

A. Additional experimental details

For NS-ada, we double the sketching dimension when $\tilde{\lambda}_f(x^{t+1}) > c_1 \tilde{\lambda}_f(x^t) \min(1, c_2 \tilde{\lambda}_f(x^t)^\tau)$. Here $c_1, c_2 > 0$ and $\tau \in [0, 1]$. For all compared methods, we use the backtracking line search method to find a step size satisfying the Armijo condition. For NS-ada, NS and NE, we stop the algorithm when $\tilde{\lambda}_f(x) < 10^{-6}$ or $\lambda_f(x) < 10^{-6}$. For first-order methods, we first compute a referenced solution \tilde{x}^* based on NS-ada. Then, we stop the algorithm when $\frac{f(x) - f(\tilde{x}^*)}{1 + f(\tilde{x}^*)} < 10^{-6}$.

The parameters of NS-ada and NS with for each dataset are summarized in Tables 2 to 4.

Dataset	m_0	c_1	τ	c_2
RCV1	100	2	0	1
MNIST	100	0.5	1	6
gisette	10	2	0	1
realsim	100	2	0	1
epsilon	100	1	0	1
a8a-kernel	10	0.5	0	1
w7a-kernel	10	0.5	0	1

Table 2. Parameters of adaptive Newton sketch with SJLT sketching.

Dataset	m_0	c_1	τ	c_2
RCV1	100	1	0	1
MNIST	100	0.5	1	6
gisette	10	2	0	1
realsim	100	2	0	1
epsilon	100	1	0	1
a8a-kernel	10	0.5	0	1
w7a-kernel	100	0.5	0	1

Table 3. Parameters of adaptive Newton sketch with RRS sketching.

Dataset	m (SJLT)	m (RRS)
RCV1	800	800
MNIST	800	1600
gisette	400	400
realsim	800	3200
epsilon	800	3200
a8a-kernel	100	800
w7a-kernel	100	800

Table 4. Sketching dimensions of Newton Sketch.

We present numerical performance of compared methods with additional details and additional numerical experiments in Figures 5 to 9. Comparatively, NS-ada-RRS tends to have larger sketching dimension than NS-ada-SJLT. This may come from that NS-RRS has stronger oscillations than NS-SJLT in the plot of $\tilde{\lambda}_f(x^t)$. Thus, NS-ada-RRS can be slower than NS-ada-SJLT in some test cases where n is not significantly larger than d .

For kernelized regularized logistic regression, the data matrices A and \tilde{A} are constructed as kernel matrices based on the original data features. Namely, it follows

$$A_{i,j} = k(\tilde{a}_i, \tilde{a}_j), \quad A_{i,j}^{\text{test}} = k(\tilde{a}_i^{\text{test}}, \tilde{a}_j),$$

where $\{\tilde{a}_i\}_{i=1}^n$ and $\{\tilde{a}_j^{\text{test}}\}_{j=1}^{n_{\text{test}}}$ are original data features from the training set and test set respectively. Here $k(x, x') : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a positive kernel function. We use the isotropic Gaussian kernel function:

$$k(x, x') = (2\pi h)^{-d/2} \exp\left(-\frac{1}{2h} \|x - x'\|_2^2\right),$$

Adaptive Newton Sketch

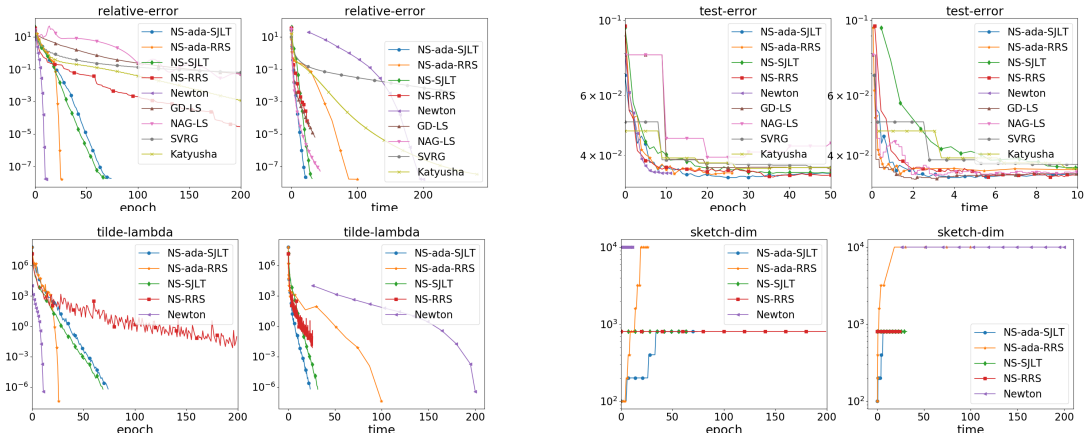


Figure 5. RCV1. $n = 10000, d = 47236, \mu = 10^{-3}$.

