From PDE-Constrained Optimization to DNNs

Ali Hamza Abidi, Syed & Andreas Mang Department of Mathematics, University of Houston, Houston, TX, USA

Teaser: Our goal was the design and analysis of effective numerical schemes for training deep neuronal networks (**DNN**s) based on optimal control formulations.

In the present work we explore numerical methods inspired by optimal control theory to train image classifiers [1]. In a first step, we consider a prototypical formulation of a variational optimization problem governed by an elliptic dynamical system [2]. We will discuss the numerical treatment and study some of the mathematical operators. Subsequently, we present the optimal control formulation for training DNNs and derive some expressions for the associated optimality conditions. In our future work, we plan to extend these optimality conditions and device a numerical scheme for the DNN training problem, similar to the scheme developed for the prototypical problem.

Optimal Control Problems

In the present work we study numerical methods for solving problems of the general form [3]

minimize $dist(u_{pred}, y^{\delta}) + reg(w)$ subject to $\mathcal{C}(w, u) = 0$.

Here, $u \in \mathcal{U}$ is the state variable, $w \in \mathcal{W}$ is the control variable, and $c: \mathcal{W} \times \mathcal{U} \to \mathcal{Q}, \ \mathcal{C}(w, u) \coloneqq \mathcal{A}(w)u - q$ is a dynamical system, i.e., a partial differential equation (PDE). The functional dist measures the discrepancy between the predicted state $u_{\text{pred}} := \mathcal{A}^{-1}(w)q$ and the observable (data) y^{δ} , with parameter-to-obervation map $\mathcal{F}(w) := \mathcal{A}^{-1}(w)q$. The functional reg is a Tikhonov regularization model introduced to alleviate issues with the ill-posedness of the optimization problem stated above [4].

Prototypical Example

We consider the following PDE-constrained optimization problem [2]

$$\begin{array}{l} \underset{u \in \mathcal{U}, \ w \in \mathcal{W}}{\text{minimize}} \quad \frac{1}{2} \int_{\Omega} (u - y^{\delta})^2 \mathrm{d}x + \frac{\alpha}{2} \|w\|_{\mathcal{W}}^2 \\ \text{subject to} \quad \mathcal{C}(w, u) = \nabla \cdot w \nabla u - q = 0 \end{array}$$

subject to Dirichlet boundary conditions u = 0, with gradient operator ∇ , divergence operator $\nabla \cdot$, and regularization parameter $\alpha > 0$. We consider q(x) = 1 if $x_1 \in (0.4, 0.6), x_2 \in (0.4.0.6)$ and q(x) = 0 otherwise, the parameter function $w_{true}(x) = 2.2 + 1$ $2\sin(\pi x_1)\sin(\pi x_2)$, and $y^{\delta} = \mathcal{A}(w_{\text{true}})^{-1}q$ (see Figure 1).



Figure 1: Visualization of the problem data. From left to right: The right hand side q, the parameter function w_{true} , and the solution of the forward problem for this input data y^{δ} .

We discretize the functions u, w, and q on a cell-centered grid of size $n_1 \times n_2$. We consider a finite-volume discretization and arrange all variables in lexicographical ordering to obtain vectors $\mathbf{u} \in \mathbb{R}^n$, $\mathbf{w} \in \mathbb{R}^n$, $\mathbf{q} \in \mathbb{R}^n$, with $n = n_1 n_2$. The discrete optimization problem is given by

$$\underset{\mathbf{u}\in\mathbb{R}^{n},\mathbf{w}\in\mathbb{R}^{n}}{\text{minimize}} \quad \frac{h}{2}\|\mathbf{u}-\mathbf{y}^{\delta}\|_{2}^{2} + \frac{\alpha h}{2}\|\mathbf{L}\mathbf{w}\|_{2}^{2} \text{ subject to } \mathbf{c}(\mathbf{u},\mathbf{w})$$

Here, $\mathbf{L} \in \mathbb{R}^{n \times n}$ denotes the regularization operator and $\mathbf{c}(\mathbf{u}, \mathbf{w}) = \mathbf{c}(\mathbf{u}, \mathbf{w})$ DSG $\mathbf{u} - \mathbf{q}$, DSG = DIV **S** GRAD $\in \mathbb{R}^{n,n}$, where

$$\mathsf{GRAD} = \begin{bmatrix} \boldsymbol{I}_2 \otimes \boldsymbol{D}_1 \\ \boldsymbol{D}_2 \otimes \boldsymbol{I}_1 \end{bmatrix} \in \mathbb{R}^{m,n}$$

