

On Obtaining Sparse Semantic Solutions for Inverse Problems, Control, and Neural Network Training

David A. B. Hyde^{a,*}, Michael Bao^{b,**}, Ronald Fedkiw^{b,c,**}

^a*UCLA Mathematics Department, Box 951555, Los Angeles, CA 90095-1555, United States*

^b*Stanford University, 353 Jane Stanford Way, Gates Computer Science Room 207, Stanford, CA 94305, United States*

^c*Epic Games, Inc., 620 Crossroads Boulevard, Cary, NC 27518, United States*

Abstract

Modern-day techniques for designing neural network architectures are highly reliant on trial and error, heuristics, and so-called best practices, without much rigorous justification. After choosing a network architecture, an energy function (or loss) is minimized, choosing from a wide variety of optimization and regularization methods. Given the ad-hoc nature of network architecture design, it would be useful if the optimization led to a sparse solution so that one could ascertain the importance or unimportance of various parts of the network architecture. Of course, historically, sparsity has always been a useful notion for inverse problems where researchers often prefer the L_1 norm over L_2 . Similarly for control, one often includes the control variables in the objective function in order to minimize their efforts. Motivated by the design and training of neural networks, we propose a novel column space search approach that emphasizes the data over the model, as well as a novel iterative Levenberg-Marquardt algorithm that smoothly converges to a regularized SVD as opposed to the abrupt truncation inherent to PCA. In the case of our iterative Levenberg-Marquardt algorithm, it suffices to consider only the linearized subproblem in order to verify our claims. However, the claims we make about our novel column space search approach require examining the impact of the solution method for the linearized subproblem on the fully nonlinear original problem; thus, we consider a complex real-world inverse problem (determining facial expressions from RGB images).

Keywords: Machine learning, Levenberg-Marquardt, principal component analysis, column space search, coordinate descent

1. Introduction

The current age of deep learning began (at least according to the Turing Award committee¹) with works addressing problems such as object classification [73, 82], reading handwritten digits and documents [80, 81, 79], and speech recognition and natural language tasks [11, 100]. Although models based on traditional scientific first principles do not exist for these sorts of problems, the underlying machine learning methods have been permeating into various scientific communities, including computational physics [69, 51, 49, 123, 48, 89, 113, 115, 132]. Perhaps the main difference between the use of machine learning for customizing advertisements [19, 55], online dating [36, 99], or self-driving cars [16, 67] and its use in computational physics is that our community has developed a fairly reasonable scientific and mathematical understanding of many of the problems of interest via a combination of theoretical, experimental, and computational approaches, especially as opposed to the ad-hoc data-driven nature of popular machine learning application areas. Unfortunately, ad-hoc approaches leave neural networks wide open to adversarial attacks [65, 3, 125], which does not bode well for predictive numerical capabilities. Therefore, one goal of our community (and perhaps contribution) would be to better understand neural network architectures in order to provide a more thorough and rigorous approach to designing them, similar to the contributions that the applied mathematics

*dabh@math.ucla.edu, UCLA

**mikebao@stanford.edu, fedkiw@cs.stanford.edu, Stanford University

¹<https://awards.acm.org/about/2018-turing>

16 community made to finite element simulation, e.g. reformulating spring and beam elements as basis functions
17 [141, 142].

18 Techniques used in modeling and training neural networks are highly related to well-studied approaches
19 for inverse problems and control. To understand some of the differences between inverse problems, control,
20 and training neural networks, consider $Y = f(X; C)$, with input X , output Y , and function parameters C .
21 In a typical inverse problem, one is given Y and aims to find an X that produces Y . Poor conditioning of the
22 function f or noise in the given/desired output Y can lead to spurious information contained in X . Thus,
23 various regularization approaches may be used to ascertain an X with a high signal-to-noise ratio, see for
24 example [26, 92, 131, 140] and the more general references [42, 39, 12]. In the control problem, X and Y are
25 both given, and the goal is to ascertain some subset of the function parameters C that allows one to coerce
26 X toward Y . Typically, most of f is a well-known function, such as the Navier-Stokes equations, and thus
27 the added controls should have a light/minimal touch; therefore, they are often included in the objective
28 function so that their magnitude/effort is minimized. This too is regularization, and needs to be done wisely
29 so that minimizing controls does not prevent one from hitting the target (while still considering signal-to-
30 noise ratio, etc.), see e.g. [70, 2, 118]. When considering neural networks, the function f is almost entirely
31 ad-hoc, and one does not know which parameters might have physicality and which are more arbitrary.
32 Thus, it becomes even more important to consider careful regularization with the hope that some of the
33 coefficients will dominate others, providing some insight into which portions of the network architecture may
34 have some basis in first principles as opposed to which may be considered for removal/replacement, see e.g.
35 [137, 95, 138, 133, 54, 90, 60, 93, 102, 134, 56, 109, 85, 4, 57, 58]. Because so little is known about f , neural
36 networks cannot proceed with one input X and one output Y as can a control problem. Instead, one requires
37 a family of given $(X; Y)$ pairs called training data, before an attempt to identify the function coefficients
38 C can be made. Methods for formulating and optimizing neural networks are typically significantly more
39 rudimentary and ad-hoc than those designed for inverse and control problems, relying on simple methods such
40 as gradient descent and stochastic gradient descent (SGD) or ordinary differential equation discretizations
41 of gradient flow, such as Adam [72], AdaGrad [37], Nesterov [103], momentum methods [114, 126], etc.
42 [117, 17, 53].

43 The process of network architecture design is often motivated by heuristics that hinder the ability to
44 subsequently train the network and find suitable coefficients. For example, the “all or none” property of bio-
45 logical neurons leads to discontinuous functions with identically zero derivatives almost everywhere, which is
46 disastrous for optimization/training [98]. The idea that biological neurons fire with increased frequency for
47 stronger signals leads to piecewise linear functions with discontinuous derivatives, also problematic for opti-
48 mization. These Heaviside and rectifier/ReLU [53] models require smoothing before they can subsequently
49 be used with numerical optimization. It seems quite dubious to design and analyze non-smooth network
50 architectures that are later smoothed in the first significant digit when deployed in practice, especially given
51 the nuances exposed by the numerical analysis community regarding the differences between continuous and
52 discrete formulations (e.g. [61]) even when such occurs only at the level of machine precision (the 7th or 15th
53 decimal place). This motivates our aim to better utilize various approaches to regularization and sparsity to
54 ascertain the importance of various components of the network architecture.

55 In Section 2, we introduce a suitably complex model problem for demonstrating the numerical methods
56 presented in the paper: determining facial expression from RGB images. This problem is both algorithmi-
57 cally challenging and grounded in physics, meaning that we can attempt to develop algorithms which find
58 semantically meaningful solutions in terms of known physical and anatomical properties. In Section 3, we
59 outline a general framework for optimization, showing how neural network training is recast as a nonlinear
60 optimization problem. We highlight various approximations made in practice during this process, such as
61 (sometimes drastic) approximations to the Hessian and Jacobian. As is typical, rank-one updates are dis-
62 cussed, which motivates the singular value decomposition (SVD) and principal component analysis (PCA),
63 both of which are used in subsequent sections. Section 4 presents a novel iterative Levenberg-Marquardt
64 [83, 97] scheme that is shown by proof and experiment to converge *smoothly* (and monotonically) to a regular-
65 ized SVD, unlike the truncation typical of a PCA approach. Section 5 presents a novel column space search
66 technique that focuses more on the data term than the model, again an improvement over PCA. Moreover,
67 we explain how column space search enables the discovery of sparse and semantically meaningful solutions
68 to fully nonlinear optimization problems. To demonstrate this experimentally, we compare to alternative
69 strategies such as Dogleg [112, 94] and BFGS with L_2 or soft L_1 regularizers.

70 **2. Tackling Complex Real-World Inverse Problems for Faces**

71 We chose a fairly complex model problem which is still cutting-edge in order to illustrate the need
 72 for robust and efficient approaches. Specifically, we consider an inverse problem where a two-dimensional
 73 RGB image of a human face is processed to determine facial expressions in terms of a three-dimensional
 74 parameterized model with a semantic, anatomical basis. This inverse problem is useful throughout industries
 75 such as medicine, surveillance, intelligence gathering, entertainment, etc. Similar to many other complex
 76 processes, one can understand the problem via a pipeline with various function layers, see Figure 1.

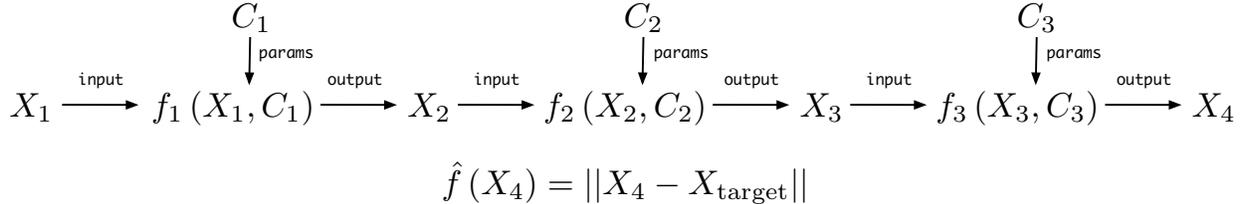


Figure 1: Multiple layers of functions f_i map an initial vector of inputs X_1 to a final output X_4 , which is evaluated with an objective function \hat{f} . Vectors of parameters C_i may either be prescribed or be determined via experimentation or neural network training.

77 The inverse problem seeks to find an X_1 that outputs an X_4 as close to X_{target} as possible, i.e. minimizing
 78 $\hat{f}(X_4)$, using regularization to combat noise and overfitting when necessary. Using classical optimization,
 79 this requires differentiation that can be expressed as

$$\frac{\partial \hat{f}}{\partial X_1} = \frac{\partial \hat{f}}{\partial X_4} \frac{\partial f_3(X_3, C_3)}{\partial X_3} \frac{\partial f_2(X_2, C_2)}{\partial X_2} \frac{\partial f_1(X_1, C_1)}{\partial X_1}, \quad (1)$$

80 implying that every function layer requires differentiability with respect to its inputs. Now suppose that
 81 $f_2(X_2; C_2)$ represented a neural network layer that needs to be trained in order to ascertain reasonable
 82 parameters C_2 . In order to do this, one would consider a large number K of known training pairs of the
 83 form $(X_1^k; X_{\text{target}}^k)$; however, notationally, one may stack all the training pairs into a single X_1 and X_{target} ,
 84 at least conceptually (for the sake of exposition). Then the required differentiation is

$$\frac{\partial \hat{f}(X_4)}{\partial C_2} = \frac{\partial \hat{f}}{\partial X_4} \frac{\partial f_3(X_3, C_3)}{\partial X_3} \frac{\partial f_2(X_2, C_2)}{\partial C_2}, \quad (2)$$

85 highlighting the notable differences as compared to an inverse problem. Firstly, any pre-process, such as f_1
 86 here, does not require differentiability and can utilize any known procedural methods. In fact, one might
 87 use an f_1 based on first principles aiming to solve the problem outright, and then supplement the results
 88 with the composition of f_2 and f_3 in order to better match real-world data. This means that data-driven
 89 neural network approaches may be added on top of any existing codebase, whether it is differentiable or
 90 not. Secondly, any post-process for the neural network, such as f_3 , only requires as much differentiability as
 91 would be required for f_3 if it were included in a typical inverse problem. Thirdly, the neural network itself,
 92 f_2 , does not require the usual differentiability inherent to inverse problems, but only requires differentiability
 93 with respect to its parameters C_2 .

94 Most facial pipelines take as input a set of parameters that govern the shape/geometry of a three-
 95 dimensional face, as given by triangle vertex positions. For example, a blendshape facial rig (see e.g. [84])
 96 describes how a face is deformed from a neutral rest state n in terms of a linear combination of basis
 97 facial shapes, e.g. semantic basis vectors which represent particular expressions such as “smile” or “yawn.”
 98 The basis facial shapes are often acquired using dense performance capture (see e.g. [9, 10, 20, 50]) or via
 99 sculpting by an artist/modeler [30, 75]. A typical high-quality blendshape rig contains hundreds of basis
 100 shapes corresponding to different expressions between which one can interpolate (see e.g. [29]). Once a
 101 blendshape model is obtained, stacking each blendshape into a column of a matrix B allows one to define
 102 the facial geometry as $n + Bb(w)$, where the vector b contains a degree of freedom for each blendshape
 103 and w represents a set of meaningful controls ($b(w)$ may be nonlinear, but should be smooth). In order to

104 avoid linearized rotation artifacts due to rotational jaw motion [32, 121, 143], one typically hybridizes the
 105 linear blendshape system with skinning/enveloping (see e.g. [96, 76]), which blends together the nonlinear
 106 six-degree-of-freedom rigid body transformation from the skull and jaw. Each triangle vertex is assigned
 107 weights that dictate the relative influence of the skull and jaw such that vertices far from the jaw move with
 108 the skull, vertices far from the skull move with the jaw, and vertices in between move in a blended fashion.
 109 This can be written compactly as a matrix $T(j(w))$, where the controls w drive the six-degree-of-freedom
 110 rigid body offset j of the jaw from the skull and T assembles all the transformations and weights so that one
 111 may write

$$x(w) = T(j(w))(n + Bb(w)), \quad (3)$$

112 where $x(w)$ are the triangle vertex positions of the face surface. Importantly, as long as the dependencies in
 113 Equation 3 are chosen carefully (in a smooth enough manner), then x is differentiable with respect to w .

114 As an alternative to blendshape approaches, one can construct an anatomically motivated finite element
 115 facial model based on soft tissue, musculature, and underlying skeletal structures (see e.g. [121, 122]). In
 116 [121], the vertex positions are differentiable as a function of the muscle activations and jaw parameters, and
 117 the authors used this differentiability to solve inverse problems. However, since anatomical facial models
 118 rely on MRI, CT scans, etc., it is difficult to make an accurate model; therefore, [121, 122] struggled to
 119 express the wide variety of shapes possible with a blendshape system. Thus, [31] augmented the results of
 120 [121] using a three-dimensional morphing process in order to derive target locations for muscles. Although
 121 the method proposed in [31] regains the expressivity of a blendshape system, their morphing process lacked
 122 differentiability. Later, [7] noted that the morphing process could be made differentiable, but that this would
 123 require a mapping from each surface vertex to all other affected vertices in the simulation mesh, which is
 124 quadratic complexity and thus impractical. Instead, [7] parameterized the morph with a standard blendshape
 125 system, so that the parameters b drive the morph, resulting in linear complexity. This was implemented
 126 (in [7]) by simulating the anatomical mesh for each blendshape (using the morphing from [31]) in order to
 127 create the muscle shapes needed in order to define a blendshape system for the muscles themselves; then,
 128 the three-dimensional target shape of each muscle can be specified by the parameters w , with each muscle
 129 tetrahedron vertex $x(w)$ defined along the lines of Equation 3. Manipulating w determines a blendshape
 130 for each muscle, which is then targeted with the anatomical finite element simulation from [31] and [121],
 131 see Figure 2. Notably, the resulting scheme is fully differentiable and hence can be used to solve inverse
 132 problems.

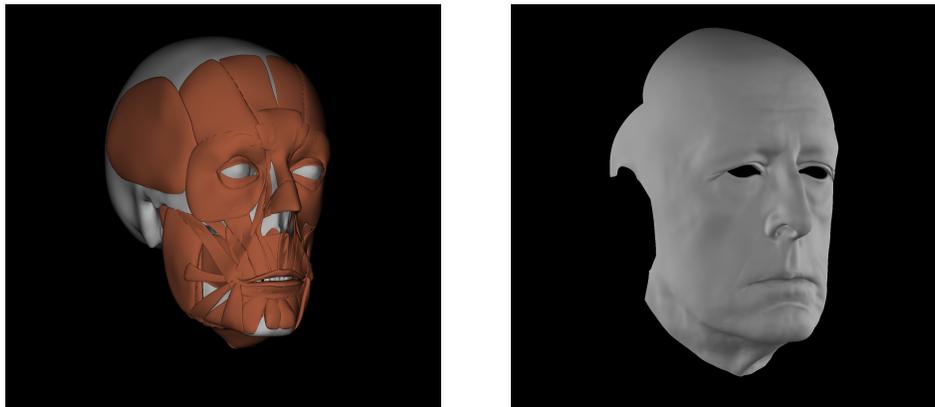


Figure 2: (Left) Skull and jaw (gray) with anatomical muscle shapes (red). (Right) Corresponding surface of the tetrahedral finite element mesh simulated from the targeted muscle shapes.

133 Given target geometry for the three-dimensional face surface, one can specify an energy that minimizes
 134 the distance between the target and the parameterized model, and then solve an inverse problem for the
 135 parameters w that drive the b and j for the muscle blendshape system, which in turn drives the quasistatic
 136 finite element simulation of [121] augmented by [31] in order to match the target (see [7] for details). In order
 137 to match a two-dimensional RGB image, one needs to render the resulting geometry with a differentiable
 138 renderer along the lines of [88, 91] and utilize an energy that considers the difference in pixel colors. Then,
 139 one can solve an inverse problem for the controls w that drive the blendshape muscles which in turn drive the

140 finite element simulation which results in the surface mesh that is rendered into pixel colors that minimize
 141 the energy. Unfortunately, as shown in Figure 3 Left, differentiable renderers don't typically have the same
 142 quality as a photorealistic renderer or photograph, so aiming to match pixel colors is overly optimistic. To
 143 overcome this limitation, [6] proposed processing both the image and the differentiable render with a pre-
 144 existing/widespread face landmark detector neural network, such as 2D/3D-FAN [24] (see Figure 3). These
 145 networks were trained with vast amounts of hand-labeled data so that they could find keypoint/landmark
 146 positions from images regardless of texture, geometry, shading, lighting, shadows, etc. As such, the poor
 147 rendering quality of a differentiable renderer is also serendipitously ignored by these neural networks. In
 148 summary, [6] utilizes an energy that computes the difference between keypoints/landmarks, and solving the
 149 inverse problem requires differentiating through the keypoint detector neural network (2D/3D-FAN), the
 150 differentiable renderer, and the quasistatic finite element muscle simulation.

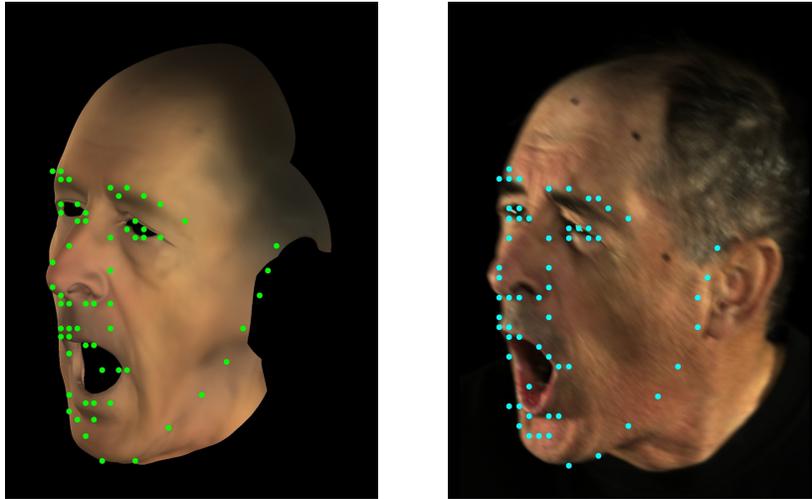


Figure 3: Results of a machine-learning based facial keypoint detector such as 2D/3D-FAN [24] on a synthetic render (left) as well as the corresponding photograph (right).