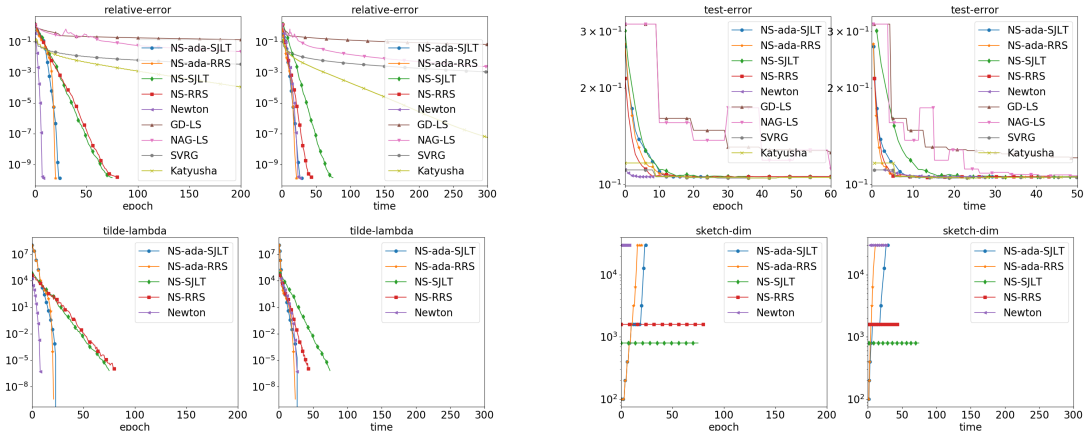


Figure 6. MNIST. $n = 30000, d = 780, \mu = 10^{-1}$.

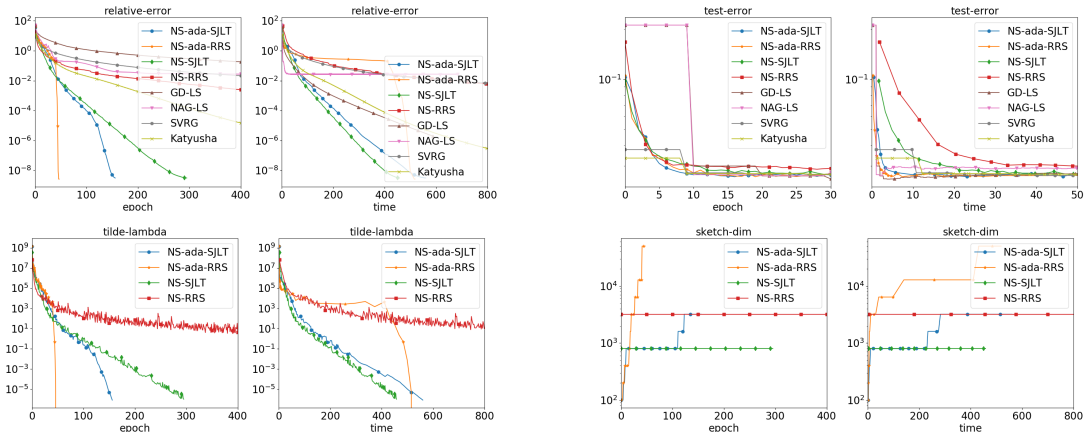


Figure 7. realsim. $n = 50000, d = 20958, \mu = 10^{-3}$.