is the discretized gradient operator with $m = (n_1+1)n_2 + n_1$ The matrix $\boldsymbol{S} = \text{diag}(\boldsymbol{s}) \in \mathbb{R}^{m,m}$, $\boldsymbol{s} = \boldsymbol{e}_m \oslash (\boldsymbol{A}(\boldsymbol{e}_n \oslash \boldsymbol{w}))$ $\boldsymbol{e}_n = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^{\mathsf{I}} \in \mathbb{R}^n, \oslash : \mathbb{R}^n \to \mathbb{R}^n, \text{ represents the discr}$ of *w* via harmonic averaging, and

$$\boldsymbol{A} = \begin{bmatrix} \boldsymbol{I}_2 \otimes \boldsymbol{A}_1 \\ \boldsymbol{A}_2 \otimes \boldsymbol{I}_1 \end{bmatrix} \in \mathbb{R}^{m,n}$$

with $\mathbf{A}_i \in \mathbb{R}^{n_i+1,n_i}$, $\mathbf{I}_i = \text{diag}(1 \dots 1) \in \mathbb{R}^{n_i,n_i}$, i = 1, 2. The matrix

$$\mathsf{DIV} = \begin{bmatrix} \mathbf{I}_2 \otimes \mathbf{D}_1 & \mathbf{D}_2 \otimes \mathbf{I}_1 \end{bmatrix} \mathbb{R}^{n,m}$$

represents the divergence operator with I_i as defined above and $D_i \in R^{n_i+1,n_i}$ is the one-dimensional derivative operator along the coordinate direction x_i , i = 1, 2. The sensitivity J(w) and its adjoint are given by

$$\boldsymbol{J} \coloneqq \boldsymbol{J}(\boldsymbol{w}) = -\boldsymbol{Q}\boldsymbol{C}_u^{-1}\boldsymbol{C}_w$$
 and $\boldsymbol{J}^\mathsf{T} \coloneqq \boldsymbol{J}(\boldsymbol{w})^\mathsf{T} = -\boldsymbol{C}_w^\mathsf{T}\boldsymbol{C}_u^{-\mathsf{T}}\boldsymbol{Q}^\mathsf{T}$,

respectively, with derivative matrices $C_w := d_w c(u, w) \in \mathbb{R}^{n,n}$ and $C_{\mu} := d_{\mu}c(\boldsymbol{u}, \boldsymbol{m}) \in \mathbb{R}^{n,n}$ (see Figures 2 & 3). Given the PDE constraint c(u, w) = 0 and a candidate model w, we can solve u = u(w). With this, we can eliminate the constraint from the optimization problem and reformulate the equality constrained problem as an unconstrained problem of the form

$$\underset{\boldsymbol{u}\in\mathbb{R}^n}{\text{minimize}} \quad \left\{ f(\boldsymbol{w}) \coloneqq \frac{h}{2} \| \boldsymbol{u}(\boldsymbol{w}) - \boldsymbol{y}^{\delta} \|_2^2 + \frac{\alpha h}{2} \| \boldsymbol{L} \boldsymbol{w} \|_2^2 \right\}.$$

A necessary condition for a minimizer of this problem is given by the non-linear system $d_w f(w^*) = h J^T (u(w^*) - y^\delta) + \alpha h L^T L w^* = 0.$

Numerical Optimization

We have developed a matrix-free Newton–Krylov method for its so-Iution. We use an *iterative line search scheme* of the form

$$\boldsymbol{w}_{k+1} = \boldsymbol{w}_k - \mu_k \boldsymbol{B}_k d_w f(\boldsymbol{w}_k), \quad k = 1, 2, \dots$$

Here, $k \in \mathbb{N}$ is the iteration index and $\mu_k \in (0, 1]$ is determined using a backtracking line search [5]. The search direction is given by $\mathbf{s}_k := -\mathbf{B}_k d_w f(\mathbf{w}_k)$. We consider Newton's method with $\boldsymbol{B}_{k} = \boldsymbol{H}^{-1}, \boldsymbol{H} := d_{ww}f(\boldsymbol{w}) = h\boldsymbol{J}(\boldsymbol{w})^{\mathsf{T}}\boldsymbol{J}(\boldsymbol{w}) + d_{w}(\boldsymbol{J}(\boldsymbol{w})^{\mathsf{T}}\boldsymbol{r}) + \alpha h\boldsymbol{L}^{\mathsf{T}}\boldsymbol{L},$ where $r := u(w) - y^{\delta}$ is constant. We invert the Hessian matrix using a matrix-free, conjugate gradient method with a superlinear forcing sequence. As a stopping criterion, we consider the relative reduction of the norm of the gradient.