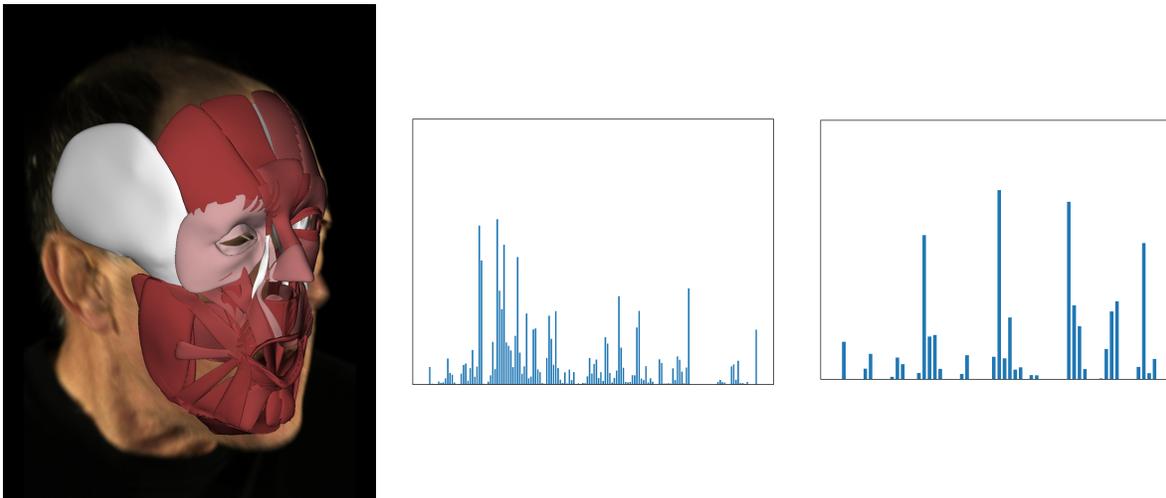


Figure 4: (Left) A pose of the hybridized muscle system of [7] depicted on top of the corresponding target RGB photograph. Whiter muscle shapes correspond to more activated muscles. (Middle) The corresponding blendshape weights. (Right) The corresponding muscle activations.

151 Figure 4 illustrates results typical of this process. Figure 4 (Middle) shows the value of b along the vertical
 152 axis for each blendshape on the horizontal axis. In spite of the expression in Figure 4 (Left) being not that

153 complex, many blendshapes have non-zero values; worse yet, successive frames in a video produce noisy
 154 uncorrelated blendshape values that are hard to interpret as meaningful semantic information. In contrast,
 155 Figure 4 (Right) illustrates that the muscle activations are sparser and more indicative of the image; in
 156 fact, successive frames tend to be highly correlated, allowing one to separate semantic information from
 157 noise. Generally speaking, sparse semantic solutions to inverse problems are obviously preferred over dense,
 158 noisy, temporally uncorrelated results. This goal of ascertaining sparser semantic information motivates our
 159 considerations throughout the rest of the paper.

160 3. Optimization Framework

161 Whether it be the search for viable inputs for an inverse problem, minimizing some measure of effort for
 162 a control problem, or the determination of network architecture parameters that allow a neural network to
 163 well-match training data, these problems all take the form of an optimization minimizing a cost function
 164 $\hat{f}(c)$ over parameters c . Importantly, one typically has certain conditions in mind to which c should be
 165 subject. For example, one might want c close to a prior/initial guess, one might desire the norm of c to
 166 be small, and/or one might want c sparse so that it carries interpretable semantic meaning. In particular,
 167 as noted above, it would be useful if neural network training resulted in a sparse c in order to identify
 168 unnecessary/unimportant components of the network architecture.

169 Either in the absence of constraints or with constraints and suitable Lagrange multipliers, the minima of
 170 \hat{f} occur at critical points where the (column vector) Jacobian $F(c) = J_{\hat{f}}^T(c) = \nabla \hat{f}(c) = 0$. Since $F(c) = 0$
 171 is (generally) a nonlinear system of equations, one typically linearizes the system by taking the first two
 172 terms of the Taylor expansion about a point c^* , $F(c) \approx F(c^*) + F'(c)(c - c^*)$ where $F'(c) = J_F(c) = H_{\hat{f}}^T(c)$
 173 is the transpose of the Hessian of \hat{f} . Newton’s method uses this relationship to write $F(c^{q+1}) - F(c^q) =$
 174 $F'(c^q) \Delta c^q$, where $\Delta c^q = c^{q+1} - c^q$ and q represents the current iteration. Then, one solves the linear system
 175 $F'(c^q) \Delta c^q = \beta F(c^q) - F(c^q)$ to update $c^{q+1} = c^q + \Delta c^q$ where $\beta \in [0, 1)$, and using $\beta \neq 0$ more slowly
 176 shrinks $F(c^q)$ towards 0. Alternatively, one can utilize Δc^q merely as a search direction and subsequently
 177 employ a number of one-dimensional approaches, e.g. bisection search, golden section search, etc.

178 While Newton’s method and similar techniques are reasonably well-justified and often converge well
 179 in practice, they depend on access to various derivatives of the cost function $\hat{f}(c)$. To compute these
 180 derivatives, one may utilize symbolic/analytic differentiation, finite differences, or automatic differentiation
 181 (e.g. backpropagation). Automatic differentiation is often preferred in the context of training neural networks
 182 both because of its ease of implementation as well as its availability via various software packages (e.g.
 183 Tensorflow [1], Caffe [71], PyTorch [110], Theano [127], etc.); however, roundoff and other errors generally
 184 accumulate proportional to the size of the network, which can turn out to be numerically catastrophic.
 185 Moreover, high dimensionality makes the computation and storage of $H_{\hat{f}}^T$ impractical, and thus practitioners
 186 typically resort to quasi-Newton methods that aim to avoid direct consideration of second derivatives.

187 Broyden’s method [21] for solving nonlinear systems, in the context of optimization, first approximates
 188 $(H_{\hat{f}}^T)^0 = I$, and then iteratively uses rank-one updates aiming for successively better estimates. Each
 189 iteration, one solves $(H_{\hat{f}}^T)^q \Delta c^q = -J_{\hat{f}}^T(c^q)$ to find a search direction Δc^q , and then uses line search to find
 190 c^{q+1} ; subsequently, Δc^q is updated via $\Delta c^q = c^{q+1} - c^q$. Given $(\Delta J_{\hat{f}}^T)^q = J_{\hat{f}}^T(c^{q+1}) - J_{\hat{f}}^T(c^q)$, the rank-one
 191 update is

$$(H_{\hat{f}}^T)^{q+1} = (H_{\hat{f}}^T)^q + \frac{1}{(\Delta c^q)^T \Delta c^q} \left((\Delta J_{\hat{f}}^T)^q - (H_{\hat{f}}^T)^q \Delta c^q \right) (\Delta c^q)^T \quad (4)$$

192 so that $(H_{\hat{f}}^T)^{q+1} \Delta c^q = (\Delta J_{\hat{f}}^T)^q$. When c is of large dimension, forming and inverting the dense $O(n^2)$
 193 $H_{\hat{f}}^T$ is undesirable, especially considering that the approximation is built from rank-one updates, and thus
 194 a matrix-free approach to the action of $H_{\hat{f}}^{-T}$ on a vector is preferred. That is, $\Delta c^q = -\left(H_{\hat{f}}^{-T}\right)^q J_{\hat{f}}^T(c^q)$ is
 195 used to find the search direction for the line search used to determine c^{q+1} , which is used to update Δc^q and

196 $(\Delta J_{\hat{f}}^T)^q$; then, rank-one update for $H_{\hat{f}}^{-T}$ is

$$(H_{\hat{f}}^{-T})^{q+1} = (H_{\hat{f}}^{-T})^q + \frac{(\Delta c^q - (H_{\hat{f}}^{-T})^q (\Delta J_{\hat{f}}^T)^q) (\Delta c^q)^T (H_{\hat{f}}^{-T})^q}{(\Delta c^q)^T (H_{\hat{f}}^{-T})^q (\Delta J_{\hat{f}}^T)^q}, \quad (5)$$

197 so that $(H_{\hat{f}}^{-T})^{q+1} (\Delta J_{\hat{f}}^T)^q = \Delta c^q$. Other low-rank update methods such as SR1 [33, 22], DFP [33, 46], and
 198 BFGS [23, 44, 52, 119] are similar in spirit. In particular, the limited-memory L-BFGS [106] only stores the
 199 past several low-rank updates making it quite efficient, see e.g. [78, 34].

200 Instead of performing rank-one updates to improve upon $(H_{\hat{f}}^T)^0 = I$ as in Broyden-style methods,
 201 gradient descent methods simply use $H_{\hat{f}}^T = I$ so that the search direction is obtained trivially via $\Delta c^q =$
 202 $-J_{\hat{f}}^T(c^q) = -\nabla \hat{f}(c^q)$. When problems have high dimensionality, practitioners often make further simplifi-
 203 cations such as evaluating only a subset of the right-hand side (mini-batch gradient descent) or even just
 204 one or a few randomly-selected entries at a time (SGD), see e.g. [117, 18]. One can even ignore the search
 205 direction equation entirely by choosing Δc^q to be various basis vectors, i.e. coordinate descent [120]. Fur-
 206 thermore, gradient descent methods can be envisioned as forward Euler approximations of gradient flow,
 207 i.e. of $\frac{dc(t)}{dt} = -\nabla \hat{f}(c(t))$, which allows for the wealth of knowledge in designing and solving ordinary differ-
 208 ential equations to be utilized. For instance, adaptive time stepping leads to such techniques as AdaGrad
 209 [37], which utilizes separate learning rates (time steps) for each parameter, or AdaDelta [139] and RMSprop
 210 [129], both of which lessen the effects of history terms in AdaGrad in order to maintain a sufficiently positive
 211 learning rate to avoid stalling. Incorporating the effects of prior search directions and state updates can
 212 be seen as utilizing momentum, which rewrites gradient flow using Newton’s Second Law [114]. The Adam
 213 method [72] combines the notion of using a moving average of gradients as in momentum methods with an
 214 adaptive learning rate for each parameter. The 52,000+ citations² of [72] indicate the success practitioners
 215 have enjoyed with Adam, often finding that it converges faster than SGD.

216 4. Iterative Levenberg-Marquardt

217 When training a neural network on data (x_i, y_i) , one seeks to find the parameters c of a generally vector-
 218 valued function $f(x, y, c)$ that minimize error over the training data, i.e. one desires $\|f(x_i, y_i, c)\|$ to be
 219 close to zero for all i . Choosing the L_2 norm leads to minimizing $\hat{f}(c) = \frac{1}{2} \sum_i f(x_i, y_i, c)^T f(x_i, y_i, c) =$
 220 $\frac{1}{2} \tilde{f}^T(c) \tilde{f}(c)$, which is a nonlinear least squares problem [14]. Critical points have $J_{\hat{f}}^T(c) = J_{\tilde{f}}^T(c) \tilde{f}(c) = 0$,
 221 which can be rewritten using the Taylor expansion of $\tilde{f}(c)$ about c^q as $J_{\tilde{f}}^T(c) \left(\tilde{f}(c^q) + J_{\tilde{f}}(c^q) \Delta c^q + \dots \right) = 0$,
 222 where $\Delta c^q = c - c^q$. Dropping high-order terms and evaluating $J_{\tilde{f}}^T$ at c^q leads to the Gauss-Newton equations
 223 $J_{\tilde{f}}^T(c^q) J_{\tilde{f}}(c^q) \Delta c^q \approx -J_{\tilde{f}}^T(c^q) \tilde{f}(c^q)$, which imply an estimate of $H_{\hat{f}}^T(c^q) \approx J_{\tilde{f}}^T(c^q) J_{\tilde{f}}(c^q)$, see e.g. [107].
 224 Notably, the Gauss-Newton approximation to the Hessian only requires first derivatives. Moreover, since the
 225 Gauss-Newton equations are the normal equations for $J_{\tilde{f}}(c^q) \Delta c^q = -\tilde{f}(c^q)$, one can obtain Δc^q via any
 226 least squares and minimum norm approach for solving this much better conditioned set of equations.

227 When $J_{\tilde{f}}$ is poorly-conditioned or rank-deficient, one can regularize the Gauss-Newton equations via
 228 $(J_{\tilde{f}}^T(c^q) J_{\tilde{f}}(c^q) + \epsilon^2 I) \Delta c^q = -J_{\tilde{f}}^T(c^q) \tilde{f}(c^q)$ with $\epsilon > 0$, which is referred to as Levenberg-Marquardt or
 229 damped nonlinear least squares, see e.g. [83, 97, 107, 14]. This makes a tradeoff between solvability and
 230 accuracy, since the unique and least squares components of the solution will be perturbed away from their
 231 correct values. To illuminate this, consider stacking a general linear system $Ac = b$ with the full-rank $\epsilon I c = 0$
 232 to obtain

$$\begin{pmatrix} A \\ \epsilon I \end{pmatrix} c = \begin{pmatrix} b \\ 0 \end{pmatrix}, \quad (6)$$

²as of September 2020, according to Google Scholar

233 which has equivalent normal equations of $(A^T A + \epsilon^2 I) c = A^T b$. Using the SVD, $A = U \Sigma V^T$, this becomes
 234 $(\Sigma^T \Sigma + \epsilon^2 I) \hat{c} = \Sigma^T \hat{b}$ where $\hat{c} = V^T c$ and $\hat{b} = U^T b$. For a general A , Σ has the form

$$\Sigma = \begin{pmatrix} \hat{\Sigma} & 0 \\ 0 & 0 \end{pmatrix}, \quad (7)$$

235 where $\hat{\Sigma}$ is diagonal and full-rank. This leads to

$$\left(\begin{pmatrix} \hat{\Sigma}^T & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \hat{\Sigma} & 0 \\ 0 & 0 \end{pmatrix} + \epsilon^2 I \right) \begin{pmatrix} \hat{c}_r \\ \hat{c}_z \end{pmatrix} = \begin{pmatrix} \hat{\Sigma}^T & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \hat{b}_r \\ \hat{b}_z \end{pmatrix}, \quad (8)$$

236 where \hat{c} and \hat{b} have been decomposed to separate out the portions that correspond to identically-zero sub-
 237 matrices of Σ^T . Equation 8 sets \hat{c}_z identically equal to zero as desired (i.e. minimum norm solution), but
 238 the entries in \hat{c}_r are determined via

$$\hat{c}_k = \frac{\sigma_k}{\sigma_k^2 + \epsilon^2} \hat{b}_k = \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k}, \quad (9)$$

239 perturbing them away from their correct unique or least squares solution $\hat{c}_k = \hat{b}_k / \sigma_k$. This perturbation is
 240 negligible for $\sigma_k \gg \epsilon$, but smaller σ_k have their associated \hat{c}_k more significantly incorrectly perturbed toward
 241 zero. One typically chooses ϵ so that it does not interfere too much with the larger (more important) singular
 242 values, while still being large enough to regularize numerical issues associated with smaller σ_k (as well as
 243 identically zero singular values). As a side note for weighted least squares, one adds the full-rank $\epsilon D c = 0$
 244 (with diagonal D) instead of $\epsilon I c = 0$ to obtain a modified version of Equation 6, which after column scaling
 245 becomes

$$\begin{pmatrix} AD^{-1} \\ \epsilon I \end{pmatrix} D c = \begin{pmatrix} b \\ 0 \end{pmatrix}. \quad (10)$$

246 Then, a simple renaming of variables results in the original Equation 6, and the above analysis applies
 247 without modification.

248 Motivated by the Broyden-style iterative methods (discussed in the previous section) which began with
 249 a simple guess for the Hessian and then corrected it after each iteration, we propose a similar strategy for
 250 Levenberg-Marquardt. That is, we start with $\epsilon I c = 0$ but subsequently update the right-hand side as the
 251 iteration proceeds, progressively removing the erroneous perturbation of the least squares solution shown in
 252 Equation 9. Our approach converges to the exact solution for larger singular values, as for example would
 253 also be achieved using PCA; however, unlike the all-or-nothing approach of PCA, our approach *smoothly*
 254 tapers between the exact solution for larger σ_k and robust regularization for smaller σ_k .³

255 We start with a guess c^* for c and stack $A c = b$ with $\epsilon I c = \epsilon c^*$ leading to the normal equations

$$(A^T A + \epsilon^2 I) c = A^T b + \epsilon^2 c^*. \quad (11)$$

256 Substituting the SVD of A leads to

$$(\Sigma^T \Sigma + \epsilon^2 I) \hat{c} = \Sigma^T \hat{b} + \epsilon^2 V^T c^* = \Sigma^T \hat{b} + \epsilon^2 \hat{c}^*. \quad (12)$$

257 where $\hat{c}^* = V^T c^*$. This modified version of Equation 8 sets \hat{c}_z equal to \hat{c}_z^* , while the entries in \hat{c}_r are
 258 determined via

$$\hat{c}_k = \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k} + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) \hat{c}_k^* \quad (13)$$

259 illustrating that \hat{c}_k is a convex combination of the exact solution \hat{b}_k / σ_k and the initial guess \hat{c}_k^* . When
 260 $\sigma_k \gg \epsilon$, the associated \hat{c}_k tend toward the correct solution as usual. When $\sigma_k \ll \epsilon$, the associated \hat{c}_k tend
 261 toward \hat{c}_k^* .

³This method/proof was derived for a CS205L lecture at Stanford in Winter quarter 2019 [41].

262 Starting with a guess of $c^* = 0$, one obtains $\hat{c}^* = 0$ and thus $\hat{c}_z^* = 0$ and $\hat{c}_z = 0$ as desired. In addition,
 263 $\hat{c}_r^* = 0$ and Equation 13 is identical to Equation 9. Multiplying by V transforms \hat{c} back to the c that would
 264 result from solving Equation 11. Setting c^* equal to this newly obtained value of c and repeating the above
 265 analysis maintains $\hat{c}_z = 0$ (as desired), while

$$\hat{c}_k^* = \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k} \quad (14)$$

266 so that Equation 13 becomes

$$\hat{c}_k = \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k} + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k} = \left(1 + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) \right) \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k}. \quad (15)$$

Repeating the entire process again results in

$$\begin{aligned} \hat{c}_k &= \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k} + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) \left(1 + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) \right) \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k} \\ &= \left(1 + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right)^2 \right) \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k}, \end{aligned} \quad (16)$$

267 and further iterations give

$$\hat{c}_k = \left(1 + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right) + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right)^2 + \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right)^3 + \dots \right) \left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2} \right) \frac{\hat{b}_k}{\sigma_k}, \quad (17)$$

268 where the term in parentheses is a geometric series with $r = \frac{\epsilon^2}{\sigma_k^2 + \epsilon^2}$.

269 Since the geometric series in Equation 17 converges to $\frac{1}{1-r} = \frac{\sigma_k^2 + \epsilon^2}{\sigma_k^2}$, Equation 17 converges to the exact
 270 solution $\hat{c}_k = \hat{b}_k/\sigma_k$. Any practical numerical method will only take q steps, leading to the partial sum

$$\frac{1 - r^q}{1 - r} = \frac{\sigma_k^2 + \epsilon^2}{\sigma_k^2} \left(1 - \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right)^q \right), \quad (18)$$

271 which yields

$$\hat{c}_k = \left(1 - \left(\frac{\epsilon^2}{\sigma_k^2 + \epsilon^2} \right)^q \right) \frac{\hat{b}_k}{\sigma_k}. \quad (19)$$

272 The scalar term premultiplying \hat{b}_k/σ_k monotonically approaches 1 as the iteration proceeds, and thus each
 273 \hat{c}_k converges monotonically to the exact solution and converges more quickly for larger σ_k as desired.