Adaptive Newton Sketch

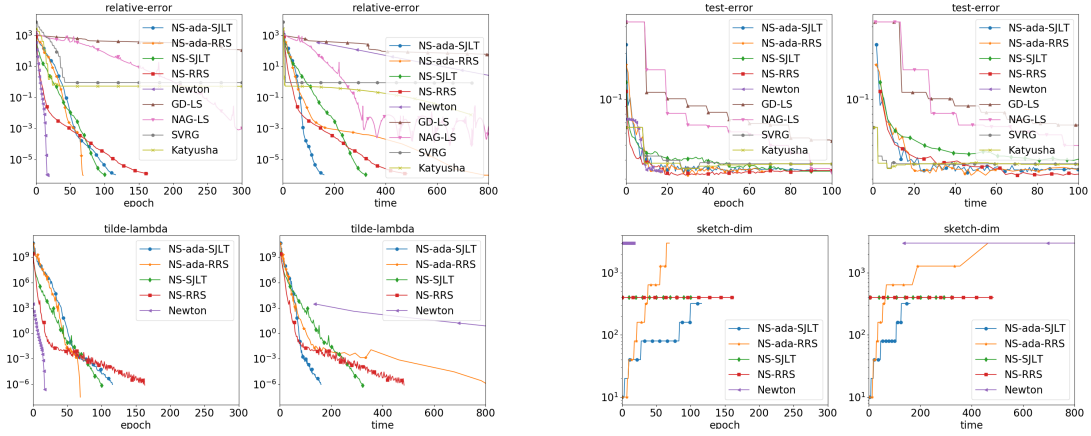


Figure 8. gisette. $n = 3000, d = 5000, \mu = 10^{-3}$.

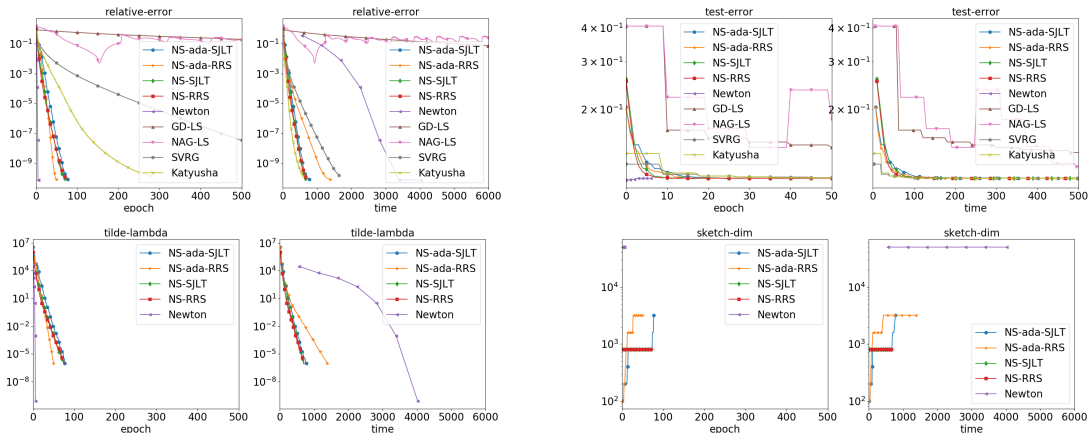


Figure 9. epsilon. $n = 50000, d = 2000, \mu = 10^{-1}$.

Adaptive Newton Sketch

where $h > 0$ is the bandwidth. We set $h = 10$ for a8a dataset and $h = 20$ for w7a dataset. For NS-ada-SJLT and NS-ada-RRS, we let $c_1 = 0.5, \tau = 0$ and $c_2 = 1$. For NS, the sketching dimensions are summarized in Table 5.

Dataset	m (SJLT)	m (RRS)
a8a-kernel	100	800
w7a-kernel	100	800

Table 5. Sketching dimensions of Newton Sketch. kernel matrix.

We present numerical results with additional details in Figures 10 and 11. We can also observe super linear convergence rate of NS-ada in the plot of $\lambda_f(x^t)$ when x^t is close to the optimum of the optimization problem. Similarly, NS-ada-RRS tends to have larger sketching dimension than NS-ada-SJLT.

