

---

# Global Convergence Newton

---

**Raffael Colonnello**  
University of Basel  
Raffael.Colonnello@unibas.ch

**Fynn Gohlke**  
University of Basel  
Fynn.Gohlke@stud.unibas.ch

**Benedikt Heuser**  
University of Basel  
ben.heuser@unibas.ch

## Abstract

1       The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and  
2       right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points.  
3       The word **Abstract** must be centered, bold, and in point size 12. Two line spaces  
4       precede the abstract. The abstract must be limited to one paragraph.

## 5   1   Introduction

6   In this paper we consider problems of the form

$$\min_{x \in \mathbb{R}^d} f(\mathbf{x}) \tag{1}$$

7   where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is a twice-differentiable function. First-order optimization methods are widely  
8   used for such problems due to their low per-iteration computational cost and their suitability for  
9   parallelization. They often suffer from slow convergence for ill-conditioned objective functions [1].  
10   Newton’s method is a popular optimization algorithm that is commonly used to solve optimization  
11   problems. It is a second-order optimization algorithm since it uses second-order information of  
12   the objective function. Newton’s method is known to have fast local convergence guarantees for  
13   convex functions. However, the global convergence properties of Newton’s method are still an  
14   active area of research [2] [3]. In contrast to first-order methods like gradient descent, second-order  
15   methods, such as Newton’s method can achieve much faster convergence when presented with ill  
16   conditioned Hessians by transferring the problem into a more isotropic optimization problem at the  
17   cost of an increase to cubic run time. Newton’s method yields local quadratic convergence if  $f$  is  
18   twice differentiable (or we have suitable regularity conditions), which degrade outside of the local  
19   regions, yielding up to sublinear global convergence guarantees, depending on the algorithm.

20   In this paper, we explore the theoretical foundations of several Newton-type methods that achieve  
21   different global convergence guarantees, compare their performance in a classification-type problem  
22   for two loss functions on four different datasets. Finally we will propose two modifications of the  
23   algorithms to achieve an increase in runtime, by either coupling the Newton-type method with a  
24   conjugate gradient method for Hessian vector multiplication or Strassen’s algorithm for fast matrix  
25   inversion.

## 2 Background

### 2.1 Loss function and Datasets

Let  $X = \begin{bmatrix} \dots x_1^\top \dots \\ \vdots \\ \dots x_i^\top \dots \\ \vdots \\ \dots x_n^\top \dots \end{bmatrix} \in \mathbb{R}^{n \times d}$  be the set of data for  $n$  datapoints with  $d$  features, i.e.  $x_i \in \mathbb{R}^d$  and labels  $y^\top = [y_1, \dots, y_n]$

For  $\sigma(x) := \frac{\exp(x)}{1+\exp(x)}$  the loss functions w.r.t. weights  $\omega$  are given by

$$L_1(\omega) = -\frac{1}{n} \sum_{i=1}^n \left( y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \right), \quad \hat{y}_i = \sigma(x_i^\top \omega) \quad (2)$$

$$L_2(\omega) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^\top \omega)) + r(\omega), \quad r(\omega) = \lambda \sum_{j=1}^d \frac{\alpha \omega_j^2}{1 + \alpha \omega_j^2} \quad (3)$$

which yields the two optimization problems

$$\min_{\omega} L_1(\omega) \quad (4)$$

$$\min_{\omega} L_2(\omega) \quad (5)$$

*Remark 1: The 0-1 loss function for logistic regression is given by*

$$-\sum_{i=1}^N \log \left[ \mu_i^{\mathbb{I}(y_i=1)} (1 - \mu_i)^{\mathbb{I}(y_i=0)} \right] = -\sum_{i=1}^N [y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)]$$

for labels  $y_i \in \{0, 1\}$  [4, Eq. 8.2–8.3]. If we instead use labels  $\tilde{y}_i \in \{-1, +1\}$ , the negative log-likelihood becomes

$$\sum_{i=1}^N \log(1 + \exp(-\tilde{y}_i \mathbf{w}^\top \mathbf{x}_i))$$

[4, Eq. 8.4]. To ensure the loss functions correspond to the correct likelihood, the label encoding must match the loss form [4, Sec. 8.3.1]. Consequently labels were adapted conditioned to meet the loss functions requirements.