274 4.1. Examples

275 Typical inverse, control, and learning problems involve numerically challenging data, where linear sub-
 276 problems may have coefficient matrices with both small and identically zero singular values and the right-
 277 hand side may not be in the range of the coefficient matrix. Accordingly, we evaluate our approach against
 278 these types of problems. Here, we compare our iterative Levenberg-Marquardt (iLM) approach to PCA
 279 because PCA is a widely applied and well-understood algorithm. Since our goal is merely to demonstrate
 280 the feasibility of iLM, we utilize straightforward Matlab implementations of both methods. For iLM, we
 281 solve Equation 11 using Matlab's `pcg` routine with no preconditioner, i.e. conjugate gradients. For PCA,
 282 we compute the largest singular values and corresponding singular vectors using Matlab's `svds` function,
 283 which finds these quantities via either Lanczos bidiagonalization [5, 77] or a computation of the full SVD
 284 depending on the number of singular vectors desired. All experiments were run on a workstation equipped
 285 with Matlab R2020a, 128GB RAM, and a 24-core Intel CPU running at 3.00GHz.

286 First we consider rather large dense matrices and compare PCA and iLM for a varying number of singular
 287 values, noting that an increased number of iLM iterations is required for increased accuracy. We generate

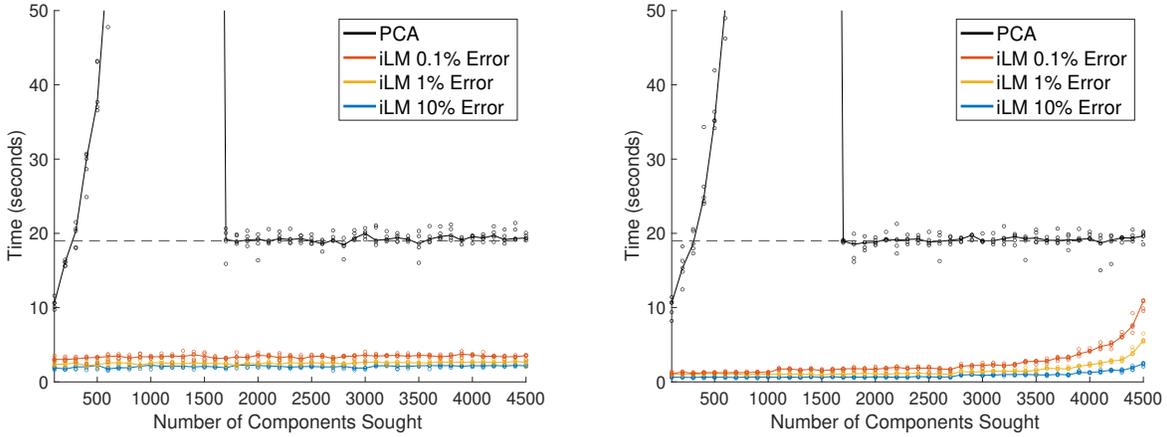


Figure 5: Performance of iLM and PCA for estimating an increasing number of \hat{c}_k given a dense $5,000 \times 5,000$ square matrix with 25 million randomly-generated entries (post-processed to have 100 zero singular values) and a randomly-generated right-hand side not in its range. We observed deleterious behavior of the Matlab software for an intermediate range of sought components⁴; however, a better-devised approach would obviously not rise above the hashed line in the figures. Five experiments were run for each number of components tested, and solid lines are drawn through the median results. (Left) $\epsilon = 0.1$. (Right) $\epsilon = 5.0$. Increased regularization slows convergence for singular values that are very small compared to ϵ , as expected (see Figure 6).

288 five random dense $5,000 \times 5,000$ matrices (each post-processed to have 100 zero singular values) as well as
 289 random right-hand sides outside the range of the coefficient matrix. Then, we estimate various numbers of
 290 components of \hat{c} , noting that in typical applications one seeks only a small number of components. While
 291 PCA estimates these components “exactly” up to numerical precision, the accuracy of iLM is limited by
 292 regularization (i.e. ϵ) and the tolerance of the CG solver. The number of iLM iterations is chosen so that
 293 the relevant \hat{c}_k are within 10%, 1%, or 0.1% of the \hat{c}_k obtained via PCA in the L_∞ norm; this required CG
 294 solver tolerances of $1e-7$, $1e-8$, and $1e-10$, respectively. Results are shown in Figure 5. Data from each of
 295 the five trials are plotted as circles, and the median results across the five trials are connected by solid lines.
 296 The hashed line represents the approximate time taken to compute the SVD of the coefficient matrix. The
 297 number of iLM iterations required for the median results for each level of accuracy are plotted in Figure 6.
 298 To help clarify the required CG tolerance for iLM, we plot in Figure 7 the tolerance required to obtain each
 299 level of accuracy for one of the five trials.

300 We also consider how iLM and PCA perform for a fixed number of desired \hat{c}_k as the size of the coefficient
 301 matrix increases. We let size vary from $1,000 \times 1,000$ to $10,000 \times 10,000$ and seek 500 components, creating a
 302 random dense matrix (post-processed to have 100 zero singular values) and random right-hand side not in
 303 the range of the coefficient matrix. For iLM, we iterate (as before) until the L_∞ norm of the vector of \hat{c}_k
 304 is within 10%, 1%, or 0.1% of that obtained via PCA, using a CG tolerance of $1e-7$, $1e-8$, or $1e-10$, respectively.
 305 The results are shown in Figure 8.

306 The aforementioned tests are unfair to iLM because they stringently require iLM to do as well as PCA
 307 on the values PCA estimates nearly exactly while ignoring the fact that PCA obtains totally inaccurate
 308 (identically zero) solutions for all the other \hat{c}_k . To illustrate the added benefit of smooth convergence
 309 obtained via iLM, we construct a small (to make the graphs easier to read) 100×100 random matrix with
 310 ten of its singular values set to zero (see Figure 9). A random right-hand side b outside the range of A is
 311 then formed. Figures 10 and 11 show the results of iLM and PCA, illustrating how well the obtained $\sigma_k \hat{c}_k$
 312 reconstruct the projected right-hand side \hat{b}_k . iLM leverages rich information about the structure of A even
 313 when ϵ is larger than the largest singular value of A (substantial regularization). In these experiments we
 314 used a CG tolerance of $1e-6$ and a maximum of 1000 CG iterations.

⁴We observed that, by default, Matlab seems to wait too long to switch from Lanczos bidiagonalization to computing the full SVD.

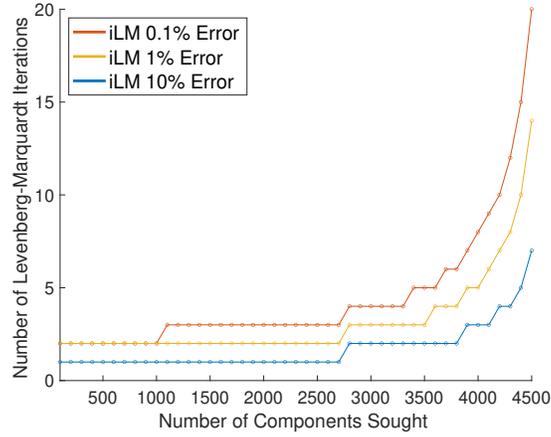


Figure 6: Number of iterations required for iLM for the medians of the trials in Figure 5 (Right) with $\epsilon = 5.0$. (When $\epsilon = 0.1$, iLM mostly converges to the desired tolerance in one iteration.)

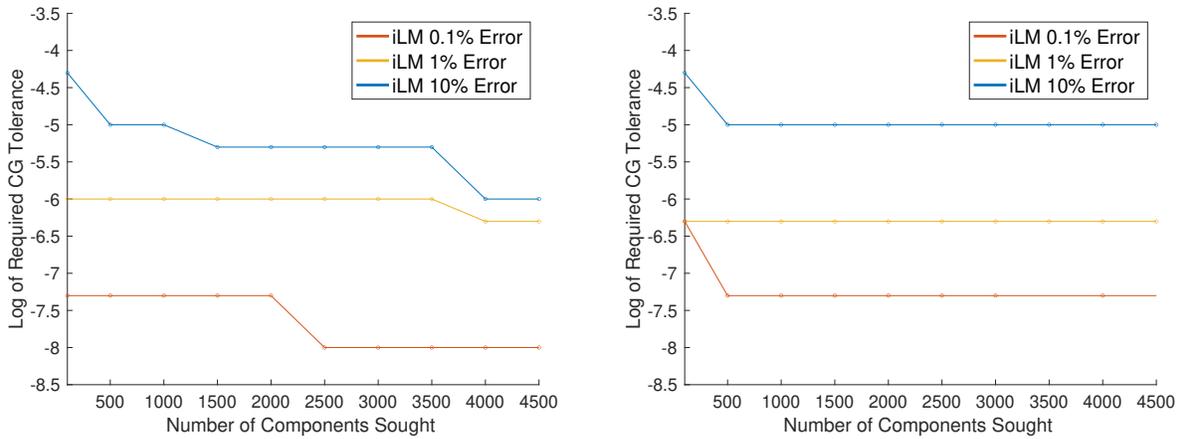


Figure 7: Using iLM, as more \hat{c}_k are sought or as more accuracy is desired, a tighter CG tolerance needs to be used to prevent convergence from stalling. Plotted are the experimentally-determined maximum CG tolerances which yielded convergent results, ranging from 100 to 4500 \hat{c}_k sought. (Left) $\epsilon = 0.1$. (Right) $\epsilon = 5.0$.

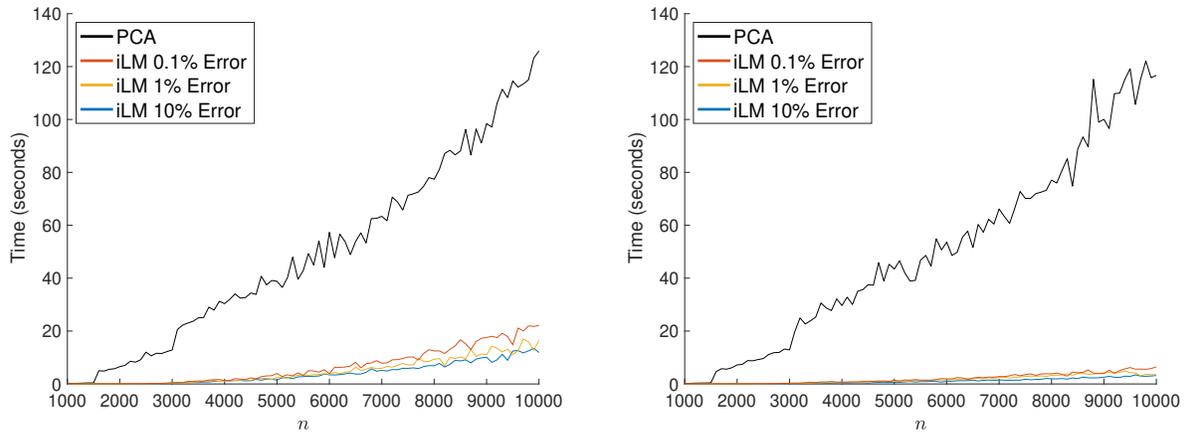


Figure 8: Performance of iLM and PCA for estimating a fixed number of \hat{c}_k (500 of them) given a dense $n \times n$ square matrix with n^2 randomly-generated entries (post-processed to have 100 zero singular values) and a randomly-generated right-hand side not in its range. (Left) $\epsilon = 0.1$. Only one iLM iteration is required for these trials, and hence the growing cost is a combination of the increased cost per CG iteration and the number of CG iterations required. (Right) $\epsilon = 5.0$. For large n , the added regularization appears to aid in the convergence of CG.

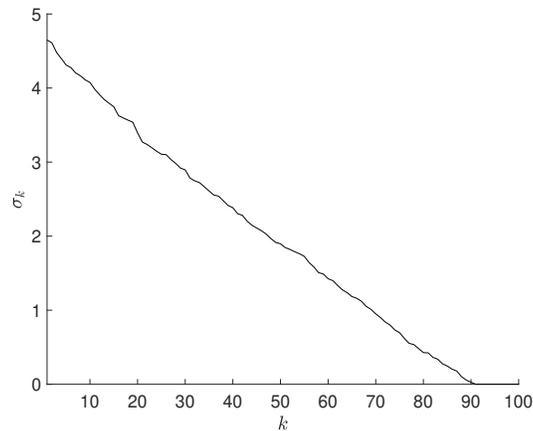


Figure 9: The singular values of the matrix used for the experiments shown in Figures 10 and 11. Ten of the singular values are identically zero.

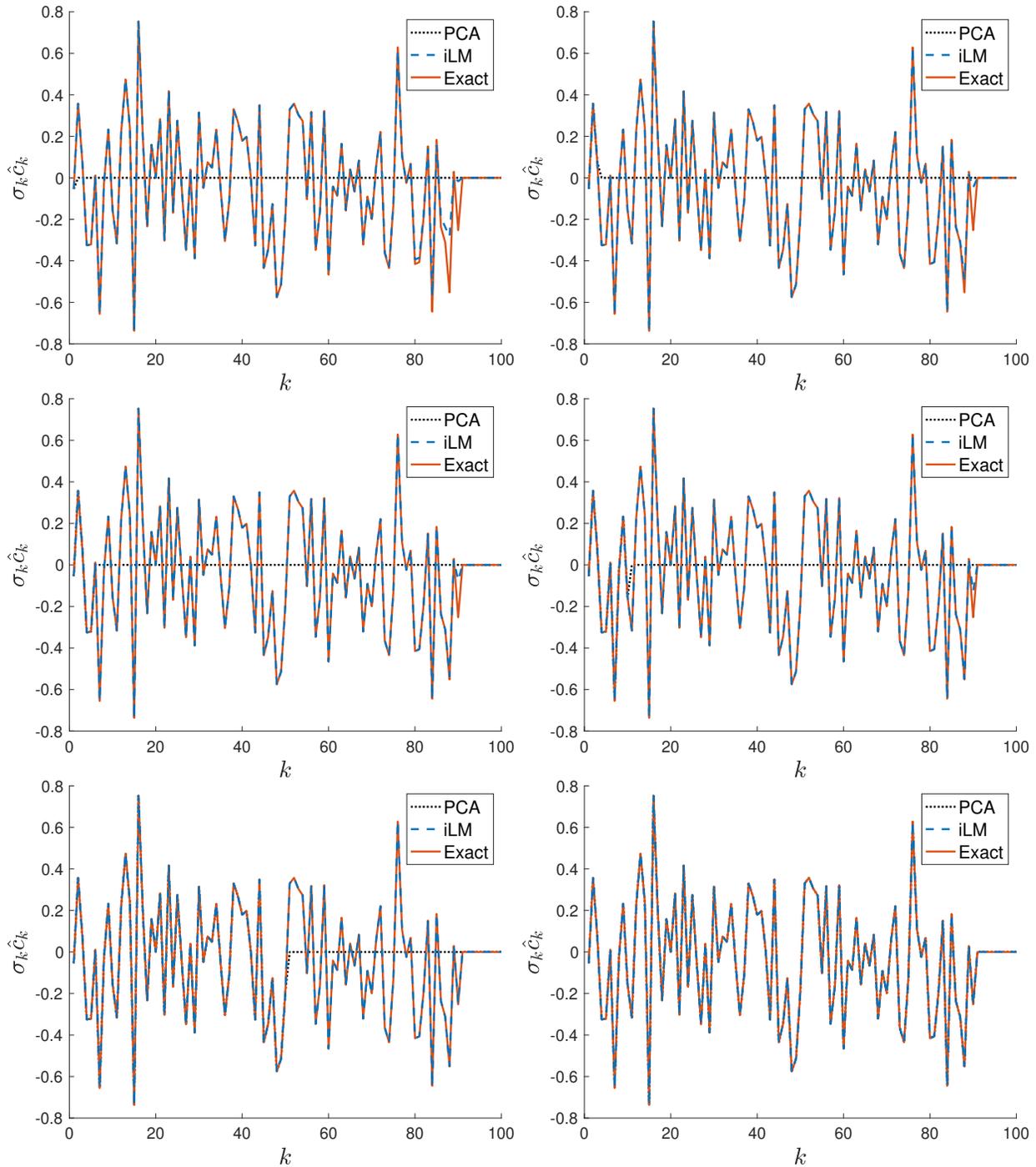


Figure 10: Convergence of iLM and PCA to the exact solutions using 1, 3, 5, 10, 50, and 100 iterations with $\epsilon = 0.1$ for iLM and 1, 3, 5, 10, 50, and 100 components for PCA, respectively. iLM converges quickly for \hat{c}_k associated with larger singular values but takes additional iterations to converge for the smallest singular values due to the regularization. iLM and PCA are both exact to numerical precision for the \hat{c}_k that should be identically zero (i.e. those in \hat{c}_z).

315 *4.1.1. Comparisons for Nonlinear Optimization Problems*

316 The above examples demonstrate the utility of iLM when solving linear problems such as those that
 317 arise on each iteration of a standard nonlinear optimization algorithm. We now explicitly consider solving
 318 nonlinear optimization problems using iLM and related approaches.

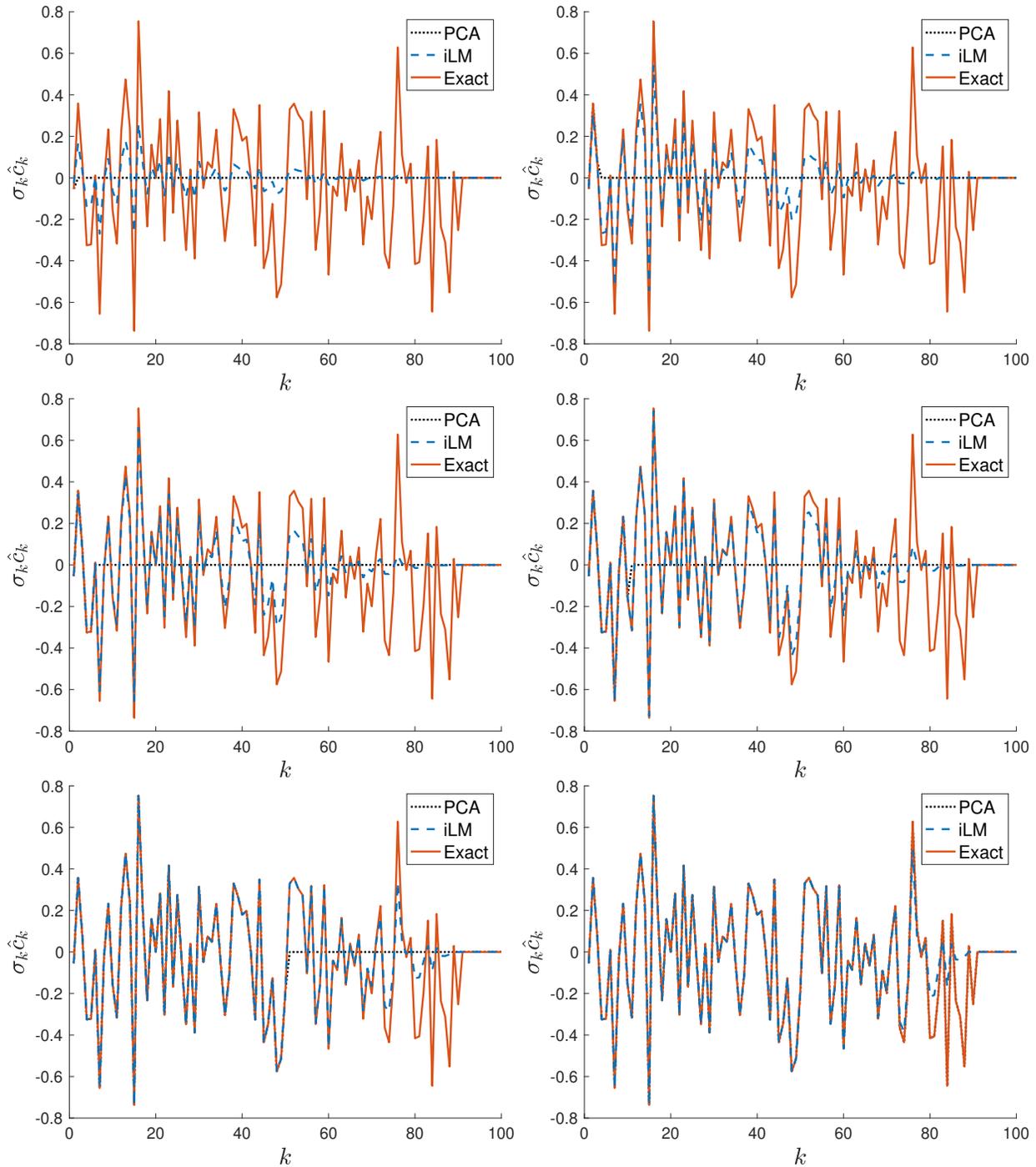


Figure 11: Same as Figure 10 except with $\epsilon = 5.0$. The increased regularization slows down iLM convergence, as expected.