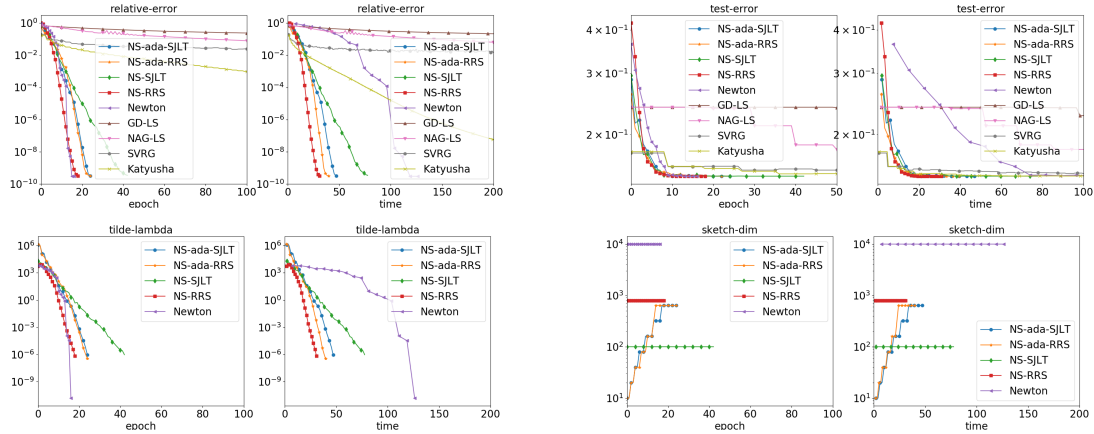


Figure 10. a8a. kernel matrix. $n = 10000, d = 10000, \mu = 10$.

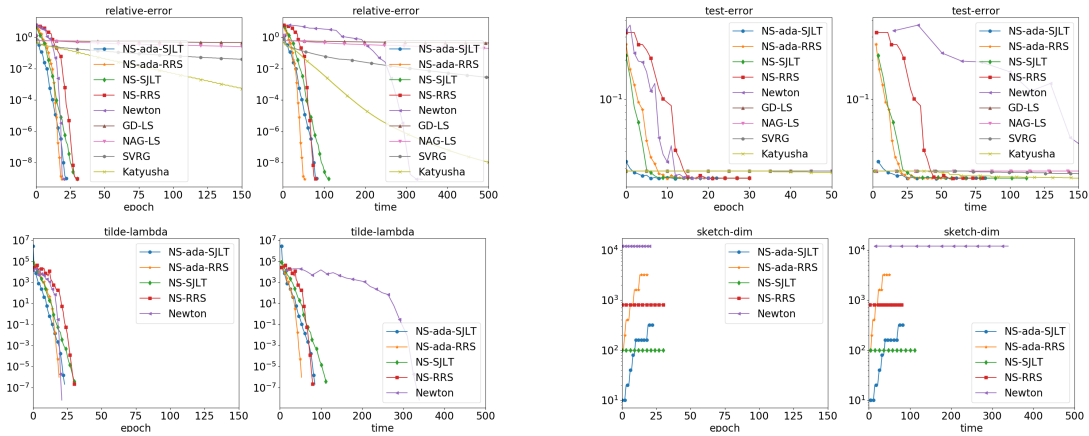


Figure 11. w7a. kernel matrix. $n = 12000, d = 12000, \mu = 10$.

B. Proof of main results

B.1. Proof of Lemma 1

Let $x \in \text{dom } f$. We use the shorthand $A := \nabla^2 f_0(x)^{1/2}$, and we let $A = U\Sigma V^\top$ be a thin SVD of A . We denote by $H^{1/2}$ an invertible square-root matrix of the Hessian $H \equiv H(x) = A^\top A + \nabla^2 g(x)$. Recall that $H_S \equiv H_S(x) =$

$A^\top S^\top SA + \nabla^2 g(x)$. Then, we have

$$\begin{aligned} C_S &= H^{-\frac{1}{2}} H_S H^{-\frac{1}{2}} = H^{-\frac{1}{2}} (H + (H_S - H)) H^{-\frac{1}{2}} \\ &= I_d + H^{-1/2} (H_S - H) H^{-1/2} \\ &= I_d + H^{-1/2} V \Sigma (U^\top S^\top S U - I_d) \Sigma V^\top H^{-1/2}. \end{aligned}$$

We use the shorthand $M := \Sigma V^\top H^{-1/2}$. Using the fact that $\nabla^2 g(x) \succeq \mu I_d$, it follows that

$$\|M\|_F^2 = \text{trace}(\Sigma V^\top H^{-1} V \Sigma) \leq \text{trace}(\Sigma V^\top (A^\top A + \mu I_d)^{-1} V \Sigma) = d_\mu(x). \quad (37)$$

It remains to control the spectral norm of $M^\top (U^\top S^\top S U - I_d) M$.