The corresponding gradients of  $L_i$  are

$$\nabla L_1(x) = \frac{1}{n} X^\top (\hat{y} - y) \quad (6)$$

$$\nabla L_2(x) = -\frac{1}{n} X^\top (y \odot \sigma(-y \odot (X\omega))) + \nabla r(x) \quad (7)$$

with  $\nabla r(\omega)^\top = \lambda \left[ \frac{2\alpha\omega_1}{(1+\alpha\omega_1^2)^2}, \dots, \frac{2\alpha\omega_d}{(1+\alpha\omega_d^2)^2} \right]$ , where  $\sigma(\cdot)$  is applied elementwise, and  $\odot$  denotes the entrywise multiplication of vectors.

Differentiating again yields the Hessians

$$\nabla^2 L_1(\omega) = \frac{1}{n} X^\top D(\omega) X \quad (8)$$

$$\nabla^2 L_2(\omega) = \frac{1}{n} X^\top D(\omega) X + \nabla^2 r(\omega), \quad \nabla^2 r(\omega) = \text{diag} \left( \lambda \frac{2\alpha(1 - 3\alpha\omega_j^2)}{(1 + \alpha\omega_j^2)^3} \right) \quad (9)$$

where the diagonal matrix  $D(\omega)$  has entries

$$D_{ii}(\omega) = \hat{y}_i(1 - \hat{y}_i) = \sigma(-y_i x_i^\top \omega)(1 - \sigma(-y_i x_i^\top \omega)), \quad (10)$$

In order to discuss the Algorithms assumptions and conditions in the later sections of this paper we will first state some properties of the given problems.

0. Symmetry of Hessian: It is easy to verify that the Hessians of both loss functions are symmetric as

$$(X^\top DX)^\top = (X^\top D^\top (X^\top)^\top) = X^\top DX$$

and for  $\nabla^2 L_2$  symmetry is even more trivial as we only add a diagonal matrix  $\nabla^2 r(\omega)$  onto the first one.

1. Invertibility:

The matrix  $\nabla^2 L_1(\omega) = X^\top D(\omega)X$  is invertible if and only if  $X$  has full rank.

*Proof:* " $\implies$ " (by contrapositive) Assume that  $X$  doesn't have full rank, then by the definition of rank it holds, that

$$\exists y \neq 0 \quad \text{s.t.} \quad Xy = 0 \implies X^\top DX \underbrace{Xy}_{=0} = 0$$

$$\implies X^\top DX \text{ rank deficient} \implies X^\top DX \text{ not invertible}$$

" $\Leftarrow$ " Assume  $X$  has full rank and let  $y \neq 0$ . Then  $Xy \neq 0 \quad \forall y \neq 0$

$$\implies \underbrace{(Xy)^\top}_{\neq 0 \text{ as full rank}} D(Xy) > 0 \implies X^\top DX \text{ is pd} \implies X^\top DX \text{ invertible because}$$

for positive definite (square) matrices  $M$  it holds

$$\det(M) = \det(U\Lambda U^\top) = \det(U)\det(\Lambda)\det(U^\top) = \underbrace{\det(UU^\top)}_{=1} \underbrace{\det(\Lambda)}_{>0} > 0$$

From the four available datasets we selected a9a, ijcnn1 and covtype and using the described criterion we proved, that the Hessian of  $L_1$  is invertible for ijcnn1 but not for a9a and covtype. Since for singular Hessians Newton's method fails due to the necessary inversion of the Hessian during the update step it is necessary to pick an algorithm with a sufficient regularization for the respective datasets with  $L_1$ .