319 We first consider an objective $f(x_1, x_2) = x_1^2 + 5x_2^2 - 4$ with an initial guess of $x^0 = (-3 \ -4)^T$. Figure
320 12 shows the results of using Newton's method, gradient descent, Levenberg-Marquardt, Fan's modified
321 Levenberg-Marquardt [40], and iLM to solve this problem. Each method is allowed to run until either the
322 objective at the current iterate is within 10^{-6} of the analytic minimum value or until the L_2 norm of the
323 iterate changes less than 10^{-6} between iterations. All CG solves use a tolerance of 10^{-6} and a maximum
324 of 1,000 iterations. Since the objective is quadratic, Newton's method converges to the unique, global
325 minimum in one iteration. Gradient descent, which lacks second-derivative information, oscillates around

326 the (non-uniformly-scaled) energy landscape before eventually reaching the minimum. A learning rate of
 327 0.15 was used. Levenberg-Marquardt can be seen as blending between the Newton and gradient descent
 328 iterates. With little regularization ($\epsilon = 0.1$), Levenberg-Marquardt looks similar to Newton’s method,
 329 while with a large regularization parameter ($\epsilon = 100.0$), the number of iterations required for convergence
 330 significantly increases. The modified Levenberg-Marquardt of [40] can offer cubic convergence rates under
 331 suitable conditions by performing essentially two Levenberg-Marquardt steps on each iteration (a standard
 332 step and a forward-looking step based on the standard step). We implemented this method using the same
 333 parameters as in Section 4 of [40], except we used an initial μ_0 of 10 in order to be similar to our regularization
 334 of the other Levenberg-Marquardt variants. Finally, we consider iLM using 1, 10, or 100 iterations, all with
 335 $\epsilon = 10.0$ and using an initial guess of $c^* = 0$. Note that iLM uses an ϵ^2 rather than an ϵ scaling of the
 336 identity term, so $\epsilon = 10.0$ is equivalent to $\epsilon = 100.0$ with Levenberg-Marquardt. iLM converges to the
 337 Newton step (when the Newton step is defined) as the number of iterations increases, so iLM has a quadratic
 338 order of convergence in the best case; though of course, like Levenberg-Marquardt, gradient descent, etc., it
 339 is possible to design parameters and scenarios which make iLM converge poorly or not at all. Moreover, we
 340 stress that various strategies for adaptive learning rates and adaptive regularization terms may improve the
 341 performance of these methods. In particular, the adaptive parameter values used for our implementation
 342 of [40] are quite useful for aiding the convergence of the method, and in practice one would want to utilize
 343 adaptive regularization schemes for Levenberg-Marquardt and iLM as well (which would likely remove many
 344 of the small steps those algorithms take as they approach the solution).

345 As a potentially greater challenge, we consider adding a third coordinate to our objective. We alter
 346 our initial guess to have a value of 1 along this direction. Since the objective function does not depend
 347 on this third component, it is possible for the solution iterate to drift along this additional axis, e.g. when
 348 regularization is perturbing the solution away from the true minimum. Newton’s method is not applicable in
 349 this case since the Hessian becomes singular, although iLM appears to converge to what Newton’s method
 350 would compute using the Hessian’s pseudoinverse. Interestingly, all methods appear to converge to the
 351 solution $(0 \ 0 \ 1)^T$, rather than e.g. the minimum norm solution at the origin. We also consider rotating
 352 the objective and initial guess by 30 degrees about the x_1 and x_3 axes in order to make the null space of
 353 the Hessian less obvious. However, the optimization methods we tested still reach the minimum in the same
 354 number of iterations as reported in Figure 12, except for Newton’s method, which remains undefined.

355 Further differences in the behavior and convergence of these optimization algorithms can be elucidated by
 356 considering the slightly modified objective $f(x) = \min(x_1^2 + 5x_2^2 - 4, (x_1 + .1)^2 + 5(x_2 - .1)^2 - 4)$, which has
 357 minima at $(0 \ 0)^T$ and $(-.1 \ .1)^T$. With the same initial guess of $x^0 = (-3 \ -4)^T$, the nearest minimum
 358 in the L_2 norm is the minimum-norm solution $(0 \ 0)^T$. However, if an algorithm does not proceed directly
 359 towards this solution, it may instead converge towards the other minimum with greater norm and less
 360 sparsity. This is demonstrated in Figure 13. Newton’s method converges in one step to $(0 \ 0)^T$. With
 361 enough iterations, iLM approximates the Newton step and also selects the minimum-norm solution. With
 362 fewer iterations, though, iLM behaves more like gradient descent and Fan’s modified Levenberg-Marquardt,
 363 which select the non-zero minimum. Standard Levenberg-Marquardt can be driven to select different minima
 364 by tuning the regularization parameter. In general, one must consider the types of solutions one seeks to
 365 an optimization problem (e.g., minimum-norm or sparse solutions) when selecting an algorithm and its
 366 parameters. Practical considerations like this and real-world performance tradeoffs can often overshadow
 367 theoretical convergence guarantees, as seen for example with the continued ubiquity of (stochastic) gradient
 368 descent.

369 5. Column Space Search

370 For the sake of motivation, consider the 2×2 linear subproblem $Ac = b$ where $A = \begin{bmatrix} 1 & -1 \\ .1 & 1 \times 10^{-6} \end{bmatrix}$ and
 371 $b = [0 \ 1]^T$. Although the right-hand side is in the range of A , it is not “easily” in the range of A ; in
 372 other words, the columns of A are mostly orthogonal to b leading to a solution that utilizes large multipliers
 373 $c_1 = 1/(1 \times 10^{-6} + .1)$ and $c_2 = 1/(1 \times 10^{-6} + .1)$ on the columns of A . See Figure 14. Even though this
 374 is the exact solution to the linear subproblem, it misleadingly heavily weights columns of A that do not
 375 correlate well with the desired b . Large values of c_1 and c_2 seemingly indicate that those columns of A are

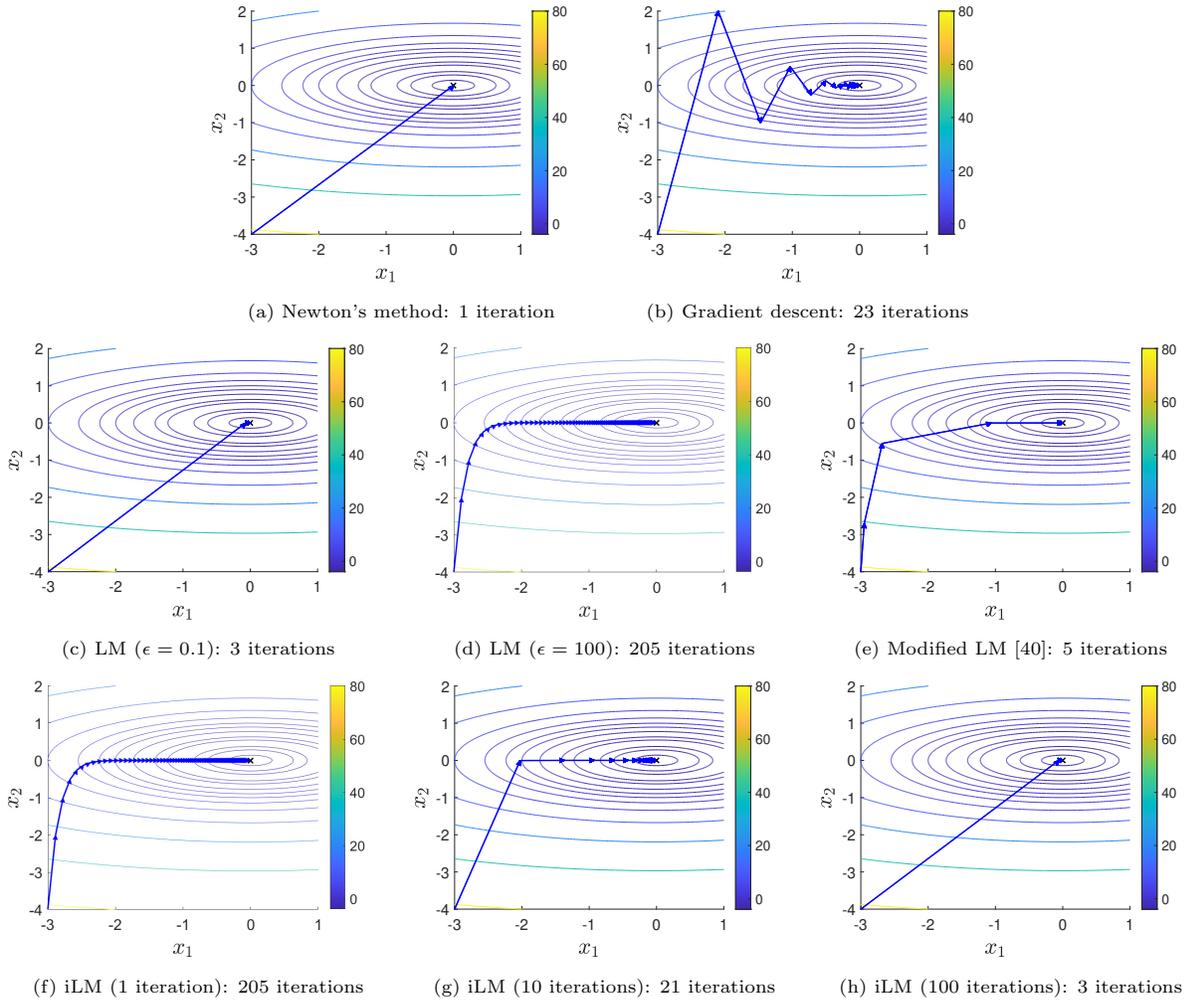


Figure 12: Different algorithms applied to minimizing $f(x_1, x_2) = x_1^2 + 5x_2^2 - 4$ with an initial guess of $x^0 = (-3 \ -4)^T$. Contours of the function are drawn and shaded by contour value. Arrows indicate steps taken on each iteration of the optimization as the algorithm is allowed to converge to $(0 \ 0)^T$ (the black x). The number of iterations required for each method to converge to a tolerance of 10^{-6} is reported. We emphasize that these methods have different computational requirements; for instance, a Levenberg-Marquardt (LM) step requires one linear solve, while an iteration of [40] requires two.

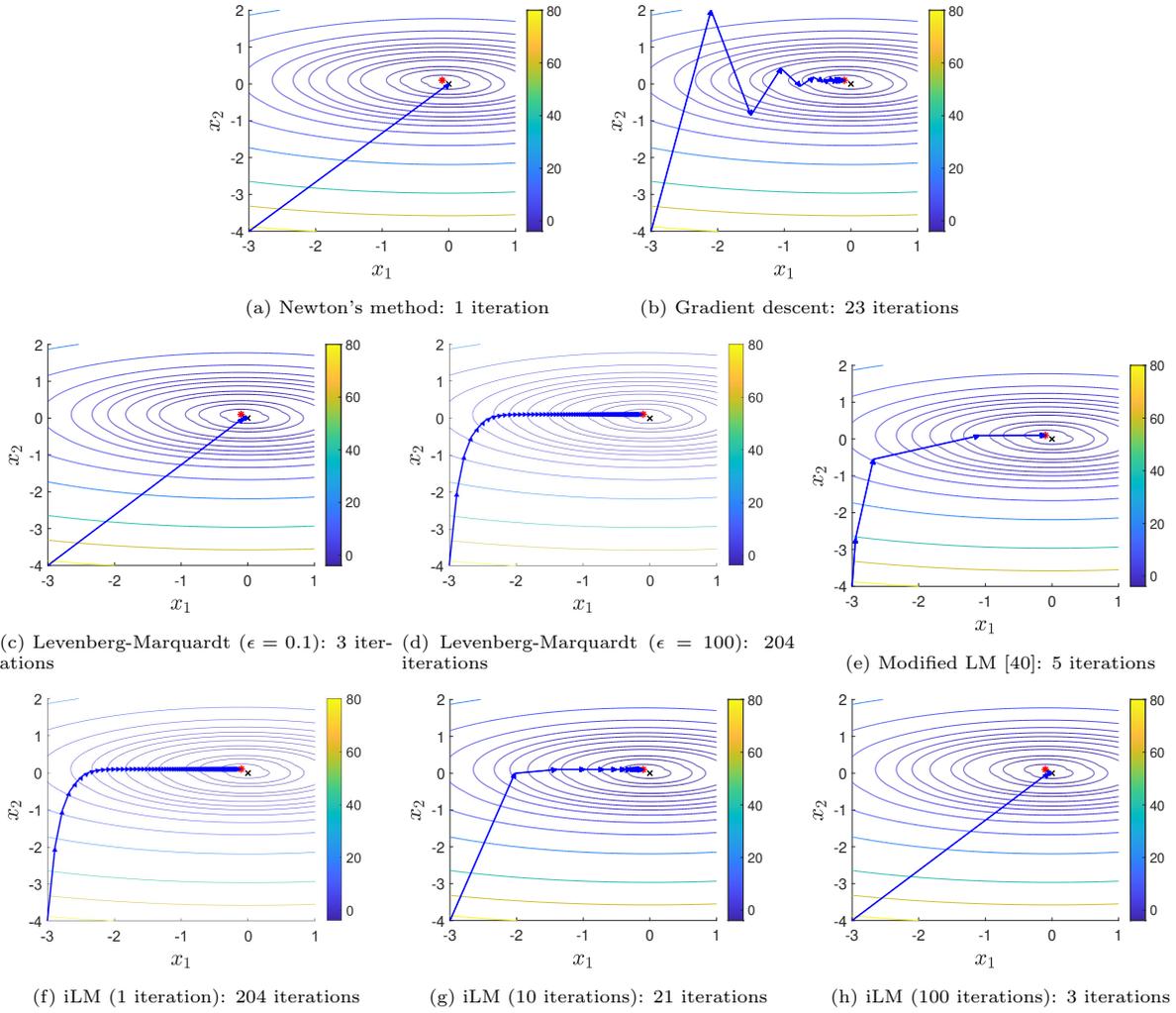


Figure 13: Repeating the experiment of Figure 12 with an objective of $f(x_1, x_2) = \min(x_1^2 + 5x_2^2 - 4, (x_1 + .1)^2 + 5(x_2 - .1)^2 - 4)$. Algorithms converge to either $(0 \ 0)^T$ (the black x) or to $(-.1 \ .1)^T$ (the red star).

376 important, even though they mostly cancel each other out being nearly orthogonal to the right-hand side.
 377 The regularized least squares problem $\min_c \|b - Ac\|_2^2 + \lambda \|c\|_2^2$ reduces the values of c_1 and c_2 , although it
 378 does not alleviate the fact that these columns mostly work to cancel each other out, making minimal progress
 379 towards b . At best, heavy regularization could drive c_1 and c_2 even further towards zero.

380 As previously discussed, the columns of the linear subproblem are often quite erroneous approximations
 381 to the Hessian, which itself is only a linearization of the nonlinear problem; yet the linear subproblem is
 382 often solved and used to increment the solution vector (i.e. via $c^{q+1} = c^q + \Delta c^q$). The original nonlinear
 383 problem may include significant noise and heavy regularization, and thus it seems more important to focus
 384 on controls that make direct progress towards energy/loss minimization than those that make only incidental
 385 progress while competing with and largely cancelling each other. Thus, we advocate dropping parameters
 386 from consideration when the gains made toward the solution by some combination of those parameters are
 387 incidental compared to the parameters' main actions. As discussed previously, in regard to neural networks,
 388 this allows one to identify and differentiate which building blocks of the neural network are more or less
 389 important than others. In order to identify the more important parameters, we make note of two common
 390 misconceptions/flaws in the pursuit of solving linear subproblems. First, solving the linear subproblem
 391 exactly is not necessarily desirable since the columns of A may be terrible approximations to those of the
 392 Hessian, which itself is a linearization. Second, the largest singular values of A do not necessarily represent
 393 the most important features of the problem (as is assumed by typical PCA approaches); oftentimes, the
 394 more important notion is which columns of A are well-correlated with the right-hand side b , allowing one to
 395 make clean, non-competitive progress toward the solution.

396 Next, consider the right-hand side $b = [5 \ 1]^T$, which is better correlated with at least one of the
 397 columns of A . See Figure 15. In this case, the exact solution in Figure 15b is an improvement over Figure
 398 14 (Right), but still contains problematic cancellation. The regularized solution shown in Figure 15c makes
 399 more progress towards the solution as compared to Figure 14 (Right), except it uses a lot more of a_2 and a lot
 400 less of a_1 than one might expect given how much better correlated a_1 is with b . Regularization damps the use
 401 of a_1 hindering its progress towards the solution; as such, a_2 ends up being utilized significantly. One could
 402 obtain a better solution for this example by changing the regularization in the least squares problem to have
 403 the form $\min_c \|b - Ac\|_2^2 + \lambda_1 c_1^2 + \lambda_2 c_2^2$ with $\lambda_1 = 0$. Figure 15d shows the result for $\lambda_2 = 1$ which is highly
 404 improved. One could do even better using only a_1 as shown in Figure 15e, obtained using $\min_{c_1} \|b - a_1 c_1\|_2^2$.
 405 For more discussion on various regularization strategies, especially pertaining to the facial expression inverse
 406 problem described in Section 2, see [15, 25, 86, 128, 13, 64, 136, 63, 87, 20, 68, 105].

