks, causing it to propose large trades that are *not* executed (giving executed utility of zero). This can be anticipated by the failed trade risk certificate. On the other hand, the L2O scheme is profitable (utility is 0.434) and is executed (consistent with the pass trade risk label).

Data availability

The datasets generated and/or analysed during the current study are available in the following repository: github.com/typal-research/xai-l2o.

Received: 27 June 2022; Accepted: 30 May 2023

Published online: 21 June 2023

References

1. Van Lent, M., Fisher, W., Mancuso, M. An explainable artificial intelligence system for small-unit tactical behavior. In *Proceedings of the National Conference on Artificial Intelligence*, 900–907. (AAAI Press; MIT Press, 1999).
2. Arrieta, A. B. *et al.* Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Inf. Fusion* **58**, 82–115 (2020).
3. Adadi, A. & Berrada, M. Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access* **6**, 52138–52160 (2018).
4. Došilović, F. K., Brčić, M., Hlupić, N. Explainable artificial intelligence: A survey. In *2018 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*, 0210–0215. (IEEE, 2018).
5. Samek, W., Müller, K.-R. Towards explainable artificial intelligence. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, 5–22. (Springer, 2019).
6. Bach, S. *et al.* On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS One* **10**(7), e0130140 (2015).
7. Montavon, G., Binder, A., Lapuschkin, S., Samek, W., Müller, K.-R. Layer-wise relevance propagation: an overview. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, 193–209 (2019).
8. Anaya, M. *Clean Code in Python: Refactor Your Legacy Code Base*. (Packt Publishing Ltd, 2018).
9. Amos, B. Tutorial on amortized optimization for learning to optimize over continuous domains (2022). arXiv preprint [arXiv:2202.00665](https://arxiv.org/abs/2202.00665).
10. Chen, T., Chen, X., Chen, W., Heaton, H., Liu, J., Wang, Z., Yin, W. Learning to optimize: A primer and a benchmark (2021). arXiv preprint [arXiv:2103.12828](https://arxiv.org/abs/2103.12828).
11. Shlezinger, N., Whang, J., Eldar, Y. C., Dimakis, A. G. Model-based deep learning (2020). arXiv preprint [arXiv:2012.08405](https://arxiv.org/abs/2012.08405).
12. Monga, V., Li, Y. & Eldar, Y. C. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Signal Process. Mag.* **38**(2), 18–44 (2021).
13. Heaton, H., McKenzie, D., Li, Q., Fung, S. W., Osher, S., Yin, W. Learn to predict equilibria via fixed point networks (2021). arXiv preprint [arXiv:2106.00906](https://arxiv.org/abs/2106.00906).
14. Koyama, Y., Murata, N., Uhlich, S., Fabbro, G., Takahashi, S., Mitsufuji, Y. Music source separation with deep equilibrium models (2021). arXiv preprint [arXiv:2110.06494](https://arxiv.org/abs/2110.06494).
15. Bai, S., Kolter, J. Z., Koltun, V. Deep equilibrium models (2019). arXiv preprint [arXiv:1909.01377](https://arxiv.org/abs/1909.01377).