For the matrix  $\nabla^2 L_2(\omega) = X^\top D(\omega)X + \nabla^2 r(\omega)$  we will show, that for finite weights there exists a sufficient choice of  $\alpha$  s.t.  $D \succ 0$  (positive definite (pd)) holds.

*Proof:* Assume for now, that  $\nabla^2 r(\omega)$  is pd (we will show this at the end of this proof) for some sufficient choice of  $\alpha$  w.r.t. finite weights  $|w_j| < \infty$ .

We first observe, that  $X^\top DX$  is positive semi-definite (psd), i.e.  $y^\top X^\top D \underbrace{Xy}_{=\xi} = \xi^\top D \xi \geq 0$  because  $D$  is pd which implies that  $y^\top Dy > 0 \quad \forall y \neq 0$  and combining this with the observation, that  $\xi = Xy$  could be equal to 0 the inequality becomes sharp. This implies

$$y^\top (X^\top DX + \nabla^2 r(\omega))y = \underbrace{y^\top X^\top DXy}_{\geq 0} + \underbrace{y^\top \nabla^2 r(\omega)y}_{> 0} > 0 \implies \nabla^2 L_2(\omega) \text{ pd} \implies \text{invertible}$$

44 Inspecting the Hessian of the non convex regularizer  $\nabla^2 r(\omega) = \text{diag} \left( \lambda \frac{2\alpha(1-3\alpha\omega_j^2)}{(1+\alpha\omega_j^2)^3} \right)$  of the cross  
 45 entropy loss function  $L_2$  we directly notice, that since the matrix is psd for  $\lambda > 0, \alpha > 0$ , it is  
 46 invertible if and only if it is pd (because if it's not pd it means that one of its eigenvalues must  
 47 be 0). The matrix is diagonal and inspecting it's entries we directly see that pd holds if and only  
 48 if  $1 - 3\alpha \cdot w_j^2 > 0$  is satisfied, which is true for  $\alpha < \frac{1}{3w_j^2}$ . Thus we can always find a feasible  
 49 choice for  $\alpha > 0$  s.t.  $\nabla^2 r(\omega)$  is pd for finite weights. In our experiments we are presented with  
 50 two practical issues, that weaken this statement. First we were given what we understood to be a  
 51 mandatory parameter choice of  $\alpha = 1$  in the project description. The more interesting observation  
 52 however, is that even when allowed to choose the regularization parameter  $\alpha > 0$  freely, the machine  
 53 precision will treat any weight entries above the machine precision number as infinite and thus  
 54 even though  $D$  is analytically pd we have that numerically the matrix degenerates for large weights.  
 55 Numerically this can be stabilized by bounding weights heuristically, but since we focused on

comparing the performance of different Newton-type methods for practical problems we refrained from doing so to not bias the results. Therefore we cannot guarantee, that the Hessian of our second loss function  $L_2$  is invertible. In the experiments we will see that in fact singular Hessians appear for this Loss functions, making it intractable to solve with non-regularized Newton-type methods.

## 2. Convexity:

Since  $\log(\hat{y}_i), \log(1 - \hat{y}_i)$  are concave on  $(0, \infty)$  it follows that  $-\log(\hat{y}_i), -\log(1 - \hat{y}_i)$  are convex and thus  $L_1$  is a linear combination of convex functions (which is again convex). Meanwhile  $L_2$  is not guaranteed to be convex due to the non-convex regularization term  $r(\omega)$ .

3. Hessian Lipschitz Both Hessians are Lipschitz, as  $\frac{1}{n}X^\top DX =: M$  satisfies

$$\|M(\omega_1) - M(\omega_2)\| = \frac{1}{n}\|X^\top(D(\omega_1) - D(\omega_2))X\| \leq \underbrace{\frac{1}{n}\|X^\top\|\|X\|}_{=:C}\|D(\omega_1) - D(\omega_2)\|$$