407 The aforementioned discussion motivates the notion of choosing only the columns of A which are most
 408 correlated with b . Such an approach can be implemented one column at a time using a basic coordinate
 409 descent algorithm [111]. Importantly, this allows one to circumvent null spaces without adding regularization,
 410 making coordinate descent an attractive option for use on ill-posed, poorly-conditioned problems. At each
 411 iteration, the column can be chosen stochastically [104] or deterministically. Popular deterministic methods
 412 for choosing the next search direction include cyclic coordinate descent [74], the Gauss-Southwell (GS) and
 413 Gauss-Southwell-Lipschitz rule [108], and the maximum block improvement (MBI) rule [28]. Instead of
 414 looking at a single column at a time, block coordinate descent can be used to update multiple columns
 415 simultaneously [130]; however, regularization may still be needed when the block of columns is poorly-
 416 conditioned or does not have full rank. See [120, 135] for more discussion. Typical coordinate descent
 417 algorithms may choose a large number of poorly correlated coordinates in place of a smaller number of more
 418 strongly correlated coordinates. Using correlation to choose the next coordinate to add to the model can
 419 alleviate this problem and is the central idea behind MBI [28], forward and backward stepwise regression [35],
 420 and LARS [38]. The latter statistical regression methods are often used to gain better prediction accuracy
 421 and interpretability of the model [59]. However, LARS converges to the least squares solution of the linear
 422 subproblem [38] because it eventually uses uncorrelated coordinates.

423 In order to facilitate our goal of obtaining sparse, semantic solutions to optimization problems, particu-
 424 larly without adding unnecessary heuristic regularization which can lead to overfitting and error, we propose⁵
 425 solving linear subproblems by first pruning away any coordinates that are geometrically uncorrelated with
 426 the right-hand side as motivated by least angle regression (LARS) [38]; then, we estimate the remaining

⁵This approach was first proposed in the following preprint: [8].

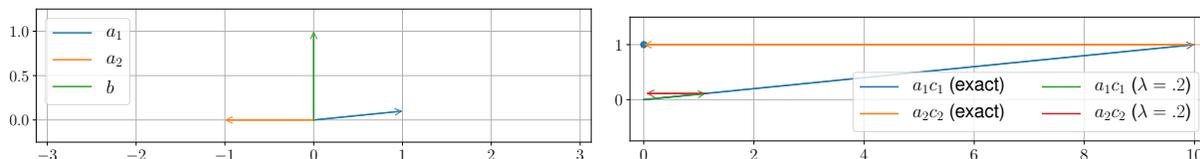


Figure 14: (Left) A visualization of the columns of A as well as b for the linear subproblem $Ac = b$ from Section 5 when $b = [0 \ 1]^T$. Note how the columns of A are mostly orthogonal to b . (Right) The exact solution utilizes quite large values of c_1 and c_2 , over-scaling largely competing columns of A in order to make progress towards b . Since the columns of A are often poor approximations to the Hessian, and the Hessian itself is only a linearization of the nonlinear problem, it seems imprudent to over-utilize controls c_1 and c_2 in order to make progress towards b . A regularized solution (with $\lambda = .2$) is also shown in the figure. It does reduce the magnitudes of c_1 and c_2 but still demonstrates the same non-desirable competitive behavior between the columns.

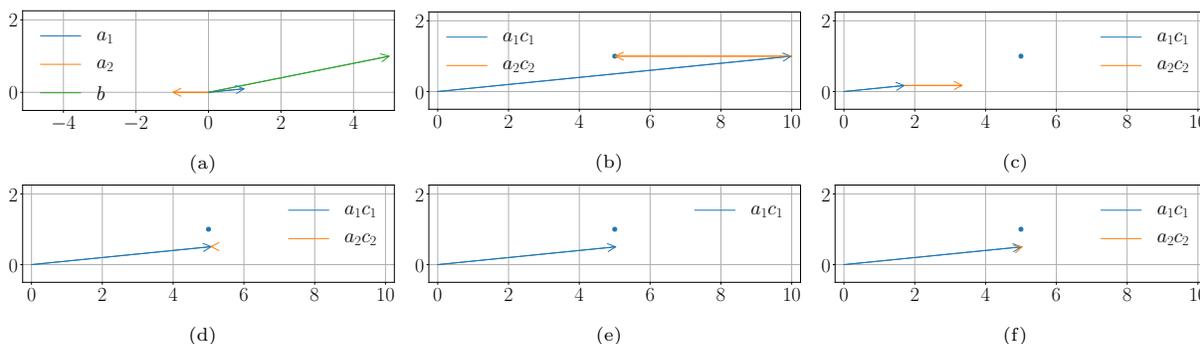


Figure 15: (a) A visualization of A 's columns and b for the linear problem $Ac = b$ from Section 5 when $b = [5 \ 1]^T$. (b) The exact solution depicted by a_1c_1 and a_2c_2 . (c) The regularized solution with $\lambda = 1$. (d) The regularized solution with $\lambda_1 = 0$ and $\lambda_2 = 1$. (e) The solution when solving for c_1 only. (f) The solution obtained after a few iterations of coordinate descent using the MBI selection rule.

427 coordinates via coordinate descent, eliminating the need to regularize for solvability.

428 5.1. Pruning Geometrically Uncorrelated Directions

429 We illustrate our approach, hereafter referred to as Column Space Search (CSS), by again consider-
 430 ing solving a generic nonlinear least squares optimization problem of the form $\min_c \|f(x, y, c)\|_2^2$. Using
 431 a Gauss-Newton based method, every iteration of the optimization requires solving the linear subproblem
 432 $J_{\tilde{f}}^T(c^q) J_{\tilde{f}}(c^q) \Delta c^q = -J_{\tilde{f}}^T(c^q) \tilde{f}(c^q)$ to find the Δc^q subsequently used to make progress towards the solution.
 433 Again, one may equivalently consider $J_{\tilde{f}}(c^q) \Delta c^q = -\tilde{f}(c^q)$.

434 We first compute the geometric correlation between each column j_i of $J_{\tilde{f}}(c^q)$ and the right-hand side
 435 $-\tilde{f}(c^q)$. Similar to LARS [38] and MBI [28], we use $|\hat{j}_i \cdot \tilde{f}(c^q)|$, where $\hat{j}_i = j_i / \|j_i\|_2$. Poorly geometrically
 436 correlated columns can only make significant progress towards the solution either when partially cancelled
 437 by other poorly geometrically correlated columns (as in Figure 14 (Right)) or as corrections to better geo-
 438 metrically correlated columns (as in Figure 15b). However, this so-called progress, while valid for the linear
 439 subproblem, may pollute the sparsity and semantics of the solution to the original nonlinear problem. See
 440 Figure 16. Thus, we prune poorly geometrically correlated columns from $J_{\tilde{f}}(c^q)$ resulting in a lower-rank
 441 J_S . Motivated by the Gauss-Southwell rule, one might instead prune using gain correlation $|j_i \cdot \tilde{f}(c^q)|$,
 442 which considers large residual decreases with smaller variable values; however, we instead prefer removing
 443 poorly geometrically correlated columns even when they may have large gains as it seems to lead to better
 444 semantic interpretation. Additionally, one could drop the absolute value and consider $\hat{j}_i \cdot \tilde{f}(c^q)$ in order to
 445 prune columns that are only semantically sensible in one direction.

446 Pruning columns of $J_{\tilde{f}}(c^q)$ to get a reduced J_S has the additional benefit of potentially eliminating
 447 portions of the null space of $J_{\tilde{f}}(c^q)$, as the pruned out columns or a combination of them with the non-
 448 pruned columns may have linear dependencies; this pruning may also improve the condition number. This
 449 is especially prudent when working with a large number of dimensions, in which case the dimension of the

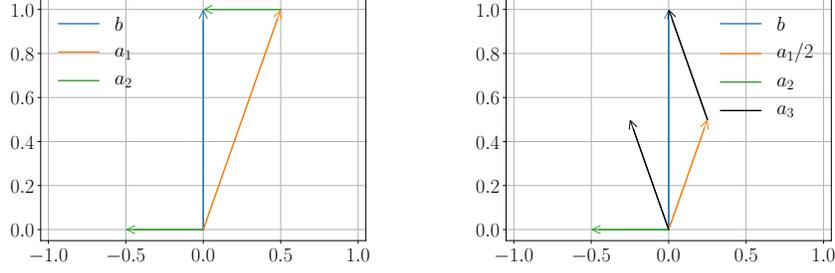


Figure 16: (Left) Here, $b = a_1 + a_2$, but a_2 is only valid for making progress towards b in conjunction with a_1 (and is otherwise orthogonal). While using a_2 may be desirable when trying to solve a truly linear system of equations, it makes less sense when solving a linearization of a high-dimensional nonlinear problem. (Right) It may be desirable to only progress in the direction of a_1 until other directions become better correlated. At that point, it would be better to find a new direction (in this case, a_3) that leads back towards b .

450 null space of $J_{\tilde{f}}(c^q)$ and the condition number of $J_{\tilde{f}}(c^q)$ may be quite large. Moreover, these difficulties are
 451 exacerbated when regularization is not used.

452 5.2. Solving the Pruned System

453 We avoid regularization entirely in order to avoid changing the solution to the problem; thus, we pursue
 454 a coordinate descent strategy to solve $J_S \Delta c_S^q = -\tilde{f}(c^q)$ where Δc_S^q is a subset of Δc^q . At each iteration, a
 455 single column j_i of J_S is used to make progress towards $-\tilde{f}(c^q)$. We generally only execute a few iterations
 456 to mimic the regularization effects of early stopping [53] and truncated-Newton methods [101], as it helps
 457 to prevent overfitting to the linearized subproblem or reaching the undesirable least squares solution as in
 458 LARS [38]. Furthermore, we also terminate early if the decrease in L_2 error is low.

459 Choosing the most geometrically correlated column j_i (as in MBI [28]) allows one to best minimize the
 460 remaining residual; however, small-magnitude columns may require large, undesirable step sizes $\alpha(j_i)$ to make
 461 progress. Instead, motivated by the Gauss-Southwell rule [108], we choose the column j_i that maximizes a
 462 discretized ratio of residual reduction to step size, i.e.

$$\frac{\Delta(r^T r)}{\Delta \alpha} = \frac{\|r(\Delta c_S^q)\|_2^2 - \|r(\Delta c_S^q) - \alpha(j_i)j_i\|_2^2}{|\alpha(j_i)|}, \quad (20)$$

where $r(\Delta c_S^q)$ is the current residual as a function of the current estimate for Δc_S^q . In addition, $\alpha(j_i)$ is the
 step size obtained when choosing column j_i . Flipping all j_i so that $r(\Delta c_S^q)^T j_i > 0$ leads to $\alpha(j_i) > 0$, which
 allows one to equivalently maximize

$$M = \frac{r(\Delta c_S^q)^T r(\Delta c_S^q) - \left(r(\Delta c_S^q)^T r(\Delta c_S^q) - 2\alpha(j_i)r(\Delta c_S^q)^T j_i + (\alpha(j_i))^2 j_i^T j_i \right)}{\alpha(j_i)} \quad (21a)$$

$$= 2r(\Delta c_S^q)^T j_i - \alpha(j_i)\|j_i\|_2^2 \quad (21b)$$

$$= r(\Delta c_S^q)^T j_i + (r(\Delta c_S^q) - \alpha(j_i)j_i)^T j_i. \quad (21c)$$

463 The greedy choice of $\alpha(j_i)$ removes as much of the residual as possible, setting $\alpha(j_i)j_i = (r(\Delta c_S^q) \cdot \hat{j}_i)\hat{j}_i$ or

$$\alpha(j_i) = (r(\Delta c_S^q) \cdot j_i) / \|j_i\|_2^2, \quad (22)$$

464 which zeros out the second term in Equation 21c leaving only $r(\Delta c_S^q)^T j_i$, i.e. gain correlation.

465 There are two subtleties to consider regarding Equations 20–22. First, we do not necessarily use columns
 466 with the largest gains because, as discussed in Section 5.1, we prune away poorly geometrically correlated
 467 columns before considering Equations 20–22. Second, one typically limits the size of $\alpha(j_i)$ when training
 468 neural networks and/or solving optimization/control problems, see e.g. trust region methods [124, 45, 107],
 469 adaptive step sizes for temporal numerical integration [43, 47], and adaptive learning rate techniques such
 470 as Adam [72], ADADELTA [139], etc. Thus, shorter j_i will not necessarily yield the greedy $\alpha(j_i)$ shown in
 471 Equation 22, leaving the second term in the last line of Equation 21c non-zero. See Figure 17.

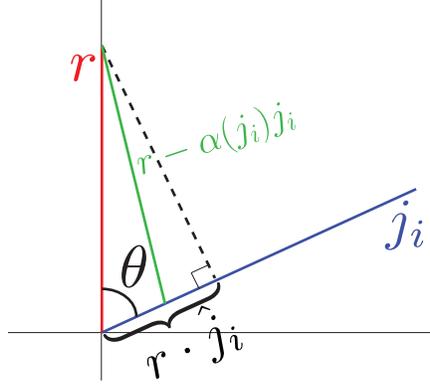


Figure 17: For longer j_i , the $\alpha(j_i)$ required to take the greedy step (Equation 22) will be small enough such that it is not clamped via various safe set or trust region considerations, resulting in the second term in Equation 21c being identically zero. However, for shorter j_i , $\alpha(j_i)$ may be clamped, resulting in an $r(\Delta c_S^q) - \alpha(j_i)j_i$ which is not perpendicular to j_i (shown in green in the figure). In this case, the second term in Equation 21c is non-zero, and the gain correlation of the remaining residual $r(\Delta c_S^q) - \alpha(j_i)j_i$ with the search direction j_i penalizes search directions that become poorly gain correlated after using them. With regard to Figure 16, this prefers a scenario using $a_1/2$ and a_3 as in Figure 16 (Right) as opposed to using a_1 and a_2 in Figure 16 (Left).

472 Consider bounding the step size $\alpha(j_i)$ from above with some α_{\max} . One can choose a reference frame
 473 such that $r(\Delta c_S^q)$ is a unit vector along the y -axis and j_i is in the first quadrant of the xy -plane, as shown
 474 in Figure 18. Referring to the greedy $\alpha(j_i)$ in Equation 22, we plot curves representing vectors j_i where the
 475 greedy $\alpha(j_i)$ is equal to $1/1.5$, 1 , and $1/0.75$ in the figure. When bounding $\alpha(j_i)$ from above by some α_{\max} ,
 476 the $\alpha_G = \alpha_{\max}$ curve represents the boundary between the tips of longer vectors that can use the greedy
 477 $\alpha(j_i)$ and the tips of shorter vectors where $\alpha(j_i)$ would be clamped. In particular, for the $\alpha_{\max} = 1$ case, the
 478 yellow region in the figure represents the tips of longer vectors and the green region represents the tips of
 479 shorter vectors. In the green region, the second term in Equation 21c is added to the usual $r(\Delta c_S^q)^T j_i$ gain
 480 correlation, increasing preference for search directions that remain well-correlated after using them. Figure
 481 18 (Right) shows the magnitude of the second term in Equation 21c. Additionally, one could multiply the
 482 second term in Equation 21c by an arbitrary scaling constant and increase its influence.

483 We also consider the case of clamping based on total progress. Figure 19 uses the same reference frame
 484 as Figure 18 but draws the boundary where $\|j_i\|_2 = \|r(\Delta c_S^q)\|_2/2$. Figure 19 (Left) draws three search
 485 directions taking the greedy step. Regardless of the length of j_i , the greedy $\alpha(j_i)$ rescales such that $\alpha(j_i)j_i$
 486 ends at the boundary between the green and yellow regions where $r(\Delta c_S^q) - \alpha(j_i)j_i$ is orthogonal to j_i .
 487 Figure 19 (Right) shows how clamping the progress limits the ability of a search direction to take the greedy
 488 step, resulting in the second term in Equation 21c being non-zero. Another way of choosing $\alpha(j_i)$ is based
 489 on the observation that $r(\Delta c_S^q) - \alpha(j_i)j_i$ is always less geometrically correlated with j_i than $r(\Delta c_S^q)$ is, since
 490 it points to the left instead of upwards. Hence, one could choose α in order to bound how much worse the
 491 gain/geometric correlation of $r(\Delta c_S^q) - \alpha(j_i)j_i$ is allowed to have compared to that of $r(\Delta c_S^q)$. In general,
 492 there are many potential strategies, but in all such cases, our methodology is to first prune so that large gain
 493 correlation does not introduce poorly geometrically correlated vectors, and then to consider correlation of
 494 the new residual $r(\Delta c_S^q) - \alpha(j_i)j_i$ in addition to correlation of the current residual in order to favor scenarios
 495 like Figure 16 (Right) over Figure 16 (Left).

496 5.3. Examples

497 We consider the problem of determining parameters w that best match a three-dimensional synthetic
 498 face model to a real image, as discussed in Section 2. Although we used CSS to generate Figure 4 and for
 499 related efforts, here we consider a slightly modified situation in order to better isolate and demonstrate the
 500 behavior of CSS. Let w represent the controls for the face blendshapes, jaw angles, and jaw translation, and
 501 $x(w)$ represent the synthetic three-dimensional face surface obtained from w . We replace the inference-based
 502 neural network keypoint detector with a more deterministic artist-drawn rotoscoping of curves for the eyes
 503 and mouth, as shown in Figure 20 (Left). In order to generate comparable keypoints on the synthetic face

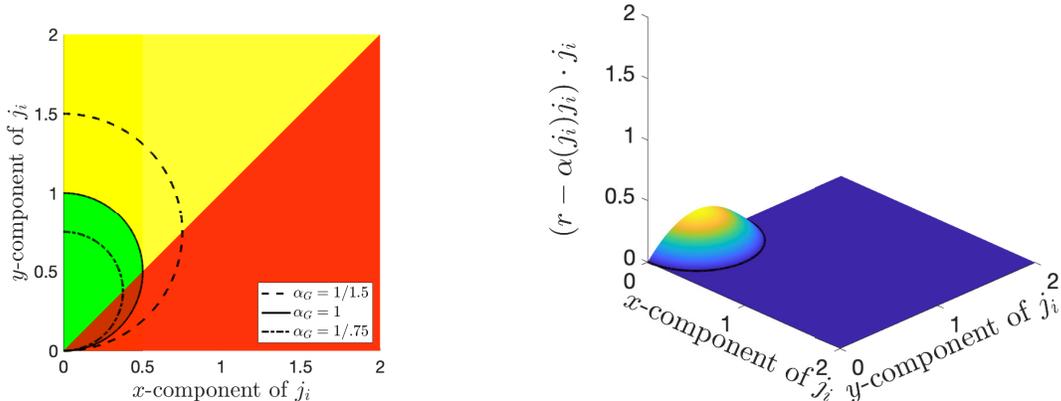


Figure 18: We choose a reference frame where $r(\Delta c_S^q)$ is a unit vector along the y -axis and j_i is in the first quadrant of the xy -plane. (Left) Poorly geometrically correlated vectors (those with tips in the red region) are pruned as in Section 5.1. Referring to the greedy $\alpha(j_i)$ in Equation 22, we plot curves representing the tips of vectors j_i where the greedy $\alpha(j_i)$ is equal to $1/1.5$, 1 , and $1/.75$. When bounding $\alpha(j_i)$ from above by α_{\max} , the $\alpha_G = \alpha_{\max}$ curve represents the boundary between the tips of longer vectors that can use the greedy $\alpha(j_i)$ and the tips of shorter vectors where $\alpha(j_i)$ would be clamped. In particular, for the $\alpha_{\max} = 1$ case, the yellow region represents the tips of longer vectors and the green region represents the tips of shorter vectors. (Right) The magnitude of the second term in Equation 21c for the $\alpha_{\max} = 1$ case. Note that it is non-zero only for shorter vectors with tips inside the $\alpha_G = 1$ curve.