(SJLT). It was shown in (Nelson & Nguyen, 2013) that for $\varepsilon > 0$ and $p \in (0, 1/2)$, it holds with probability at least $1 - p$ that $\|M^\top (U^\top S^\top S U - I_d) M\|_2 \leq \varepsilon$ provided that $m \geq c_0 \frac{\|M\|_F^4}{\varepsilon^2 p}$, where $c_0 > 0$ is a universal constant. Note that this lower bound on the sketch size is increasing as a function of $\|M\|_F^2$. From inequality (37), it is then sufficient to have $m \geq c_0 \frac{d_\mu(x)^2}{\varepsilon^2 p}$ for the above inequality to hold with probability at least $1 - p$.

(SRHT). According to Theorems 1 and 9 in (Cohen et al., 2015), it holds with probability at least $1 - p$ that $\|M^\top (U^\top S^\top S U - I_d) M\|_2 \leq \varepsilon$ provided that $m \geq c_0 \varepsilon^{-2} \left(\|M\|_F^2 + \log(\frac{1}{\varepsilon p}) \log(\|M\|_F^2/p) \right)$, where c_0 is a universal constant. Note that this lower bound on the sketch size is increasing as a function of $\|M\|_F^2$. From inequality (37), it is then sufficient to have $m \geq c_0 \varepsilon^{-2} \left(d_\mu(x) + \log(\frac{1}{\varepsilon p}) \log(d_\mu(x)/p) \right)$ for the above inequality to hold with probability at least $1 - p$. \square

B.2. Proof of Theorem 1

Let $x \in \text{dom } f$. Plugging-in the definitions of v_{ne} and v_{nsk} , we have

$$\begin{aligned} \|v_{\text{ne}} - v_{\text{nsk}}\|_{H(x)} &= \|H^{1/2}(v_{\text{ne}} - v_{\text{nsk}})\|_2 = \|H^{1/2}(H_S^{-1} \nabla f(x) - H^{-1} \nabla f(x))\|_2 \\ &= \|(H^{1/2} H_S^{-1} H^{1/2} - I_d) H^{-1/2} \nabla f(x)\|_2 \\ &\leq \|C_S^{-1} - I_d\|_2 \|H^{-1/2} \nabla f(x)\|_2. \end{aligned}$$

Using that $\|H^{-1/2} \nabla f(x)\|_2 = \|v_{\text{ne}}\|_{H(x)}$, we further obtain

$$\|v_{\text{ne}} - v_{\text{nsk}}\|_{H(x)} \leq \|C_S^{-1} - I_d\|_2 \|v_{\text{ne}}\|_{H(x)}.$$

Under the event $\mathcal{E}_{x,m,\varepsilon}$, it holds for $\varepsilon \in (0, 1/4)$ that $(1 + \varepsilon/2)^{-1} I_d \preceq C_S^{-1} \preceq (1 - \varepsilon/2)^{-1} I_d$. Using the facts that $(1 + \varepsilon/2)^{-1} \geq 1 - \varepsilon$ and $(1 - \varepsilon/2)^{-1} \leq 1 + \varepsilon$, we obtain the inequality $\|C_S^{-1} - I_d\|_2 \leq \varepsilon$, whence

$$\|v_{\text{ne}} - v_{\text{nsk}}\|_{H(x)} \leq \varepsilon \|v_{\text{ne}}\|_{H(x)},$$

which proves the first inequality of Theorem 1. On the other hand, we have

$$\begin{aligned} \tilde{\lambda}_f(x)^2 &= \langle \nabla f(x), H_S^{-1} \nabla f(x) \rangle = \left\langle H^{-\frac{1}{2}} \nabla f(x), H^{\frac{1}{2}} H_S^{-1} H^{\frac{1}{2}} H^{-\frac{1}{2}} \nabla f(x) \right\rangle \\ &= \|C_S^{-\frac{1}{2}} H^{-\frac{1}{2}} \nabla f(x)\|_2^2. \end{aligned}$$

It follows that

$$\frac{1}{\sigma_{\max}(C_S)} \lambda_f(x)^2 \leq \tilde{\lambda}_f(x)^2 \leq \frac{1}{\sigma_{\min}(C_S)} \lambda_f(x)^2.$$

Conditional on the event $\mathcal{E}_{x,m,\varepsilon}$ and using that $(1 + \varepsilon/2)^{-1} \geq 1 - \varepsilon$ and $(1 - \varepsilon/2)^{-1} \leq 1 + \varepsilon$, we obtain the claimed result, i.e.,

$$(1 - \varepsilon) \lambda_f(x)^2 \leq \tilde{\lambda}_f(x)^2 \leq (1 + \varepsilon) \lambda_f(x)^2.$$

\square

B.3. Proof of Lemma 2

Our proof of this result closely follows the steps of the proof of Lemma 3(a) in (Pilanci & Wainwright, 2017): the core arguments are the same, but we adapt the proof to our technical framework, that is, conditional on the event $\mathcal{E}_{x,m,\varepsilon}$.