16. Bai, S., Koltun, V., Kolter, J. Z. Multiscale deep equilibrium models (2020). arXiv preprint [arXiv:2006.08656](https://arxiv.org/abs/2006.08656).
17. Heaton, H., Fung, S. W., Gibali, A., Yin, W. Feasibility-based fixed point networks (2021). arXiv preprint [arXiv:2104.14090](https://arxiv.org/abs/2104.14090).
18. Gilton, D., Ongie, G., Willett, R. Deep equilibrium architectures for inverse problems in imaging (2021). arXiv preprint [arXiv:2102.07944](https://arxiv.org/abs/2102.07944).
19. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., Spitzer, E., Raji, I. D., Gebru, T. Model cards for model reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 220–229 (2019).
20. Morik, K., Kotthaus, H., Heppel, L., Heinrich, D., Fischer, R., Pauly, A., Piatkowski, N. The care label concept: A certification suite for trustworthy and resource-aware machine learning (2021). arXiv preprint [arXiv:2106.00512](https://arxiv.org/abs/2106.00512).
21. Morik, K., Kotthaus, H., Heppel, L., Heinrich, D., Fischer, R., Mücke, S., Pauly, A., Jakobs, M., Piatkowski, N. Yes We Care!—Certification for machine learning methods through the care label framework (2021). arXiv preprint [arXiv:2105.10197](https://arxiv.org/abs/2105.10197).
22. Arnold, M. *et al.* FactSheets: Increasing trust in AI services through supplier's declarations of conformity. *IBM J. Res. Dev.* **63**(4/5), 1–6 (2019).
23. Deng, W. & Yin, W. On the global and linear convergence of the generalized alternating direction method of multipliers. *J. Sci. Comput.* **66**(3), 889–916 (2016).
24. Siddamal, K., Bhat, S. P., Saroja, V. A survey on compressive sensing. In *2015 2nd International Conference on Electronics and Communication Systems (ICECS)*, 639–643. (IEEE, 2015).
25. Gemmeke, J. F., Van Hamme, H., Cranen, B. & Boves, L. Compressive sensing for missing data imputation in noise robust speech recognition. *IEEE J. Sel. Top. Signal Process.* **4**(2), 272–287 (2010).
26. Liu, J., Chen, X. ALISTA: Analytic weights are as good as learned weights in LISTA. In *International Conference on Learning Representations (ICLR)* (2019).
27. Gregor, K., LeCun, Y. Learning fast approximations of sparse coding. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, 399–406 (2010).
28. Chen, X., Liu, J., Wang, Z., Yin, W. Theoretical linear convergence of unfolded ISTA and its practical weights and thresholds (2018). arXiv preprint [arXiv:1808.10038](https://arxiv.org/abs/1808.10038).
29. Krasnosel'skii, M. Two remarks about the method of successive approximations. *Uspekhi Mat. Nauk* **10**, 123–127 (1955).
30. Combettes, P. L. & Pesquet, J.-C. Lipschitz certificates for layered network structures driven by averaged activation operators. *SIAM J. Math. Data Sci.* **2**(2), 529–557 (2020).
31. Gao, B., Pavel, L. On the properties of the softmax function with application in game theory and reinforcement learning (2017). arXiv preprint [arXiv:1704.00805](https://arxiv.org/abs/1704.00805).
32. Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural. Inf. Process. Syst.* **32**, 8026–8037 (2019).
33. Kingma, D. P., Ba, J. Adam: A method for stochastic optimization (2014). arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
34. Ryu, E., Yin, W. *Large-Scale Convex Optimization: Algorithm Designs via Monotone Operators*. (Cambridge University Press, 2022).
35. Fung, S. W., Heaton, H., Li, Q., McKenzie, D., Osher, S., Yin, W. JFB: Jacobian-free backpropagation for implicit networks (2021). arXiv preprint [arXiv:2103.12803](https://arxiv.org/abs/2103.12803).
36. Bai, S., Koltun, V., Kolter, Z. Stabilizing equilibrium models by jacobian regularization. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research* (eds. Meila, M., Zhang, T.) 554–565 (PMLR, 2021).
37. Geng, Z., Zhang, X.-Y., Bai, S., Wang, Y., Lin, Z. On training implicit models. In *Thirty-Fifth Conference on Neural Information Processing Systems* (2021).
38. Huang, Z., Bai, S., Kolter, J. Z. Implicit²: Implicit layers for implicit representations. *Adv. Neural Inf. Process. Syst.* **34** (2021).
39. Osher, S., Shi, Z. & Zhu, W. Low dimensional manifold model for image processing. *SIAM J. Imag. Sci.* **10**(4), 1669–1690 (2017).
40. Zhang, Z., Xu, Y., Yang, J., Li, X. & Zhang, D. A survey of sparse representation: Algorithms and applications. *IEEE Access* **3**, 490–530 (2015).
41. Carlsson, G., Ishkhanov, T., De Silva, V. & Zomorodian, A. On the local behavior of spaces of natural images. *Int. J. Comput. Vis.* **76**(1), 1–12 (2008).
42. Lee, A. B., Pedersen, K. S. & Mumford, D. The nonlinear statistics of high-contrast patches in natural images. *Int. J. Comput. Vis.* **54**(1–3), 83–103 (2003).
43. Peyré, G. Image processing with nonlocal spectral bases. *Multiscale Model. Simul.* **7**(2), 703–730 (2008).
44. Peyré, G. Manifold models for signals and images. *Comput. Vis. Image Underst.* **113**(2), 249–260 (2009).
45. Jin, K. H., McCann, M. T., Froustey, E. & Unser, M. Deep convolutional neural network for inverse problems in imaging. *IEEE Trans. Image Process.* **26**(9), 4509–4522 (2017).
46. Fan, J. & Li, R. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Stat. Assoc.* **96**(456), 1348–1360 (2001).
47. Candès, E. J., Romberg, J. & Tao, T. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inf. Theory* **52**(2), 489–509 (2006).
48. Leuschner, J., Schmidt, M., Baguer, D. O., Maaß, P. The LoDoPaB-CT dataset: A benchmark dataset for low-dose CT reconstruction methods (2019). arXiv preprint [arXiv:1910.01113](https://arxiv.org/abs/1910.01113).
49. Jiang, C., Zhang, Q., Fan, R. & Hu, Z. Super-resolution CT image reconstruction based on dictionary learning and sparse representation. *Sci. Rep.* **8**(1), 1–10 (2018).
50. Xu, Q. *et al.* Low-dose X-ray CT reconstruction via dictionary learning. *IEEE Trans. Med. Imaging* **31**(9), 1682–1697 (2012).
51. Ring, W. Structural properties of solutions to total variation regularization problems. *ESAIM Math. Model. Numer. Anal.* **34**(4), 799–810 (2000).
52. Chan, T., Marquina, A. & Mulet, P. High-order total variation-based image restoration. *SIAM J. Sci. Comput.* **22**(2), 503–516 (2000).
53. Kamvar, S., Olszewski, M., Reinsberg, R. Celos: A multi-asset cryptographic protocol for decentralized social payments. *White Paper* (2019). [storage.googleapis.com/celo-whitepapers/Celo A Multi Asset Cryptographic Protocol for Decentralized Social Payments.pdf](https://storage.googleapis.com/celo-whitepapers/Celo%20A%20Multi%20Asset%20Cryptographic%20Protocol%20for%20Decentralized%20Social%20Payments.pdf).
54. Project, M. The Maker Protocol: MakerDAO's Multi-Collateral Dai (MCD) System (2020). *White Paper*. [storage.googleapis.com/celo-whitepapers/Celo A Multi Asset Cryptographic Protocol for Decentralized Social Payments.pdf](https://storage.googleapis.com/celo-whitepapers/Celo%20A%20Multi%20Asset%20Cryptographic%20Protocol%20for%20Decentralized%20Social%20Payments.pdf).
55. Zhang, Y., Chen, X., Park, D. Formal specification of constant product (xy = k) market maker model and implementation. *White Paper* (2018).
56. Warren, W., Bandeau, A. 0x: An open protocol for decentralized exchange on the Ethereum blockchain. *White Paper* (2017). github.com/0xProject/whitepaper.
57. Hertzog, E., Benartzi, G., Benartzi, G. Bancor protocol. *White Paper* (2017). storage.googleapis.com/website-bancor/2018/04/01ba8253-bancor_protocol_whitepaper_en.pdf (accessed 24 Apr 2022).
58. Werner, S. M., Perez, D., Gudgeon, L., Klages-Mundt, A., Harz, D., Knottenbelt, W. J. Sok: Decentralized finance (defi) (2021). arXiv preprint [arXiv:2101.08778](https://arxiv.org/abs/2101.08778).
59. Schär, F. *Decentralized finance: On blockchain-and smart contract-based financial markets* (FRB of St. Louis Review, 2021).
60. Angeris, G., Chitra, T. Improved price oracles: Constant function market makers. In *Proceedings of the 2nd ACM Conference on Advances in Financial Technologies*, 80–91 (2020).