504 model, we draw corresponding curves barycentrically embedded on the three-dimensional face geometry.
 505 Then, $x(w)$ determines the three-dimensional location of these barycentrically embedded curves, which
 506 subsequently are projected into the image plane using calibrated camera intrinsic and extrinsic parameters
 507 [62]; this simple projection replaces the differentiable renderer. See Figure 20 (Right). In order to obtain
 508 comparable keypoints, we label easily-identifiable locations on both sets of curves (i.e. those drawn on the
 509 synthetic model and those drawn on the image), e.g. corners of the mouth and eyes, middles of the lips,
 510 etc. To increase the number of comparable keypoints, we uniformly sample between the projected (into the
 511 image plane) locations of the easily-identifiable keypoints. Letting C^* be the two-dimensional keypoints on
 512 the real image and $C(x(w))$ be the corresponding projected keypoints determined by the parameters w of
 513 the synthetic model, we then solve

$$\min_w \|C^* - C(x(w))\|_2^2 \quad (23)$$

514 in order to recover the parameters w that best match the two sets of keypoints together.

515 For comparison against CSS, we consider solving Equation 23 using Dogleg [112, 94] with no prior, Dogleg
 516 with a prior weight of $\lambda = 3600$, and BFGS [107] with a soft- L_1 prior with a weight of 3600 (i.e. with an
 517 extra term $3600 \sum_i 2(\sqrt{1 + w_i^2} - 1)$ [27]). When solving with CSS, we first prune all columns whose angle
 518 to the residual has an absolute cosine less than 0.3. Then, $\alpha(j_i)$ is set to a fixed size of 0.01 and coordinate
 519 descent is run until the linear L_2 error no longer sufficiently decreases or when over 10 coordinates are
 520 used. We limit all four methods to at most 10 Gauss-Newton linearization iterations. Figure 21 shows the
 521 results. CSS and methods using regularization give the most reasonable geometric results. A major benefit
 522 of CSS is the resulting sparsity of the weights: while Dogleg with and without regularization sets nearly all
 523 the parameters to a non-zero value, CSS generally uses only a small number of non-zero weights. The soft
 524 L_1 regularized solution is sparser than the L_2 regularized solution; however, due to approximations in the
 525 chosen optimization approach (BFGS, Soft L_1), it produces many small (i.e. $< 1 \times 10^{-3}$) weights instead of
 526 identically zero values. While one could clamp small values to zero, care must be taken to not accidentally
 527 clamp weights that contribute significantly to the overall solution.

528 To further elucidate the performance of these approaches, we construct a known exact solution and
 529 subsequently add an increasing amount of noise. First, an exact set of keypoints is determined by subsampling
 530 the contours on the three-dimensional face geometry. Then, a smile expression is created by setting two
 531 specific components of w to 1 while the other components remain set to zero. This known w^* determines

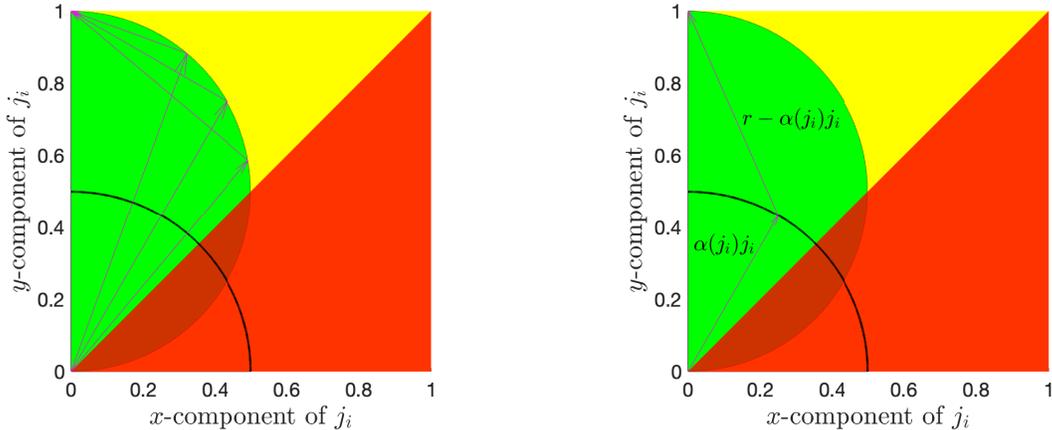


Figure 19: We use the same reference frame as in Figure 18 but draw the boundary where $\|j_i\|_2 = \|r(\Delta c_S^q)\|_2/2$. The yellow and green regions are shaded as in the $\alpha_{\max} = 1$ case from Figure 18. (Left) Three search directions are drawn taking the greedy step. Regardless of the length of j_i , the greedy $\alpha(j_i)$ rescales such that $\alpha(j_i)j_i$ ends at the boundary between the green and yellow regions where $r(\Delta c_S^q) - \alpha(j_i)j_i$ is orthogonal to j_i . (Right) Clamping progress limits the ability of a search direction to take the greedy step, resulting in the second term in Equation 21c being non-zero.

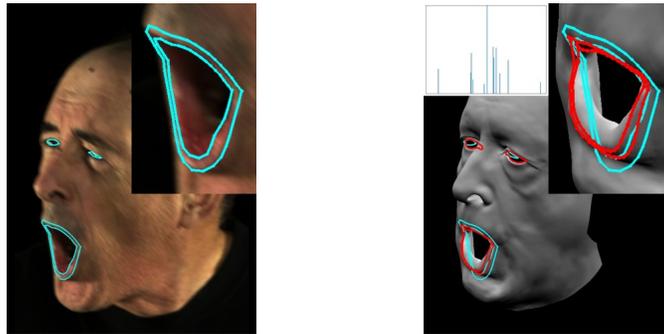


Figure 20: (Left) Hand-drawn rotscope curves on a real image. (Right) Barycentrically embedded curves on the three-dimensional geometric facial model are projected into the image plane (red) and compared to curves drawn/rotscoped on the real image (blue). The inset shows the facial parameters w for this pose.

532 face geometry $x(w^*)$ along with barycentrically embedded keypoints that can be projected into the image
 533 plane to determine C^* . The results obtained are shown in Figures 22a and 22d and are similar to those
 534 obtained using the real image data in Figure 21. Next, we add an increasing amount of uniformly distributed
 535 noise to C^* in the image plane. As expected, Dogleg with no regularization produces reasonable results when
 536 the noise is low but begins to overfit to the erroneous data as the amount of noise increases, see Figure 22
 537 (first row). Both of the regularized approaches as well as CSS are able to target the noisy curves without
 538 overfitting, producing more reasonable geometry. The right half of Figure 22 shows that CSS yields sparser,
 539 more semantic solutions even with the added noise. Tables 1 and 2 demonstrate quantitative results for these
 540 examples. As seen in Table 1, CSS performs the best in all cases as measured by various metrics. Table 2
 541 uses two sparsity measures: the l^0 metric counts how many facial parameters are strictly 0, and the Gini
 542 metric is $1 - 2 \sum_{i=1}^{\hat{w}_i/\|\hat{w}\|_1} (\frac{N-i+0.5}{N})$, where \hat{w} are the sorted parameters with \hat{w}_i the i^{th} largest [66]. Note
 543 how regularization improves the Gini metric, but does not necessarily improve the l^0 metric.

544 5.4. Parameter Study

545 *Column Choice.* Returning to the solution of Equation 23 for the real image data, we compare our approach
 546 for choosing the next coordinate descent column (Section 5.2) to using Gauss-Southwell and MBI. For each
 547 approach, we linearize and solve with no thresholding for the relative decrease in L_2 error, an upper limit
 548 of 10 unique coordinates used, and a fixed step size of 0.01; in these examples, we remove the eye rotscope

Table 1: Comparing the accuracy of estimating the facial parameters in the synthetic tests under various metrics. CSS produces the best results regardless of noise and metric.

Method	L_2 Error			L_1 Error			EMD [116] Error		
	No Noise	0.005	0.01	No Noise	0.005	0.01	No Noise	0.005	0.01
Dogleg	2.578	8.815	20.325	19.227	70.852	157.22	0.128	0.485	1.07
Dogleg+ L_2	0.972	0.952	1.209	4.954	5.324	5.704	0.034	0.036	0.039
BFGS+Soft L_1	0.923	0.91	1.023	3.139	3.057	3.359	0.0215	0.021	0.023
CSS	0.741	0.392	0.509	2.208	0.99	1.08	0.015	0.007	0.007

Table 2: The sparsity of the results of the synthetic tests using common sparsity metrics (a larger number is better).

Method	l^0 Metric			Gini Metric		
	No Noise	0.005	0.01	No Noise	0.005	0.01
Dogleg	21	21	21	0.628	0.580	0.607
Dogleg+ L_2	21	21	21	0.807	0.745	0.739
BFGS+Soft L_1	21	21	21	0.913	0.905	0.916
CSS	128	137	140	0.949	0.974	0.978

549 curves from the energy function and only consider curves drawn around the mouths on the image and model.
550 Results are shown in Figure 23. MBI overfits and overuses mouth blendshapes, e.g. the two most heavily
551 weighted shapes have magnitudes of 85.78 and 63.12. On the other hand, Gauss-Southwell and CSS keep
552 the parameters within a reasonable range while maintaining the sparsity of the solution. We note that with
553 coordinate descent it is generally a matter of when, not if, the algorithm chooses a coordinate that will be
554 overused/overweighted; our examples demonstrate that MBI chooses those coordinates more quickly than
555 Gauss-Southwell and CSS.

556 *Step Size & Convergence.* Since the problem has been normalized so that the $\alpha(j_i)$ generally make most
557 sense between 0 and 1, here we compare fixed step sizes of $\alpha(j_i) = 0.01, 0.02, 0.1, 0.5$ and 1.0 to the full,
558 greedy step in Equation 22. Without pruning, we run 10 Gauss-Newton iterations with no thresholding for
559 the relative decrease in L_2 error and an upper limit of 10 unique coordinates used. We find that smaller step
560 sizes achieve better overall facial shapes and less overused parameters (see Figure 24). In particular, the
561 greedy step sets 7 parameters to be greater than 1 while step sizes of 0.02 and 0.01 only set 4. Removing the
562 eye rotopscope curves causes the overused parameters to disappear; however, as seen in Figure 25, the greedy
563 step causes the mouth to move unnaturally. This would seem to indicate that always taking the greedy step
564 will result in some overfitting.

565 We also compare the effect of using fixed step sizes in Equation 20 versus the full, greedy step size
566 equivalent to Gauss-Southwell without pruning. To isolate this variable, we run 10 Gauss-Newton iterations
567 with no thresholding for the relative decrease in L_2 error and an upper limit of 10 unique coordinates used.
568 We vary $\alpha(j_i)$ in Equation 20 but set the actual step size taken to be fixed at 0.01. As shown in Figure 26,
569 while the resulting geometry and weights are all similar, our approach of allowing the step size to influence
570 the chosen coordinate allows the optimization to more quickly reduce the error in earlier Gauss-Newton
571 iterations than when using Gauss-Southwell (see Figure 27). Therefore, it may be beneficial to use CSS with
572 a fixed size step when only a few Gauss-Newton iterations are desired.

573 *Pruning.* We rescale r to $\hat{r} = r/||r||$ and then compare different threshold values for pruning: 0.0 (no
574 pruning), 0.2, 0.3, and 0.5. We run 10 Gauss-Newton iterations with a step size of 0.01 with no thresholding
575 for the relative decrease in L_2 error. To emphasize the effect of pruning, we allow up to 50 unique coordinates
576 per linearization, and focus only on the rotopscope curves around the mouth. With little to no pruning the
577 model overfits and the geometry around the mouth deforms unreasonably. As the pruning threshold increases,

578 the geometry becomes more regularized and the facial parameters are sparser, as the optimization is forced
579 to use only the most correlated directions. See Figure 28. However, we caution that too much pruning causes
580 MBI style column choices.

581 **6. Conclusions**

582 In difficult nonlinear problems such as the one described in Section 2, one often solves linear subproblems
583 to make progress. Although PCA is quite popular for solving such problems, especially when there are issues
584 with null spaces and the right-hand side not being in the range of the linearized system, we showed that our
585 iLM method not only efficiently monotonically converges to the exact solution of the linearized subproblem,
586 but does so more smoothly. We subsequently pointed out that the larger singular values of the coefficient
587 matrix can be less important than considering which controls are optimal for obtaining the right-hand side.
588 These considerations motivated our column space search (CSS) approach. We chose a complex real-world
589 problem, estimating three-dimensional facial expressions from a mere eight contours drawn on a single two-
590 dimensional RGB image, that allows even non-experts to simply glance at an image and comprehend the
591 effects of noise, overfitting, and regularization. We were able to robustly estimate clean sparse parameter
592 values with good semantic meaning in a highly underconstrained situation where one would typically need
593 significant regularization. In fact, the standard approach without regularization was wildly inaccurate, and
594 although regularization helped to moderate the overall face shape, it excited almost every parameter in the
595 model, clouding semantic interpretation.

596 **Acknowledgements**

597 Research supported in part by ONR N00014-13-1-0346, ONR N00014-17-1-2174, ARL AHPCRC W911NF-
598 07-002, and generous gifts from Amazon and Toyota. In addition, we would like to thank both Reza and
599 Behzad at ONR for supporting our efforts into computer vision and machine learning, as well as Cary Phillips
600 and Industrial Light & Magic for supporting our efforts into facial performance capture. M.B. was supported
601 in part by The VMWare Fellowship in Honor of Ole Agesen. We would also like to thank Paul Huston for
602 his acting.

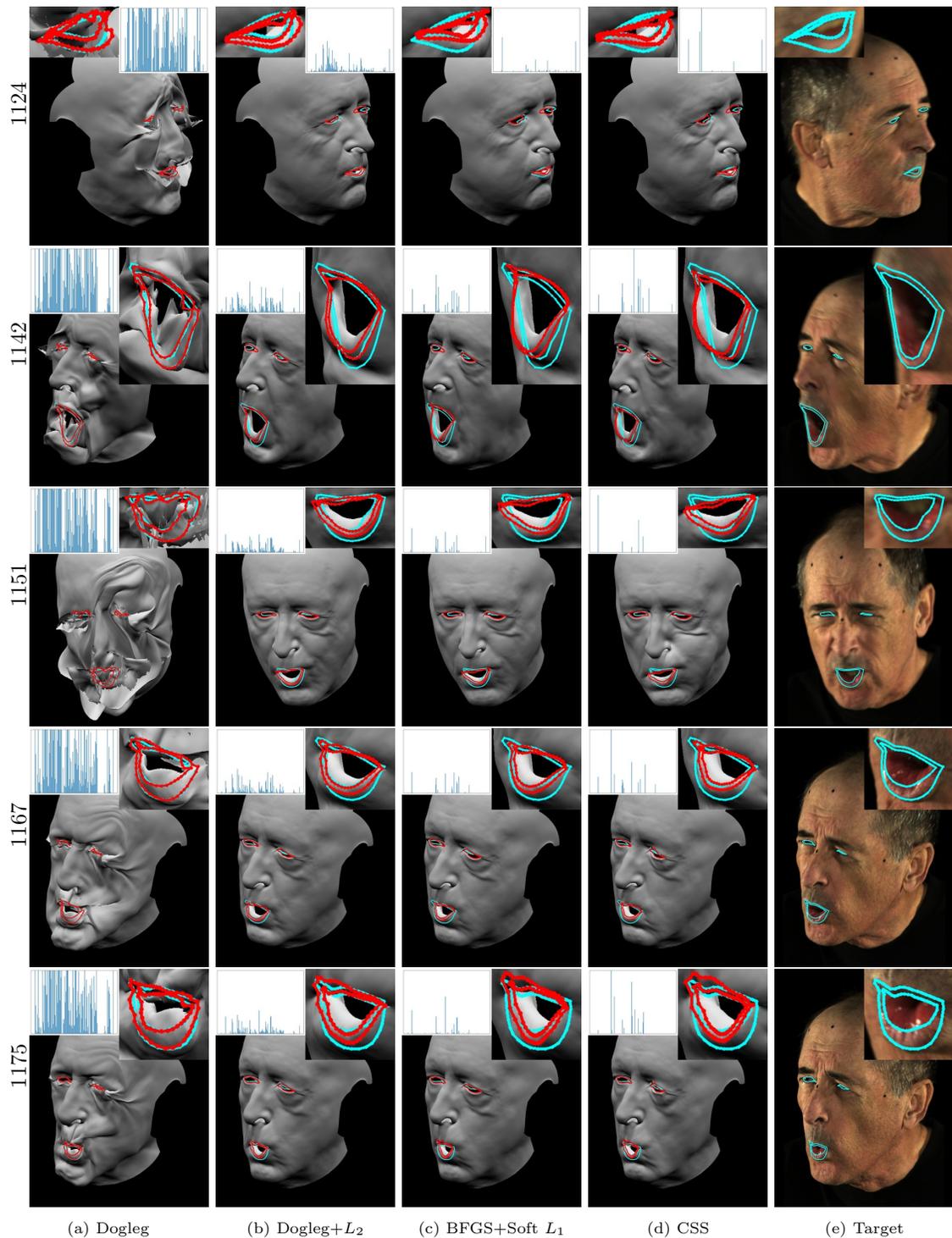


Figure 21: Dogleg without regularization clearly overfits to the curves, producing highly unrealistic face shapes. Dogleg with regularization performs better but sometimes overfits as well. This could be tuned by increasing the regularization weight at the cost of potentially damping out the performance. Our approach produces facial expressions that are reasonably representative of the captured image. The inset bar plots demonstrate the sparsity of the weights for each of the methods. Our method generally produces the sparsest set of weights; e.g. in frame 1142, our method has 12 non-zero parameter values while L_2 regularization produces fully dense results and soft L_1 regularization has 49 significant parameter values (i.e. $> 1 \times 10^{-3}$).