The strategy of the proof is to show that the backtracking line search leads to a step size $s > 0$ such that $f(x_{\text{nsk}}) - f(x) \leq -\nu$. We define the univariate function $g(u) := f(x + uv_{\text{nsk}})$ and we set $\varepsilon' = \frac{2\varepsilon}{1-\varepsilon}$. We first show that $\hat{u} = \frac{1}{1+(1+\varepsilon')\tilde{\lambda}_f(x)}$ satisfies the bound

$$g(\hat{u}) \leq g(0) - a\hat{u}\tilde{\lambda}_f^2(x), \quad (38)$$

which implies that \hat{u} satisfies the exit condition of backtracking line search. Therefore, the step size s must be lower bounded as $s \geq b\hat{u}$, which further implies that the new iterate $x_{\text{nsk}} = x + sv_{\text{nsk}}$ satisfies the decrement bound

$$f(x_{\text{nsk}}) - f(x) \leq -ab \frac{\tilde{\lambda}_f(x)^2}{1 + (1 + \frac{2\varepsilon}{1-\varepsilon})\tilde{\lambda}_f(x)}.$$

By assumption, $\tilde{\lambda}_f(x) > \eta$. Using the fact that the function $u \mapsto \frac{u^2}{1+(1+\frac{2\varepsilon}{1-\varepsilon})u}$ is monotone increasing, we get that

$$f(x_{\text{nsk}}) - f(x) \leq -ab \frac{\eta^2}{1 + (1 + \frac{2\varepsilon}{1-\varepsilon})\eta} = \nu,$$

which is exactly the claimed result. It remains to prove the claims (38).

According to Lemma 4 in (Pilanci & Wainwright, 2017), we have for any $u \geq 0$ and $\gamma \geq 0$ that

$$g(u) \leq g(0) - u\tilde{\lambda}_f(x)^2 - \gamma - \log(1 - \gamma), \quad (39)$$

provided that $u\|v_{\text{nsk}}\|_{H(x)} \leq \gamma < 1$. By assumption, the event $\mathcal{E}_{x,m,\varepsilon}$ holds true. As a consequence of Theorem 1, we have that

$$\|v_{\text{nsk}}\|_{H(x)} \leq (1 + \varepsilon)\lambda_f(x) \leq \frac{1 + \varepsilon}{1 - \varepsilon} \tilde{\lambda}_f(x) = (1 + \varepsilon')\tilde{\lambda}_f(x).$$

It follows that $\hat{u}\|v_{\text{nsk}}\|_{H(x)} \leq \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x) < 1$. Plugging-in $u = \hat{u}$ and $\gamma = \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x)$ into (39), we obtain that

$$\begin{aligned} g(\hat{u}) &\leq g(0) - \hat{u}\tilde{\lambda}_f(x)^2 - \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x) - \log(1 - \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x)) \\ &= g(0) - \left\{ \hat{u}(1 + \varepsilon')^2\tilde{\lambda}_f(x)^2 + \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x) + \log(1 - \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x)) - \hat{u}((1 + \varepsilon')^2 - 1)\tilde{\lambda}_f(x)^2 \right\}. \end{aligned}$$

Using that $\hat{u}(1 + \varepsilon')^2\tilde{\lambda}_f(x)^2 + \hat{u}(1 + \varepsilon')\tilde{\lambda}_f(x) = (1 + \varepsilon')\tilde{\lambda}_f(x)$ and $\hat{u}((1 + \varepsilon')^2 - 1)\tilde{\lambda}_f(x)^2 = \frac{(\varepsilon'^2 + 2\varepsilon')\tilde{\lambda}_f(x)^2}{1 + (1 + \varepsilon')\tilde{\lambda}_f(x)}$, we find that

$$g(\hat{u}) \leq g(0) - (1 + \varepsilon')\tilde{\lambda}_f(x) + \log(1 + (1 + \varepsilon')\tilde{\lambda}_f(x)) + \frac{(\varepsilon'^2 + 2\varepsilon')\tilde{\lambda}_f(x)^2}{1 + (1 + \varepsilon')\tilde{\lambda}_f(x)}.$$