Figure 2: Left: Sparsity pattern of the Jacobians C_u (1210 non-zero entries) and C_u (1220 non-zero entries) for d = 2, $n_i = 16$. Middle: Sparsity pattern (for d = 2, $n_i = 16$) and visualization of the entries (for d = 2, $n_i = 8$) of the inverse of C_u (65 500 non-zero entries; the matrix is dense). Right: Visualization of the entries of the sensitivities **J** and $J^{\dagger}J$ for d = 2, $n_i = 8.$

=**0**

$$(n_2+1).$$

 $)) \in \mathbb{R}^m,$ retization

(1)





 $i \in \{1, 2, 4, 32, 64, 128\}.$

Optimal Control Formulation for DNNs The **optimal control formulation** for training a deep neural network is given by [1]

minimize dist(
$$\boldsymbol{C}_{pred}, \boldsymbol{C}$$
) + α reg($\boldsymbol{W}, \boldsymbol{\mu}, \{\boldsymbol{K}_i\}_{i=1}^n, \{b_i\}$
subject to $\boldsymbol{Y}_{j+1} = \boldsymbol{Y}_j + h\sigma(\boldsymbol{Y}_j\boldsymbol{K}_j + \boldsymbol{b}_j)$,

 $j = 0, 1, \dots, n-1$. Here, dist : $\mathbb{R}^{s,m} \times \mathbb{R}^{s,m} \to \mathbb{R}$ measures the discrepancy between the predicted classification $C_{pred} \in \mathbb{R}^{s,m}$ and the labels (data) $\boldsymbol{C} \in \{0,1\}^{s,m}$. The unknowns $\boldsymbol{\Phi}$ of the optimization problem (2) are the weights $K_i \in \mathbb{R}^{n,n}$ and biases b_i of the 'ResNET' forward propagation and the weights $\boldsymbol{W} \in \mathbb{R}^{n,m}$ and biases $\boldsymbol{\mu} \in \mathbb{R}^m$ that parameterize the classifier. Consequently, $\mathbf{\Phi} := \{ \mathbf{W}, \mathbf{\mu}, \{ \mathbf{K}_i \}_{i=1}^n, \{ b_i \}_{i=1}^n \}$. The prediction \mathbf{C}_{pred} is computed according to $C_{pred} = g(Y_n W + e_s \otimes \mu)$, where Y_n is the final state computed by solving the forward propagation, $e_s := (1, \ldots, 1)^{\top} \in \mathbb{R}^s$, and $g: \mathbb{R}^{s,m} \to \mathbb{R}^{s,m}$ is the so called hypothesis function. The partial derivatives of the loss function $\|g(\mathbf{Y}_n\mathbf{W} + \mathbf{e}_s \otimes \mu) - \mathbf{C}\|_F^2/s$ are given by $d_W \|g(\mathbf{Y}_n \mathbf{W}) - \mathbf{C}\|_F^2 = 2\mathbf{Y}_n^T (g'(\mathbf{Y}_n \mathbf{W}) \odot (g(\mathbf{Y}_n \mathbf{W}) - \mathbf{C}))/s$ and $d_Y \|g(\mathbf{Y}_n \mathbf{W}) - \mathbf{C}\|_F^2 = 2(g'(\mathbf{Y}_n \mathbf{W}) \odot (g(\mathbf{Y}_n \mathbf{W}) - \mathbf{C}))\mathbf{W}^\top / s$. If we consider the regularization operators $reg(\mathbf{K}) = \|\mathbf{K}\|_F^2/2 = tr(\mathbf{K}\mathbf{K}^{\mathsf{T}})/h$ and $\operatorname{reg}(\boldsymbol{b}) = \|\boldsymbol{b}\|_2^2/h$ we obtain the derivatives $d_K = \boldsymbol{K}$ and $d_b = \boldsymbol{b}$, respectively. The derivatives of the dynamical system (derivation of the sensitivities) is more complicated and forms the basis of our current work.

Conclusions

We have explored the implementation of an optimization framework for optimal control problems governed by PDEs. We have developed a matrix-free, Newton-Krylov method globalized by an Armijo *linesearch* for numerical optimization. We have started to derive optimality conditions for the optimal control problem for training DNNs described in [1]. The derivation of the sensitivities for this problem forms the basis of our current work.

References

- **1.** E. Haber & L. Ruthotto, Stable Architectures for Deep Neural Networks, Inverse Problems 34 014004, 2017.
- **2.** E. Haber, Computational Methods in Geophysical Electromagnetics, SIAM, 2015.
- **3.** F. Tröltzsch, Optimal Control of Partial Differential Equations: Theory, Methods, and Applications, American Mathematical Society, 2010.
- 4. H. W. Engl, M. Hanke, A. Neubauer, Regularization of Inverse Problems, Springer, 2000
- 5. J. Nocedal & S. Wright, Numerical Optimization, Springer Science, 1999.