61. Angeris, G., Agrawal, A., Evans, A., Chitra, T., Boyd, S. Constant function market makers: Multi-asset trades via convex optimization (2021). arXiv preprint [arXiv:2107.12484](https://arxiv.org/abs/2107.12484).
62. Makarov, I. & Schoar, A. Trading and arbitrage in cryptocurrency markets. *J. Financ. Econ.* **135**(2), 293–319 (2020).
63. Angeris, G., Chitra, T., Evans, A., Boyd, S. Optimal routing for constant function market makers (2021).
64. Daian, P., Goldfeder, S., Kell, T., Li, Y., Zhao, X., Bentov, I., Breidenbach, L., Juels, A. Flash boys 2.0: Frontrunning, transaction reordering, and consensus instability in decentralized exchanges (2019). arXiv preprint [arXiv:1904.05234](https://arxiv.org/abs/1904.05234).
65. Davis, D. & Yin, W. A three-operator splitting scheme and its optimization applications. *Set-Valued Var. Anal.* **25**(4), 829–858 (2017).
66. Raissi, M., Perdikaris, P. & Karniadakis, G. E. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *J. Comput. Phys.* **378**, 686–707 (2019).
67. Ruthotto, L., Osher, S. J., Li, W., Nurbekyan, L. & Fung, S. W. A machine learning framework for solving high-dimensional mean field game and mean field control problems. *Proc. Natl. Acad. Sci.* **117**(17), 9183–9193 (2020).
68. Lin, A. T., Fung, S. W., Li, W., Nurbekyan, L., Osher, S. J. Alternating the population and control neural networks to solve high-dimensional stochastic mean-field games. *Proce. Natl. Acad. Sci.* **118**(31) (2021).

Acknowledgements

The authors thank Wotao Yin, Stanley Osher, Daniel McKenzie, Qiuwei Li, and Luis Tenorio for many fruitful discussions. Howard Heaton and Samy Wu Fung were supported by AFOSR MURI FA9550-18-1-0502 and ONR grants: N00014-18-1-2527, N00014-20-1-2093, and N00014-20-1-2787. Samy Wu Fung was also partially funded by the National Science Foundation award number DMS-2309810.

Author contributions

H.H. and S.W.F. performed the research and wrote the manuscript. H.H. created each figure except for Fig. 6. S.W.F. created Fig. 6.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-36249-3>.

Correspondence and requests for materials should be addressed to H.H. or S.W.F.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023