Now consider that  $d(\sigma) := \sigma(z_k)(1 - \sigma(z_k))$  where  $z_k = y_i x_i^\top \omega_k$  and observe, that  $\frac{d}{d\sigma} = 1 - 2\sigma$  for  $\sigma \in (0, 1)$  then it follows, that  $\sigma'(z) = \sigma(z)(1 - \sigma(z))$  (by Remark 2) and by mean value theorem (MVT) we can conclude, that for

$$d'(z) = \sigma(z)(1 - \sigma(z))(1 - 2\sigma(z))$$

we have

$$|d(z_1) - d(z_2)| \leq \sup_z |d'(z)| \cdot |z_1 - z_2| \text{ where}$$

$$|z_1 - z_2| = |x_i^\top(\omega_1 - \omega_2)| \leq \underbrace{|y_i|}_{\leq 1} \|x_i\| \|\omega_1 - \omega_2\|$$

$$\implies |D_{ii}(\omega_1) - D_{ii}(\omega_2)| \leq \sup_z |d'(z)| \cdot \|x_i\| \|\omega_1 - \omega_2\|$$

$$\|D(\omega_1) - D(\omega_2)\| \leq \sup_z |d'(z)| \cdot \max_i \|x_i\| \|\omega_1 - \omega_2\|$$

$$\sup_z |d'(z)| = \frac{1}{4}$$

$$\|D(\omega_1) - D(\omega_2)\| \leq \underbrace{\frac{1}{4} \max_i \|x_i\|}_{=:L'} \|\omega_1 - \omega_2\|$$

$$\implies \|M(\omega_1) - M(\omega_2)\| \leq C \|D(\omega_1) - D(\omega_2)\| \leq \underbrace{CL'}_{=:L} \|\omega_1 - \omega_2\|$$

## 4. Positive definiteness of Hessians

Since  $\nabla^2 L_1(\omega)$  is psd we know, that the Hessian is pd for a dataset, if and only if the Hessian is invertible (because psd Hessians with nonzero eigenvalues only have positive eigenvalues and thus they're pd). It follows that the Hessian of `ijcnn1` is pd while the Hessians of `a9a` and `covtype` are not pd. Since the Hessian of the regularization term  $\nabla^2 r(\omega)$  potentially has negative diagonal entries  $\nabla^2 L_2(\omega)$  is not guaranteed to be pd. [TODO: I could analyze when exactly that happens, but it feels pretty useless tbh, as the point was made I wanted to make i.e. that we cannot assume pd for AICN].

## 2.2 Classic Newton's Method

The classical origin of Newton's method is as an algorithm for finding the roots of functions. In this paper it is used to find the roots  $x^*$  of  $\nabla(f(x))$  s.t.  $\nabla(f(x^*)) = 0$  and  $x^*$  a local minimum of  $f$ . Newton's method combined with a stepsize  $\eta$  uses the update rule [1]:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k) \quad (11)$$

The inverse Hessian can be interpreted as transforming the gradient landscape to be more isotropic, thereby improving the conditioning of the problem.

## 87 2.3 Cubic Newton

88 AICN gibt sich zwar als regularized method aus, kann in Wirklichkeit aber nicht umgehen die Matrix  
 89 trotzdem zur Berechnung des Faktors Alpha invertieren zu müssen. Es kämpft deshalb für singulare  
 90 oder illcondiitoned matrizen mit genau denselben problemen, wie unregularisierte Methoden. Kann  
 91 man das sicher nicht umgehen, dass man für das Alpha das Skalarprodukt invertieren muss

## 92 2.4 Cubic Newton

93 The cubic Newton method was one of the first to achieve a good complexity guarantee globally  
 94 [REFERENCE TO DO: What convergence rate exactly?]. It is based on cubic regularization and uses  
 95 the update rule:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k) + H \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \mathbf{I})^{-1} \nabla f(\mathbf{x}_k) \quad (12)$$

## 96 2.5 Levenberg and Marquardt method

97 The Levenberg-Marquardt's algorithm [REFERENCE] is an early form of regularized Newton's  
 98 method that modifies the Hessian. For ill conditioned (or singular) H regularization can increase  
 99 the conergence (or make the problem solvable as  $H + \lambda I$  is always invertible for sufficiently large  
 100  $\text{eig}(H) > -\lambda, \lambda > 0$ ). The update rule is:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k) + \lambda_k \mathbf{I})^{-1} \nabla f(\mathbf{x}_k) \quad (13)$$