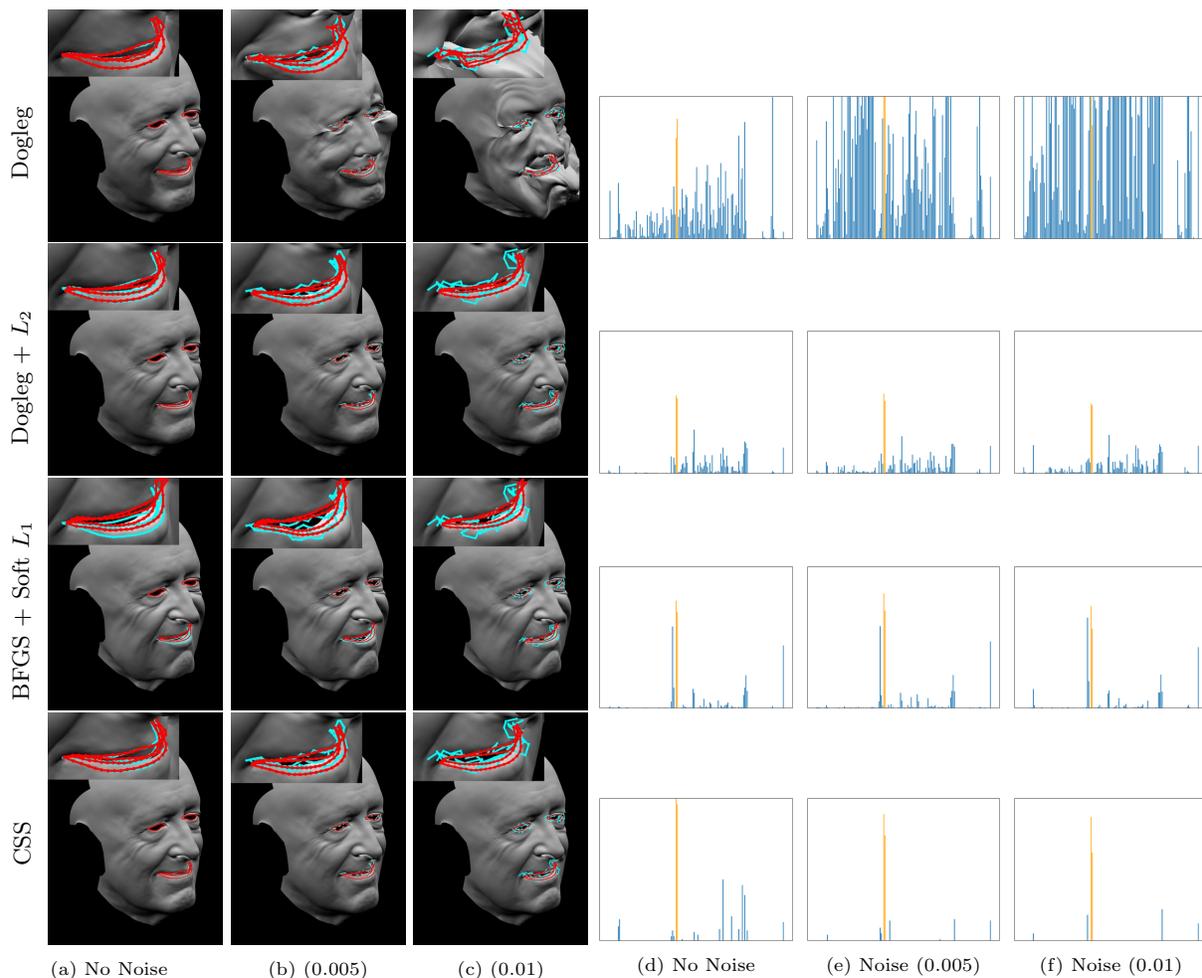


Figure 22: A synthetic test where a known w^* is used to create blue target curves. (*Left*) As we increase the amount of noise added to the points on the blue target curve, the Dogleg method without regularization overfits causing the mesh to “explode” in spite of having the smallest error as measured by Equation 23 (typical of overfitting). On the other hand, both standard regularization and our approach prevent the model from overfitting to the noisy curves. (*Right*) The corresponding facial parameters. The target solution was generated by setting the two orange columns to one and the blue columns to zero. The figure heights are clipped at 1.0 and many parameter values exceed that. Though the regularized solves have smaller, spurious weights than the non-regularized version (second and third row vs. first row), our approach (last row) produces a much sparser solution with more semantic meaning even in the presence of noise.

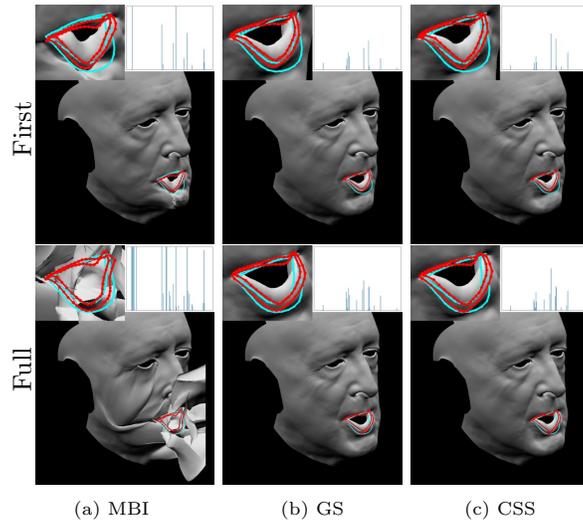


Figure 23: A comparison of the coordinates chosen by the MBI rule, the Gauss-Southwell (GS) rule, and CSS when solving without the eye rosetope curves. The top row are the results after a single Gauss-Newton iteration, and the bottom row are the results after 10 Gauss-Newton iterations.

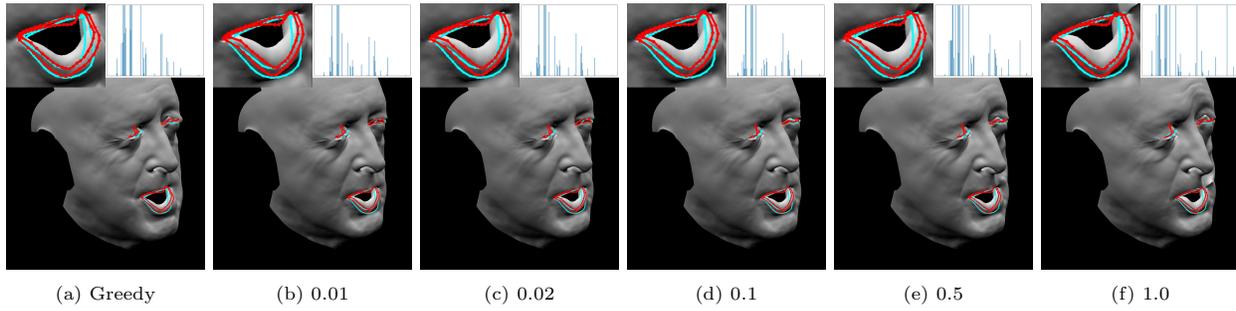


Figure 24: We compare the behavior of the geometry and the facial parameters when using different step sizes.

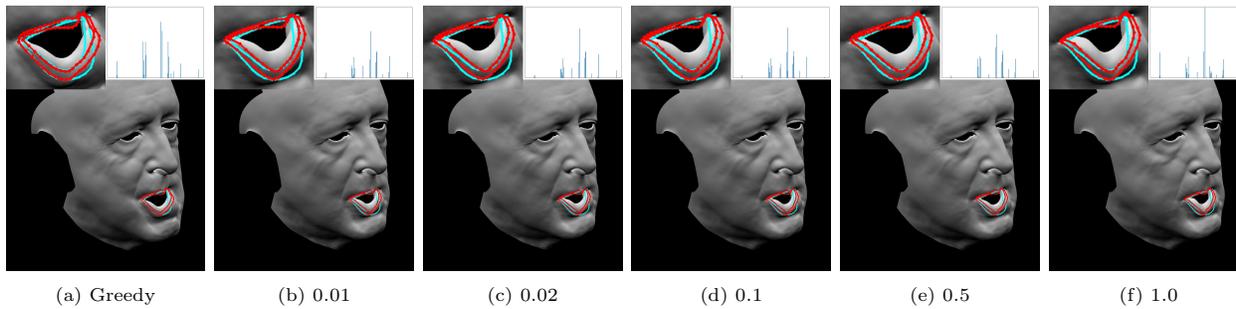


Figure 25: A comparison of the behavior of the geometry and the facial parameters when using different step sizes without the eye rosetope curves.



Figure 26: A comparison of the geometry and parameter results from varying the step size used for choosing the next coordinate in CSS.

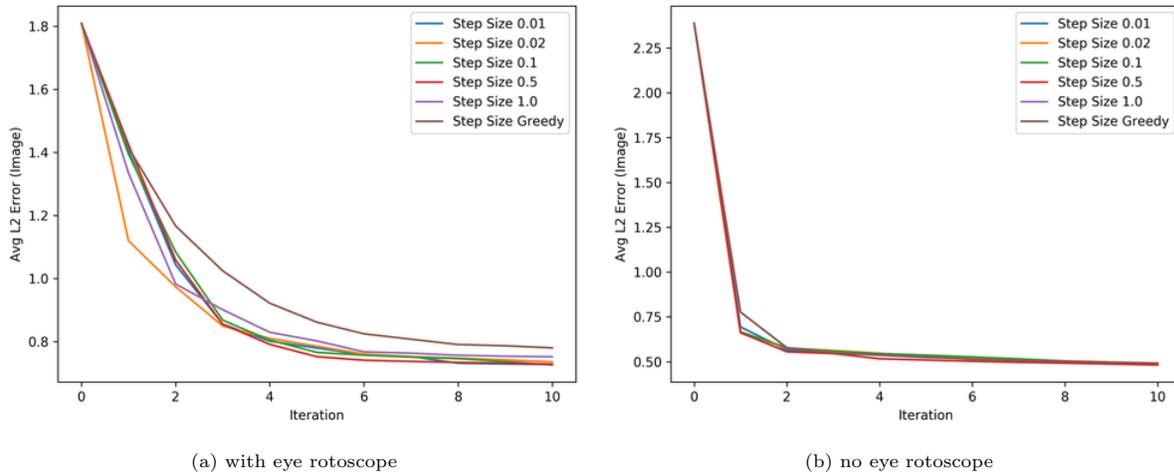


Figure 27: A comparison of the average L_2 errors plotted before every Gauss-Newton iteration when varying the step size used to choose the next coordinate direction in CSS. The brown lines plot the average L_2 errors when using the Gauss-Southwell approach; notice how CSS allows for a faster reduction in error.

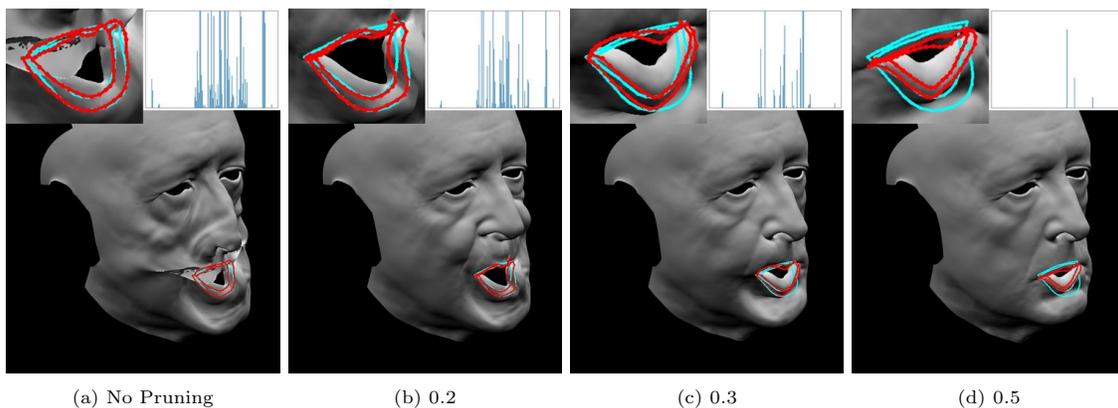


Figure 28: As we increase the threshold for pruning, the resulting solution becomes sparser and more regularized.

603 **References**

- 604 [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S.
605 Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp,
606 Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh
607 Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster,
608 Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay
609 Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan
610 Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
611 Software available from tensorflow.org.
- 612 [2] Frédéric Abergel and Roger Temam. On some control problems in fluid mechanics. *Theoretical and*
613 *Computational Fluid Dynamics*, 1:303–325, 1990.
- 614 [3] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision:
615 A survey. *IEEE Access*, 6:14410–14430, 2018.
- 616 [4] Jose M. Alvarez and Mathieu Salzmann. Learning the number of neurons in deep networks. In D. D.
617 Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information*
618 *Processing Systems 29*, pages 2270–2278. Curran Associates, Inc., 2016.
- 619 [5] James Baglama and Lothar Reichel. Augmented implicitly restarted Lanczos bidiagonalization meth-
620 ods. *SIAM J. Sci. Comput.*, 27(1):19–42, July 2005.
- 621 [6] Michael Bao, Jane Wu, Xinwei Yao, and Ronald Fedkiw. Deep energies for estimating three-dimensional
622 facial pose and expression, 2018, 1812.02899.
- 623 [7] Michael H. Bao, Matthew D. Cong, Stéphane Grabli, and Ronald Fedkiw. High-quality face capture us-
624 ing anatomical muscles. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
625 *(CVPR)*, pages 10794–10803, 2019.
- 626 [8] Michael H. Bao, David Hyde, Xinru Hua, and Ronald Fedkiw. Improved search strategies with appli-
627 cation to estimating facial blendshape parameters, 2020, 1812.02897.
- 628 [9] Thabo Beeler, Bernd Bickel, Paul Beardsley, Bob Sumner, and Markus Gross. High-quality single-shot
629 capture of facial geometry. *ACM Trans. Graph.*, 29(4), July 2010.
- 630 [10] Thabo Beeler, Fabian Hahn, Derek Bradley, Bernd Bickel, Paul Beardsley, Craig Gotsman, Robert W.
631 Sumner, and Markus Gross. High-quality passive facial performance capture using anchor frames.
632 *ACM Trans. Graph.*, 30(4), July 2011.
- 633 [11] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic lan-
634 guage model. *Journal of Machine Learning Research*, 3:1137–1155, 2003.
- 635 [12] Mario Bertero and Patrizia Boccacci. *Introduction to inverse problems in imaging*. CRC press, 1998.
- 636 [13] Kiran S. Bhat, Rony Goldenthal, Yuting Ye, Ronald Mallet, and Michael Koperwas. High fidelity
637 facial animation capture and retargeting with contours. In *Proceedings of the 12th ACM SIG-*
638 *GRAPH/Eurographics Symposium on Computer Animation*, pages 7–14. ACM, 2013.
- 639 [14] Åke Björck. *Numerical methods for least squares problems*. SIAM, 1996.
- 640 [15] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *Proceedings*
641 *of the 26th annual conference on computer graphics and interactive techniques*, pages 187–194. ACM
642 Press/Addison-Wesley Publishing Co., 1999.
- 643 [16] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon
644 Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and
645 Karol Zieba. End to end learning for self-driving cars, 2016, 1604.07316.

- 646 [17] Léon Bottou. *Stochastic Gradient Descent Tricks*, pages 421–436. Springer Berlin Heidelberg, Berlin,
647 Heidelberg, 2012.
- 648 [18] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine
649 learning, 2016, 1606.04838.
- 650 [19] Léon Bottou, Jonas Peters, Joaquin Quiñero Candela, Denis X. Charles, D. Max Chickering, Elon
651 Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning
652 systems: The example of computational advertising. *J. Mach. Learn. Res.*, 14(1):3207–3260, January
653 2013.
- 654 [20] Sofien Bouaziz, Yangang Wang, and Mark Pauly. Online modeling for realtime facial animation. *ACM*
655 *Transactions on Graphics (TOG)*, 32(4):40, 2013.
- 656 [21] Charles G. Broyden. A class of methods for solving nonlinear simultaneous equations. *Mathematics of*
657 *computation*, 19(92):577–593, 1965.
- 658 [22] Charles G. Broyden. Quasi-Newton methods and their application to function minimisation. *Mathe-*
659 *matics of Computation*, 21(99):368–381, 1967.
- 660 [23] Charles G. Broyden. A new double-rank minimisation algorithm. preliminary report. In *Notices of*
661 *the American Mathematical Society*, volume 16, page 670. Amer. Mathematical Soc., 201 Charles St.,
662 Providence, RI 02940-2213, 1969.
- 663 [24] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2D 3D face alignment
664 problem? (and a dataset of 230,000 3D facial landmarks). In *2017 IEEE International Conference on*
665 *Computer Vision (ICCV)*, pages 1021–1030, 2017.
- 666 [25] Chen Cao, Yanlin Weng, Shun Zhou, Yiying Tong, and Kun Zhou. Facewarehouse: A 3D facial ex-
667 pression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*,
668 20(3):413–425, 2014.
- 669 [26] Tony F. Chan and Xue-Cheng Tai. Level set and total variation regularization for elliptic inverse
670 problems with discontinuous coefficients. *Journal of Computational Physics*, 193(1):40–66, 2004.
- 671 [27] Pierre Charbonnier, Laure Blanc-Féraud, Gilles Aubert, and Michel Barlaud. Deterministic edge-
672 preserving regularization in computed imaging. *IEEE Transactions on image processing*, 6(2):298–311,
673 1997.
- 674 [28] Bilian Chen, Simai He, Zhening Li, and Shuzhong Zhang. Maximum block improvement and polynomial
675 optimization. *SIAM Journal on Optimization*, 22(1):87–107, 2012.
- 676 [29] Matthew Cong, Michael Bao, Jane L. E, Kiran S. Bhat, and Ronald Fedkiw. Fully automatic generation
677 of anatomical face simulation models. In *Proceedings of the 14th ACM SIGGRAPH / Eurographics*
678 *Symposium on Computer Animation*, SCA '15, page 175–183, New York, NY, USA, 2015. Association
679 for Computing Machinery.
- 680 [30] Matthew Cong, Lana Lan, and Ronald Fedkiw. Muscle simulation for facial animation in Kong: Skull
681 Island. In *ACM SIGGRAPH 2017 Talks*, SIGGRAPH '17, New York, NY, USA, 2017. Association for
682 Computing Machinery.
- 683 [31] Matthew D. Cong, Kiran S. Bhat, and Ronald Fedkiw. Art-directed muscle simulation for high-end
684 facial animation. In *Proceedings of the 15th ACM SIGGRAPH / Eurographics Symposium on Computer*
685 *Animation (SCA'16)*, pages 119–127, 2016.
- 686 [32] Matthew D. Cong, Lana Lan, and Ronald Fedkiw. Local geometric indexing of high resolution data
687 for facial reconstruction from sparse markers. *arXiv preprint arXiv:1903.00119*, 2019.
- 688 [33] William C Davidon. Variable metric method for minimization. Technical Report ANL-5990, Argonne
689 National Laboratory, 5 1959.

- 690 [34] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc'aurelio Ran-
691 zato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks. In *Advances*
692 *in neural information processing systems*, pages 1223–1231, 2012.
- 693 [35] Shelley Derksen and Harvey J. Keselman. Backward, forward and stepwise automated subset selection
694 algorithms: Frequency of obtaining authentic and noise variables. *British Journal of Mathematical and*
695 *Statistical Psychology*, 45(2):265–282, 1992.
- 696 [36] Fernando Diaz, Donald Metzler, and Sihem Amer-Yahia. Relevance and ranking in online dating
697 systems. In *Proceedings of the 33rd International ACM SIGIR Conference on Research and Develop-*
698 *ment in Information Retrieval, SIGIR '10*, page 66–73, New York, NY, USA, 2010. Association for
699 Computing Machinery.
- 700 [37] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
701 stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- 702 [38] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *The*
703 *Annals of statistics*, 32(2):407–499, 2004.
- 704 [39] Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. *Regularization of inverse problems*, volume
705 375. Springer Science & Business Media, 1996.
- 706 [40] Jinyan Fan. The modified Levenberg-Marquardt method for nonlinear equations with cubic conver-
707 gence. *Mathematics of Computation*, 81(277):447–466, 2012.
- 708 [41] Ronald Fedkiw, Yilin Zhu, Winnie Lin, and Jane Wu. Continuous mathematical methods, emphasizing
709 machine learning, 2020. Stanford CS205L Winter 2020 Lecture Slides.
- 710 [42] Ronald P. Fedkiw, Guillermo Sapiro, and Chi-Wang Shu. Shock capturing, level sets, and PDE based
711 methods in computer vision and image processing: a review of Osher's contributions. *Journal of*
712 *Computational Physics*, 185(2):309 – 341, 2003.
- 713 [43] Erwin Fehlberg. Low-order classical Runge-Kutta formulas with stepsize control and their application
714 to some heat transfer problems. Technical Report Technical Report 315, NASA, 1969.
- 715 [44] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322,
716 1970.
- 717 [45] Roger Fletcher. *Practical Methods Of Optimization*. John Wiley & Sons, 1980.
- 718 [46] Roger Fletcher and Michael J. D. Powell. A rapidly convergent descent method for minimization. *The*
719 *computer journal*, 6(2):163–168, 1963.
- 720 [47] Jessica G. Gaines and Terry J. Lyons. Variable step size control in the numerical solution of stochastic
721 differential equations. *SIAM Journal on Applied Mathematics*, 57(5):1455–1484, 1997.
- 722 [48] Nicholas Geneva and Nicholas Zabaras. Modeling the dynamics of PDE systems with physics-
723 constrained deep auto-regressive networks. *Journal of Computational Physics*, 403:109056, 2020.
- 724 [49] Zhenglin Geng, Daniel Johnson, and Ronald Fedkiw. Coercing machine learning to output physically
725 accurate results. *Journal of Computational Physics*, 406:109099, 2020.
- 726 [50] Abhijeet Ghosh, Graham Fyffe, Borom Tunwattanapong, Jay Busch, Xueming Yu, and Paul Debevec.
727 Multiview face capture using polarized spherical gradient illumination. In *Proceedings of the 2011 SIG-*
728 *GRAPH Asia Conference, SA '11*, New York, NY, USA, 2011. Association for Computing Machinery.
- 729 [51] Frederic Gibou, David Hyde, and Ron Fedkiw. Sharp interface approaches and deep learning techniques
730 for multiphase flows. *Journal of Computational Physics*, 380:442–463, 2019.
- 731 [52] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of*
732 *computation*, 24(109):23–26, 1970.

