

Algorithms

Jimmy Zhan

August 2025

Idea 1

One of the ideas I came up with is to replace the local quadratic Taylor approximation with a truncated Fourier expansion of the per-sample loss.

Logistic-regression negative log-likelihood has the form

$$F(\omega) = \sum_{i=1}^n l(z_i), \quad z_i = \alpha + \beta^\top x_i, \quad (1)$$

where

$$\omega = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$

My idea is to approximate $l(z)$ by a truncated Fourier series on a compactified interval and then use that approximation to compute gradient/Hessian and a Newton-like update.

Compactification

Fourier series are for functions on a compact interval (often $[-\pi, \pi]$). $l(z)$ (per-sample negative log-likelihood) is smooth but defined on \mathbb{R} .

Two practical options:

A. Truncate domain: pick an interval $[L, U]$ that contains the likely range of $z = \alpha + \beta^\top x$ during optimization (e.g. use current iterate to estimate).

Rescale z linearly to $\theta \in [-\pi, \pi]$:

$$\theta(z) = \pi \cdot \frac{2(z - L)}{U - L} - \pi \quad (2)$$

B. Smooth compactifier (preferred if z can be large): use a smooth monotone map $\theta(z) = \arctan(\lambda(z - z_0))$ scaled into $[-\pi, \pi]$. This avoids abrupt periodic

extension and reduces Gibbs phenomena.

Denote $\theta = \theta(z)$. We approximate $l(z)$ by

$$l_F(z) = c_0 + \sum_{m=1}^M (a_m \cos(m\theta(z)) + b_m \sin(m\theta(z))) \quad (3)$$

Coefficients a_m, b_m, c_m are computed once (numerically) by sampling $l \circ \theta^{-1}$ on a dense grid of θ and doing a real FFT / discrete Fourier transform. (If you use the arctan compactifier include Jacobian when transforming integrals; for coefficient estimation you only need sampled values.)

Closed-form

The Fourier approximation gives closed-form expressions for the derivative and second derivative of l_F with respect to z . Use chain rule via $\theta = \theta(z)$

Define

$$S_1(\theta) = \sum_{m=1}^M m(-a_m \sin(m\theta) + b_m \cos(m\theta)) \quad (4)$$

and

$$S_2(\theta) = - \sum_{m=1}^M m^2 (a_m \cos(m\theta) + b_m \sin(m\theta)) \quad (5)$$

Then

$$\begin{aligned} l'_F(z) &= \theta'(z) S_1(\theta(z)), \\ l''_F(z) &= \theta''(z) S_1(\theta(z)) + [\theta'(z)]^2 S_2(\theta(z)) \end{aligned} \quad (6)$$

Now for the full parameter vector $\omega = [\alpha; \beta]$, $z_i = \omega^\top \tilde{x}_i$ where $\tilde{x}_i = [1; x_i]$. So the approximated objective

$$F_F(\omega) = \sum_{i=1}^n l_F(z_i) \quad (7)$$

has

$$\nabla F_F(\omega) = \sum_{i=1}^n l'_F(z_i) \tilde{x}_i \quad (8)$$

and the Hessian

$$H_F(\omega) = \sum_{i=1}^n l_F''(z_i) \tilde{x}_i \tilde{x}_i^\top \quad (9)$$

Newton-Fourier update:

$$\omega_{k+1} = \omega_k - (H_F(\omega_k) + \lambda I)^{-1} \nabla F_F(\omega_k) \quad (10)$$

optionally adding damping $\lambda > 0$ (Levenberg-Marquardt style) for stability.

Practical Implementation and Parameter Selection

The performance of the Newton-Fourier method depends on several key parameters.

1. Compactification Parameters $[L, U]$ or z_0, λ :

These parameters should be chosen to capture the "active" region of the loss function. A robust approach is to perform a pre-analysis on the dataset. For instance, one could train a simple model (e.g., a few steps of gradient descent) to get a rough estimate of the distribution of z_i values and set the interval $[L, U]$ to cover, for example, three standard deviations around the mean. For the arctan mapping, z_0 can be set to the mean of z_i and λ can be chosen such that the interval of interest is mapped to a significant portion of $[\pi, \pi]$. Importantly, these parameters should be fixed before optimization begins to avoid the costly re-computation of Fourier coefficients at each iteration.

2. Number of Fourier Terms M:

The choice of M represents a trade-off between approximation accuracy and computational cost.

Small M: Leads to a coarse approximation but faster computations per iteration.

Large M: Provides a more accurate approximation of the loss function but increases the cost of evaluating the gradient and Hessian.

Selection Strategy: One can analyze the decay of the computed Fourier coefficients $|a_m|$ and $|b_m|$. Typically, these decay rapidly for smooth functions.

M can be chosen such that coefficients beyond this point are negligible (e.g., below a certain machine-epsilon tolerance). Alternatively, M can be treated as a hyperparameter and tuned using a validation set

Computational Cost Analysis

The main overhead of this method lies in the pre-computation of Fourier coefficients and the per-iteration cost of evaluating the series.

For Preprocessing cost, the Fourier coefficients are computed once by sampling the loss function at K points and running an FFT. The complexity is $O(K \log K)$, where K is typically chosen to be much larger than M. This is a one-time, fixed cost.

Pre-Iteration Cost

Let n be the number of samples and d be the number of features.

For each of the n samples, we evaluate series S_1 with M terms, which is an $O(M)$ operation. This is scaled and summed up for d features, therefore $O(n \cdot d \cdot M)$.

For each sample, we evaluate S_1 and $S_2(O(M))$ and then form a $d \times d$ outer product, costing $O(d^2)$. Therefore $O(n \cdot d^2 \cdot M)$ in total.

Matrix Inversion: $O(d^3)$ for solving the linear system.

The Newton-Fourier method is computationally more expensive per iteration than the standard Newton method by a factor of M. Its potential advantage must therefore come from a significant reduction in the total number of iterations required for convergence.

Core Advantage: Global Approximation

The fundamental distinction of this method is its use of a global approximation.

A second-order Taylor expansion is a purely local approximation. It models the function as a parabola that is accurate only in an infinitesimal neighborhood of the current point ω_k . If the function's true shape deviates significantly from this parabola (which is common far from the optimum), the resulting Newton

step can be too large, too small, or even point in a wrong direction, requiring line searches or trust regions to stabilize.

A truncated Fourier series provides a global approximation across the entire compactified interval. By capturing the overall shape, including non-quadratic features of the loss function, the resulting Hessian H_F can provide a more accurate model of the function's curvature over a much wider range. This could lead to significantly better and more direct update steps, especially during the early stages of optimization when the iterates are far from the solution. This improved global perspective has the potential to drastically reduce the number of iterations needed to reach convergence.

Idea 1 Ver 2.0

The Newton-Fourier method has the following iteration formula

$$\omega_{k+1} = \omega_k - (H_F(\omega_k) + \lambda I)^{-1} \nabla F_F(\omega_k) \quad (11)$$