## 101 2.6 Regularized Newton

102 In their 2023 article Michenko presents a variation of Newton's method that uses the update rule [2]:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k) + \sqrt{H \|\nabla f(\mathbf{x}_k)\|} \mathbf{I})^{-1} \nabla f(\mathbf{x}_k) \quad (14)$$

103 where  $H > 0$  is a constant. The convergence rate of this algorithm is  $\mathcal{O}(\frac{1}{k^2})$ . This method uses an  
 104 adaptive variant of the Levenberg-Marquardt regularization.

## 105 2.7 Appendix

106 Remark 2:

$$\begin{aligned} \sigma(z) &= \frac{1}{1 + e^{-z}} \\ \implies \frac{d}{dz} \sigma(z) &= \frac{d}{dz} (1 + e^{-z})^{-1} = -(1 + e^{-z})^{-2} \cdot (-e^{-z}) = \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} = \sigma(z)(1 - \sigma(z)) \end{aligned}$$

107

$$\begin{aligned} L_1(\omega) &= -\frac{1}{n} \sum_{i=1}^n \left[ \underbrace{y_i \log \hat{y}_i}_{=: A_i} + \underbrace{(1 - y_i) \log(1 - \hat{y}_i)}_{=: B_i} \right] \\ \hat{y}_i &= \sigma(x_i^\top \omega) = \frac{1}{1 + e^{-x_i^\top \omega}} \end{aligned}$$

108 and applying Remark 2 to  $\hat{y}$  we get, that

$$\begin{aligned}
\frac{\partial}{\partial \omega} A_i &= \frac{\partial}{\partial \omega} (-y_i \log \hat{y}_i) = -y_i \frac{1}{\hat{y}_i} (\hat{y}_i (1 - \hat{y}_i)) x_i = -y_i (1 - \hat{y}_i) x_i \\
\frac{\partial}{\partial \omega} B_i &= \frac{\partial}{\partial \omega} (-(1 - y_i) \log(1 - \hat{y}_i)) = (1 - y_i) \frac{1}{1 - \hat{y}_i} \hat{y}_i (1 - \hat{y}_i) x_i = (1 - y_i) \hat{y}_i x_i \\
\frac{\partial}{\partial \omega} A + \frac{\partial}{\partial \omega} B &= -y_i (1 - \hat{y}_i) x_i + (1 - y_i) \hat{y}_i x_i = (-y_i + y_i \hat{y}_i + \hat{y}_i - y_i \hat{y}_i) x_i \\
&= (-y_i + \hat{y}_i) x_i = (\hat{y}_i - y_i) x_i \\
\Rightarrow \nabla L_1(\omega) &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \omega} A_i + \frac{\partial}{\partial \omega} B_i = \frac{1}{n} \sum_{i=1}^n [\hat{y}_i - y_i] x_i = \frac{1}{n} X^\top (\hat{y} - y)
\end{aligned}$$

109 For the Hessian it then follows

$$\begin{aligned}
\nabla_\omega^2 L_1(\omega) &= \nabla_\omega \frac{1}{n} X^\top (\hat{y} - y) = \frac{1}{n} X^\top = \nabla_\omega (\hat{y} - y) = \frac{1}{n} X^\top \nabla_\omega \hat{y} \\
\frac{\partial}{\partial \omega} (\hat{y}_i x_i) &= \hat{y}_i (1 - \hat{y}_i) x_i x_i^\top \\
\Rightarrow \frac{\partial \hat{y}}{\partial \omega} &= \text{diag}(\sigma(X\omega) \odot (1 - \sigma(X\omega))) X \\
\Rightarrow \nabla^2 L_1(\omega) &= \frac{1}{n} X^\top \text{diag}(\hat{y} \odot (1 - \hat{y})) X \\
\Rightarrow \nabla^2 L_1(\omega) &= \frac{1}{n} X^\top D X \\
D &= \text{diag}(\hat{y}_i (1 - \hat{y}_i)).
\end{aligned}$$