- 733 [53] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 734 [54] Ariel Gordon, Elad Eban, Ofir Nachum, Bo Chen, Hao Wu, Tien-Ju Yang, and Edward Choi. MorphNet: Fast simple resource-constrained structure learning of deep networks. In *2018 IEEE/CVF*
735 *Conference on Computer Vision and Pattern Recognition*, pages 1586–1595, 2018.
- 737 [55] Thore Graepel, Joaquin Quiñero Candela, Thomas Borchert, and Ralf Herbrich. Web-scale Bayesian
738 click-through rate prediction for sponsored search advertising in Microsoft’s Bing search engine. In
739 *Proceedings of the 27th International Conference on International Conference on Machine Learning*,
740 ICML’10, page 13–20, Madison, WI, USA, 2010. Omnipress.
- 741 [56] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. In *Advances*
742 *in neural information processing systems (NIPS)*, 2016.
- 743 [57] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks
744 with pruning, trained quantization and Huffman coding, 2015, 1510.00149.
- 745 [58] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for
746 efficient neural network. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett,
747 editors, *Advances in Neural Information Processing Systems 28*, pages 1135–1143. Curran Associates,
748 Inc., 2015.
- 749 [59] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning*. Springer
750 series in statistics New York, NY, USA:, 2001.
- 751 [60] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks.
752 In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1389–1397, 2017.
- 753 [61] Michael T. Heath. *Scientific computing: an introductory survey*. SIAM, 2002.
- 754 [62] Janne Heikkila and Olli Silven. A four-step camera calibration procedure with implicit image correction.
755 In *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*,
756 pages 1106–1112. IEEE, 1997.
- 757 [63] Pei-Lun Hsieh, Chongyang Ma, Jihun Yu, and Hao Li. Unconstrained realtime facial performance
758 capture. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
759 1675–1683, 2015.
- 760 [64] Haoda Huang, Jinxiang Chai, Xin Tong, and Hsiang-Tao Wu. Leveraging motion capture and 3D
761 scanning for high-fidelity facial performance acquisition. In *ACM Transactions on Graphics (TOG)*,
762 volume 30, page 74. ACM, 2011.
- 763 [65] Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks
764 on neural network policies, 2017, 1702.02284.
- 765 [66] Niall Hurley and Scott Rickard. Comparing measures of sparsity. *IEEE Transactions on Information*
766 *Theory*, 55(10):4723–4741, 2009.
- 767 [67] Brody Huval, Tao Wang, Sameep Tandon, Jeff Kiske, Will Song, Joel Pazhayampallil, Mykhaylo
768 Andriluka, Pranav Rajpurkar, Toki Migimatsu, Royce Cheng-Yue, Fernando Mujica, Adam Coates,
769 and Andrew Y. Ng. An empirical evaluation of deep learning on highway driving, 2015, 1504.01716.
- 770 [68] Alexandru Eugen Ichim, Sofien Bouaziz, and Mark Pauly. Dynamic 3d avatar creation from hand-held
771 video input. *ACM Transactions on Graphics (ToG)*, 34(4):45, 2015.
- 772 [69] Ameya D. Jagtap, Kenji Kawaguchi, and George Em Karniadakis. Adaptive activation functions ac-
773 celerate convergence in deep and physics-informed neural networks. *Journal of Computational Physics*,
774 404:109136, 2020.

- 775 [70] Antony Jameson, Luigi Martinelli, and Niles A. Pierce. Optimum aerodynamic design using the
776 Navier-Stokes equations. *Theoretical and Computational Fluid Dynamics*, 10:213–237, 1998.
- 777 [71] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio
778 Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In
779 *Proceedings of the 22nd ACM International Conference on Multimedia*, MM '14, pages 675–678, New
780 York, NY, USA, 2014. Association for Computing Machinery.
- 781 [72] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014, 1412.6980.
- 782 [73] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convo-
783 lutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors,
784 *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc.,
785 2012.
- 786 [74] Igor' Anatol'yevich Krylov and Feliks Leonidovich Chernous'ko. Solution of problems of optimal control
787 by the method of local variations. *USSR Computational Mathematics and Mathematical Physics*,
788 6(2):12–31, 1966.
- 789 [75] Lana Lan, Matthew Cong, and Ronald Fedkiw. Lessons from the evolution of an anatomical facial
790 muscle model. In *Proceedings of the ACM SIGGRAPH Digital Production Symposium*, DigiPro '17,
791 New York, NY, USA, 2017. Association for Computing Machinery.
- 792 [76] Jeff Lander. Skin them bones: Game programming for the web generation. *Game Developer Magazine*,
793 5(1):10–18, 1998.
- 794 [77] Rasmus Larsen. Lanczos bidiagonalization with partial reorthogonalization. *DAIMI Report Series*,
795 27(537), Dec. 1998.
- 796 [78] Quoc V. Le, Jiquan Ngiam, Adam Coates, Ahbik Lahiri, Bobby Prochnow, and Andrew Y Ng. On
797 optimization methods for deep learning. 2011.
- 798 [79] Yann LeCun, Bernhard E. Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne E.
799 Hubbard, and Lawrence D. Jackel. Handwritten digit recognition with a back-propagation network.
800 In D. S. Touretzky, editor, *Advances in Neural Information Processing Systems 2*, pages 396–404.
801 Morgan-Kaufmann, 1990.
- 802 [80] Yann LeCun, Léon Bottou, and Yoshua Bengio. Reading checks with multilayer graph transformer net-
803 works. In *1997 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1,
804 pages 151–154 vol.1, 1997.
- 805 [81] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
806 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 807 [82] Yann LeCun, Fu Jie Huang, and Léon Bottou. Learning methods for generic object recognition with
808 invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on*
809 *Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 2, pages II–104 Vol.2, 2004.
- 810 [83] Kenneth Levenberg. A method for the solution of certain non-linear problems in least squares. *Quarterly*
811 *of applied mathematics*, 2(2):164–168, 1944.
- 812 [84] John P. Lewis, Ken Anjyo, Taehyun Rhee, Mengjie Zhang, Fred Pighin, and Zhigang Deng. Prac-
813 tice and Theory of Blendshape Facial Models. In Sylvain Lefebvre and Michela Spagnuolo, editors,
814 *Eurographics 2014 - State of the Art Reports*. The Eurographics Association, 2014.
- 815 [85] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient
816 convnets, 2016, 1608.08710.
- 817 [86] Hao Li, Thibaut Weise, and Mark Pauly. Example-based facial rigging. *ACM Transactions on Graphics*
818 *(TOG)*, 29(4):32, 2010.

- 819 [87] Hao Li, Jihun Yu, Yuting Ye, and Chris Bregler. Realtime facial animation with on-the-fly correctives.
820 *ACM Trans. Graph.*, 32(4):42–1, 2013.
- 821 [88] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing
822 through edge sampling. *ACM Trans. Graph.*, 37(6), December 2018.
- 823 [89] Julia Ling, Reese Jones, and Jeremy Templeton. Machine learning strategies for systems with invariance
824 properties. *Journal of Computational Physics*, 318:22 – 35, 2016.
- 825 [90] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning
826 efficient convolutional networks through network slimming. In *2017 IEEE International Conference
827 on Computer Vision (ICCV)*, pages 2755–2763, 2017.
- 828 [91] Matthew M. Loper and Michael J. Black. Opendr: An approximate differentiable renderer. In *European
829 Conference on Computer Vision*, pages 154–169. Springer, 2014.
- 830 [92] Yifei Lou, Xiaoqun Zhang, Stanley Osher, and Andrea Bertozzi. Image recovery via nonlocal operators.
831 *Journal of Scientific Computing*, 42:185–197, 2010.
- 832 [93] Christos Louizos, Karen Ullrich, and Max Welling. Bayesian compression for deep learning. In I. Guyon,
833 U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances
834 in Neural Information Processing Systems 30*, pages 3288–3298. Curran Associates, Inc., 2017.
- 835 [94] Manolis I. A. Lourakis and Antonis A. Argyros. Is Levenberg-Marquardt the most efficient optimization
836 algorithm for implementing bundle adjustment? In *Tenth IEEE International Conference on Computer
837 Vision (ICCV’05)*, volume 2, pages 1526–1531. IEEE, 2005.
- 838 [95] Rongrong Ma, Jianyu Miao, Lingfeng Niu, and Peng Zhang. Transformed ℓ_1 regularization for learning
839 sparse deep neural networks. *Neural Networks*, 119:286 – 298, 2019.
- 840 [96] Nadia Magnenat-Thalmann, Richard Laperrière, and Daniel Thalmann. Joint-dependent local defor-
841 mations for hand animation and object grasping. In *Proceedings on Graphics Interface ’88*, pages
842 26–33, CAN, 1989. Canadian Information Processing Society.
- 843 [97] Donald W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal of
844 the society for Industrial and Applied Mathematics*, 11(2):431–441, 1963.
- 845 [98] Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity.
846 *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- 847 [99] Brian McFee and Gert Lanckriet. Metric learning to rank. In *Proceedings of the 27th International
848 Conference on International Conference on Machine Learning, ICML’10*, page 775–782, Madison, WI,
849 USA, 2010. Omnipress.
- 850 [100] Abdel-Rahman Mohamed, George E. Dahl, and Geoffrey Hinton. Acoustic modeling using deep belief
851 networks. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1):14–22, 2012.
- 852 [101] Stephen G. Nash. A survey of truncated-Newton methods. *Journal of computational and applied
853 mathematics*, 124(1-2):45–59, 2000.
- 854 [102] Kirill Neklyudov, Dmitry Molchanov, Arsenii Ashukha, and Dmitry P Vetrov. Structured Bayesian
855 pruning via log-normal multiplicative noise. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach,
856 R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing
857 Systems 30*, pages 6775–6784. Curran Associates, Inc., 2017.
- 858 [103] Yurii Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$.
859 *Soviet Mathematics Doklady*, 27:372–376, 1983.
- 860 [104] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM
861 Journal on Optimization*, 22(2):341–362, 2012.

- 862 [105] Thomas Neumann, Kiran Varanasi, Stephan Wenger, Markus Wacker, Marcus Magnor, and Chris-
863 tian Theobalt. Sparse localized deformation components. *ACM Transactions on Graphics (TOG)*,
864 32(6):179, 2013.
- 865 [106] Jorge Nocedal. Updating quasi-Newton matrices with limited storage. *Mathematics of computation*,
866 35(151):773–782, 1980.
- 867 [107] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2006.
- 868 [108] Julie Nutini, Mark Schmidt, Issam Laradji, Michael Friedlander, and Hoyt Koepke. Coordinate descent
869 converges faster with the Gauss-Southwell rule than random selection. In *International Conference on*
870 *Machine Learning*, pages 1632–1641, 2015.
- 871 [109] Jongsoo Park, Sheng Li, Wei Wen, Ping Tak Peter Tang, Hai Li, Yiran Chen, and Pradeep Dubey.
872 Faster CNNs with direct sparse convolutions and guided pruning, 2016, 1608.01409.
- 873 [110] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
874 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang,
875 Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,
876 Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning
877 library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors,
878 *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.,
879 2019.
- 880 [111] John C. Platt. Sequential minimal optimization: A fast algorithm for training support vector machines.
881 Technical Report MSR-TR-98-14, Microsoft Research, April 1998.
- 882 [112] Michael J. D. Powell. A hybrid method for nonlinear equations. In Philip Rabinowitz, editor, *Numerical*
883 *Methods for Nonlinear Algebraic Equations*, pages 87–114. Gordon and Breach, 1970.
- 884 [113] Yinghe Qi, Jiakai Lu, Ruben Scardovelli, Stéphane Zaleski, and Grétar Tryggvason. Computing curva-
885 ture for volume of fluid methods using machine learning. *Journal of Computational Physics*, 377:155–
886 161, 2019.
- 887 [114] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*,
888 12(1):145–151, 1999.
- 889 [115] Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Physics-informed neural networks: A
890 deep learning framework for solving forward and inverse problems involving nonlinear partial differen-
891 tial equations. *Journal of Computational Physics*, 378:686 – 707, 2019.
- 892 [116] Yossi Rubner, Carlo Tomasi, and Leonidas J. Guibas. A metric for distributions with applications to
893 image databases. In *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)*,
894 pages 59–66, 1998.
- 895 [117] Sebastian Ruder. An overview of gradient descent optimization algorithms, 2016, 1609.04747.
- 896 [118] Moktar A. Salama, John A. Garba, Laura A. Demsetz, and Firdaus E. Udwadia. Simultaneous opti-
897 mization of controlled structures. *Computational Mechanics*, 3:275–282, 1988.
- 898 [119] David F. Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of*
899 *computation*, 24(111):647–656, 1970.
- 900 [120] Hao-Jun Michael Shi, Shenyinying Tu, Yangyang Xu, and Wotao Yin. A primer on coordinate descent
901 algorithms, 2016, 1610.00040.
- 902 [121] Eftychios Sifakis, Igor Neverov, and Ronald Fedkiw. Automatic determination of facial muscle activa-
903 tions from sparse motion capture marker data. In *ACM SIGGRAPH 2005 Papers*, SIGGRAPH ’05,
904 page 417–425, New York, NY, USA, 2005. Association for Computing Machinery.

- 905 [122] Eftychios Sifakis, Andrew Selle, Avram Robinson-Mosher, and Ronald Fedkiw. Simulating speech
906 with a physics-based facial muscle model. In *Proceedings of the 2006 ACM SIGGRAPH/Eurographics*
907 *Symposium on Computer Animation*, SCA '06, page 261–270, Goslar, DEU, 2006. Eurographics Asso-
908 ciation.
- 909 [123] Justin Sirignano and Konstantinos Spiliopoulos. DGM: A deep learning algorithm for solving partial
910 differential equations. *Journal of Computational Physics*, 375:1339 – 1364, 2018.
- 911 [124] Danny C. Sorensen. Newton’s method with a model trust region modification. *SIAM Journal on*
912 *Numerical Analysis*, 19(2):409–426, 1982.
- 913 [125] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural
914 networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, 2019.
- 915 [126] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization
916 and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147,
917 2013.
- 918 [127] Theano Development Team. Theano: A Python framework for fast computation of mathematical
919 expressions. *arXiv e-prints*, abs/1605.02688, May 2016.
- 920 [128] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner.
921 Face2Face: Real-time face capture and reenactment of RGB videos. In *Proceedings of the IEEE*
922 *Conference on Computer Vision and Pattern Recognition*, pages 2387–2395, 2016.
- 923 [129] Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5: rmsprop: Divide the gradient by a running average
924 of its recent magnitude. *Coursera: Neural networks for machine learning*, 4(2):26–31, 2012.
- 925 [130] Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization.
926 *Journal of optimization theory and applications*, 109(3):475–494, 2001.
- 927 [131] Marko Vauhkonen, Dénes Vadász, Pasi A. Karjalainen, Erkki Somersalo, and Jari P. Kaipio. Tikhonov
928 regularization and prior information in electrical impedance tomography. *IEEE Transactions on Med-*
929 *ical Imaging*, 17(2):285–293, 1998.
- 930 [132] Min Wang, Siu Wun Cheung, Wing Tat Leung, Eric T. Chung, Yalchin Efendiev, and Mary Wheeler.
931 Reduced-order deep learning for flow dynamics. the interplay between deep learning and model reduc-
932 tion. *Journal of Computational Physics*, 401:108939, 2020.
- 933 [133] Wei Wen, Yiran Chen, Hai Li, Yuxiong He, Samyam Rajbhandari, Minjia Zhang, Wenhan Wang, Fang
934 Liu, and Bin Hu. Learning intrinsic sparse structures within long short-term memory. In *ICLR 2018*
935 *Conference*, February 2018.
- 936 [134] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in
937 deep neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors,
938 *Advances in Neural Information Processing Systems 29*, pages 2074–2082. Curran Associates, Inc.,
939 2016.
- 940 [135] Stephen J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- 941 [136] Chenglei Wu, Derek Bradley, Markus Gross, and Thabo Beeler. An anatomically-constrained local
942 deformation model for monocular face capture. *ACM Transactions on Graphics (TOG)*, 35(4):115,
943 2016.
- 944 [137] Huanrui Yang, Wei Wen, and Hai Li. Deepfayer: Learning sparser neural network with differentiable
945 scale-invariant sparsity measures. In *International Conference on Learning Representations*, 2020.

- 946 [138] Jihun Yun, Peng Zheng, Eunho Yang, Aurelie Lozano, and Aleksandr Aravkin. Trimming the ℓ_1
947 regularizer: Statistical analysis, optimization, and applications to deep learning. In Kamalika Chaud-
948 huri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine*
949 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7242–7251, Long Beach,
950 California, USA, 09–15 Jun 2019. PMLR.
- 951 [139] Matthew D. Zeiler. ADADELTA: An adaptive learning rate method, 2012, 1212.5701.
- 952 [140] Xiaoqun Zhang, Martin Burger, Xavier Bresson, and Stanley Osher. Bregmanized nonlocal regulariza-
953 tion for deconvolution and sparse reconstruction. *SIAM Journal on Imaging Sciences*, 3(3):253–276,
954 2010.
- 955 [141] Olgierd Cecil Zienkiewicz and Robert Leroy Taylor. *The finite element method, volume 1: the basis*,
956 volume 1. Butterworth-Heinemann, 2000.
- 957 [142] Olgierd Cecil Zienkiewicz and Robert Leroy Taylor. *The finite element method, volume 2: solid*
958 *mechanics*, volume 2. Butterworth-Heinemann, 2000.
- 959 [143] Gaspard Zoss, Derek Bradley, Pascal Bérard, and Thabo Beeler. An empirical rig for jaw animation.
960 *ACM Transactions on Graphics (TOG)*, 37(4):1–12, 2018.