110 For  $L_2$  we have

$$L_2(\omega) = \frac{1}{n} \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i x_i^\top \omega))}_{f_i(\omega)} + \lambda \underbrace{\sum_{j=1}^d \frac{\alpha \omega_j^2}{1 + \alpha \omega_j^2}}_{r(\omega)}$$

111 For the gradient we then get

$$\begin{aligned}
\frac{\partial}{\partial \omega_j} r(\omega) &= 2\lambda \alpha \frac{\omega_j}{(1 + \alpha \omega_j^2)^2} \Rightarrow \nabla r(\omega) = 2\lambda \alpha \frac{\omega}{(1 + \alpha \omega^2)^2} \\
\nabla f_i(\omega) &= \frac{\partial}{\partial \omega} \log(1 + e^{-y_i x_i^\top \omega}) \\
&= \frac{1}{\underbrace{1 + e^{y_i x_i^\top \omega}}_{\sigma(-y_i x_i^\top \omega)}} \cdot (-y_i x_i) = \sigma(-y_i x_i^\top \omega) \cdot (-y_i x_i) = -y_i x_i \sigma(-y_i x_i^\top \omega) \\
\nabla f(\omega) &= -\frac{1}{n} \sum_{i=1}^n y_i x_i \sigma(-y_i x_i^\top \omega) = -\frac{1}{n} X^\top (y \odot \sigma(-y \odot (X\omega))) \\
\nabla L_2(\omega) &= \nabla f(\omega) + \nabla r(\omega) \\
&= -\frac{1}{n} X^\top (y \odot \sigma(-y \odot (X\omega))) + 2\lambda \alpha \frac{\omega}{(1 + \alpha \omega^2)^2}
\end{aligned}$$

112 For the Hessians we first observe two remarks:

113 Remark 3: By chain rule we have

$$\begin{aligned}
z_i(\omega) &:= -y_i x_i^\top \omega \\
\Rightarrow \nabla_\omega z_i(\omega) &= -y_i x_i \\
\Rightarrow \nabla_\omega \sigma(z_i(\omega)) &= \sigma'(z_i(\omega)) \nabla_\omega z_i(\omega) \\
&= \sigma(-y_i x_i^\top \omega) (1 - \sigma(-y_i x_i^\top \omega)) (-y_i x_i)
\end{aligned}$$

114 From the gradient we have

$$\nabla_{\omega}^2 f(\omega) = \nabla_{\omega} \left( -\frac{1}{n} X^{\top} (y \odot \sigma(-y \odot (X\omega))) \right) = -\frac{1}{n} X^{\top} \nabla_{\omega} (y \odot \sigma(-y \odot (X\omega)))$$

115 Now notice, that

$$y \odot \sigma(-y \odot (X\omega)) = \begin{pmatrix} y_1 \sigma(-y_1 x_1^{\top} \omega) \\ y_2 \sigma(-y_2 x_2^{\top} \omega) \\ \vdots \\ y_n \sigma(-y_n x_n^{\top} \omega) \end{pmatrix}$$

116 and applying Remark 3 yields

$$\begin{aligned} \nabla_{\omega} \sigma(-y_i x_i^{\top} \omega) &= \sigma(-y_i x_i^{\top} \omega) (1 - \sigma(-y_i x_i^{\top} \omega)) (-y_i x_i) \\ \implies \nabla_{\omega} (y_i \sigma(-y_i x_i^{\top} \omega)) &= - \underbrace{y_i^2}_{=1 \text{ by Remark 1}} \sigma(-y_i x_i^{\top} \omega) (1 - \sigma(-y_i x_i^{\top} \omega)) x_i = -\sigma(-y_i x_i^{\top} \omega) (1 - \sigma(-y_i x_i^{\top} \omega)) x_i \end{aligned}$$

117

$$\begin{aligned} \implies \nabla_{\omega} (y \odot \sigma(-y \odot (X\omega))) &= - \begin{pmatrix} \overbrace{\sigma(-y_1 x_1^{\top} \omega) (1 - \sigma(-y_1 x_1^{\top} \omega))}^{=D_{1,1}} x_1 \\ \vdots \\ \underbrace{\sigma(-y_n x_n^{\top} \omega) (1 - \sigma(-y_n x_n^{\top} \omega))}_{D_{n,n}} x_n \end{pmatrix} \\ &= - \begin{bmatrix} D_{1,1} & 0 & \cdots & 0 \\ 0 & D_{2,2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & D_{n,n} \end{bmatrix} \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,d} \end{bmatrix} \\ &= - \begin{bmatrix} D_{1,1} x_{1,1} & D_{1,1} x_{1,2} & \cdots & D_{1,1} x_{1,d} \\ D_{2,2} x_{2,1} & D_{2,2} x_{2,2} & \cdots & D_{2,2} x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ D_{n,n} x_{n,1} & D_{n,n} x_{n,2} & \cdots & D_{n,n} x_{n,d} \end{bmatrix} = - \begin{bmatrix} D_{1,1} x_1^{\top} \\ D_{2,2} x_2^{\top} \\ \vdots \\ D_{n,n} x_n^{\top} \end{bmatrix} = -DX \end{aligned}$$

118 where we factored out the  $x_i$  in the last step to rewrite it as matrix-vector product. Deriving the entire  
119 expression we conclude:

$$\begin{aligned} \nabla^2 f(\omega) &= -\frac{1}{n} X^{\top} \nabla_{\omega} (y \odot \sigma(-y \odot (X\omega))) = \frac{1}{n} X^{\top} DX \\ D_{ii} &= \sigma(-y_i x_i^{\top} \omega) (1 - \sigma(-y_i x_i^{\top} \omega)) \end{aligned}$$

120 The hessian of the non-convex regularization term is derived by

$$\begin{aligned} \nabla_{\omega}^2 r(\omega) &= \nabla_{\omega} \left( 2\lambda \alpha \frac{\omega_j}{(1 + \alpha \omega_j^2)^2} \right) \\ \frac{\partial^2}{\partial \omega_j^2} r(\omega) &= 2\lambda \alpha \frac{\partial}{\partial \omega_j} \left( \frac{\omega_j}{(1 + \alpha \omega_j^2)^2} \right) = 2\lambda \alpha \frac{(1 + \alpha \omega_j^2)^2 - 4\alpha \omega_j^2 (1 + \alpha \omega_j^2)}{(1 + \alpha \omega_j^2)^4} = 2\lambda \alpha \frac{1 - 3\alpha \omega_j^2}{(1 + \alpha \omega_j^2)^3} \\ \implies \nabla^2 r(\omega) &= \text{diag} \left( 2\lambda \alpha \frac{1 - 3\alpha \omega_j^2}{(1 + \alpha \omega_j^2)^3} \right)_{j=1, \dots, d} \end{aligned}$$

121 Combining the steps we derive the Hessian

$$\begin{aligned} \nabla^2 L_2(\omega) &= \nabla^2 f(\omega) + \nabla^2 r(\omega) = \frac{1}{n} X^{\top} DX + \text{diag} \left( 2\lambda \alpha \frac{1 - 3\alpha \omega_j^2}{(1 + \alpha \omega_j^2)^3} \right) \\ D_{ii} &= \sigma(-y_i x_i^{\top} \omega) (1 - \sigma(-y_i x_i^{\top} \omega)) \end{aligned}$$

## 122 2.8 Inexact Newton Method

Given that Newton has cubic complexity we now outline how we aim to reduce the runtime by extending CG and MINRES methods to the Newton-type methods described in our paper. In order for the modified algorithms to inherit the convergence guarantees of the algorithms we want to approximate  $p$  s.t.

$$\|Hp + \nabla f\| \leq \epsilon \text{ (absolute tolerance)} < \epsilon = 10^{-8}$$

Since  $H_{1,2} = \nabla^2 L_{1,2}$  are clearly symmetric (as both  $X^\top DX$  and  $\nabla^2 r(x)$  are) we can apply the conjugate gradient method if the  $H$  is positive definite or have to fall back on MINRES if it is not pd. Positive definiteness depends on the data matrix and the regularizer curvature. [TODO: runtime for MINRES and CG]

Every iteration of Vanilla Newton takes  $O(n^3)$  per iteration because inversion of the Hessian costs  $O(n^3)$ .

for symmetric applying CG to newton drops the effort for inversion down to

$$O(k \cdot n^2) = O(\sqrt{\kappa} \log(1/\epsilon) \cdot n^2)$$

123 where  $\kappa(H) = \frac{\lambda_{\max}(H)}{\lambda_{\min}(H)}$

124 Precondition with SSOR to reduce condition number.

## 125 References

- 126 [1] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.
- 127 [2] Konstantin Mishchenko. Regularized newton method with global convergence. *SIAM Journal on*  
128 *Optimization*, 33(3):1440–1462, 2023.
- 129 [3] Slavomír Hanzely, Dmitry Kamzolov, Dmitry Pasechnyuk, Alexander Gasnikov, Peter Richtárik,  
130 and Martin Takáč. A damped newton method achieves global  $(o)(\frac{1}{k^2})$  and local quadratic  
131 convergence rate. *Advances in Neural Information Processing Systems*, 35:25320–25334, 2022.
- 132 [4] Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, Cambridge, MA,  
133 2012.

## 134 Checklist

135 The checklist follows the references. Please read the checklist guidelines carefully for information on  
136 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
137 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
138 the appropriate section of your paper or providing a brief inline description. For example:

- 139 • Did you include the license to the code and datasets? **[Yes]** See Section
- 140 • Did you include the license to the code and datasets? **[No]** The code and the data are  
141 proprietary.
- 142 • Did you include the license to the code and datasets? **[N/A]**

143 Please do not modify the questions and only use the provided macros for your answers. Note that the  
144 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
145 block and only keep the Checklist section heading above along with the questions/answers below.

146 1. For all authors...

- 147 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
148 contributions and scope? **[TODO]**
- 149 (b) Did you describe the limitations of your work? **[TODO]**
- 150 (c) Did you discuss any potential negative societal impacts of your work? **[TODO]**
- 151 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
152 them? **[TODO]**



- 153 2. If you are including theoretical results...
- 154 (a) Did you state the full set of assumptions of all theoretical results? **[TODO]**
- 155 (b) Did you include complete proofs of all theoretical results? **[TODO]**
- 156 3. If you ran experiments...
- 157 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 158 mental results (either in the supplemental material or as a URL)? **[TODO]**
- 159 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 160 were chosen)? **[TODO]**
- 161 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 162 ments multiple times)? **[TODO]**
- 163 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 164 of GPUs, internal cluster, or cloud provider)? **[TODO]**
- 165 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 166 (a) If your work uses existing assets, did you cite the creators? **[TODO]**
- 167 (b) Did you mention the license of the assets? **[TODO]**
- 168 (c) Did you include any new assets either in the supplemental material or as a URL?
- 169 **[TODO]**
- 170 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 171 using/curating? **[TODO]**
- 172 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 173 information or offensive content? **[TODO]**
- 174 5. If you used crowdsourcing or conducted research with human subjects...
- 175 (a) Did you include the full text of instructions given to participants and screenshots, if
- 176 applicable? **[TODO]**
- 177 (b) Did you describe any potential participant risks, with links to Institutional Review
- 178 Board (IRB) approvals, if applicable? **[TODO]**
- 179 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 180 spent on participant compensation? **[TODO]**

## 181 **A Appendix**

182 Optionally include extra information (complete proofs, additional experiments and plots) in the

183 appendix. This section will often be part of the supplemental material.