Applying the inequality $-z + \log(1 + z) \leq -\frac{1}{2}\frac{z^2}{(1+z)}$ with $z = (1 + \varepsilon')\tilde{\lambda}_f(x)$, we further obtain that

$$\begin{aligned} g(\hat{u}) &\leq g(0) - \frac{\frac{1}{2}(1 + \varepsilon')^2\tilde{\lambda}_f(x)^2}{1 + (1 + \varepsilon')\tilde{\lambda}_f(x)} + \frac{(\varepsilon'^2 + 2\varepsilon')\tilde{\lambda}_f(x)^2}{1 + (1 + \varepsilon')\tilde{\lambda}_f(x)} \\ &= g(0) - \left(\frac{1}{2} - \frac{\varepsilon'^2}{2} - \varepsilon' \right) \tilde{\lambda}_f(x)^2 \hat{u} \\ &\leq g(0) - a\tilde{\lambda}_f(x)^2 \hat{u}, \end{aligned}$$

where the final inequality follows by the assumption that $a \leq 1 - \frac{1}{2}\left(\frac{1+\varepsilon}{1-\varepsilon}\right)^2$, that is, $a \leq \frac{1}{2} - \frac{\varepsilon'^2}{2} - \varepsilon'$. This concludes the proof. \square

B.4. Proof of Lemma 3

We recall Theorem 4.1.6 of (Nesterov, 2003) (see, also, Exercise 9.17 in (Boyd & Vandenberghe, 2004)): it guarantees that for a step size $s > 0$ such that $|1 - s\|v_{\text{nsk}}\|_{H(x)}| < 1$, we have

$$(1 - s\|v_{\text{nsk}}\|_{H(x)})^2 H(x) \preceq H(x + sv_{\text{nsk}}) \preceq \frac{1}{(1 - s\|v_{\text{nsk}}\|_{H(x)})^2} H(x). \quad (40)$$

By assumption, the event $\mathcal{E}_{x,m,\varepsilon'}$ holds. As a consequence of Theorem 1, we have $\|v_{\text{nsk}}\|_{H(x)} \leq (1 + \varepsilon')\|v_{\text{ne}}\|_{H(x)}$. Plugging this bound into (40) and using that $\|v_{\text{ne}}\|_{H(x)} = \lambda_f(x)$, we obtain

$$(1 - s(1 + \varepsilon')\lambda_f(x))^2 H(x) \preceq H(x + sv_{\text{nsk}}) \preceq \frac{1}{(1 - s(1 + \varepsilon')\lambda_f(x))^2} H(x), \quad (41)$$

for $s > 0$ such that $s(1 + \varepsilon')\lambda_f(x) < 1$. Denote by s_{nsk} the step size obtained by backtracking line search. It satisfies $s_{\text{nsk}} = 1$. Then, it holds that

$$\begin{aligned} s_{\text{nsk}}(1 + \varepsilon')\lambda_f(x) &\leq (1 + \varepsilon')\lambda_f(x) \stackrel{(i)}{\leq} \frac{1 + \varepsilon'}{\sqrt{1 - \varepsilon'}} \tilde{\lambda}_f(x) \\ &\stackrel{(ii)}{\leq} \frac{1 + \varepsilon'}{\sqrt{1 - \varepsilon'}} \eta \\ &\stackrel{(iii)}{<} 1, \end{aligned}$$

where inequality (i) follows from the assumption that $\mathcal{E}_{x,m,\varepsilon'}$ holds and from Theorem 1; inequality (ii) follows from the assumption that $\tilde{\lambda}_f(x) \leq \eta$. Furthermore, we have $\varepsilon' \leq \varepsilon < 1/4$, as well as $\eta < 1/16$ (see Lemma 7) and this yields inequality (iii).

Using (41), we then obtain that

$$\begin{aligned} \lambda_f(x_{\text{nsk}}) &= \|H(x_{\text{nsk}})^{-1/2} \nabla f(x_{\text{nsk}})\|_2 \\ &\leq \frac{1}{(1 - (1 + \varepsilon')\lambda_f(x))} \|H(x)^{-1/2} \nabla f(x_{\text{nsk}})\|_2 \\ &= \frac{1}{(1 - (1 + \varepsilon')\lambda_f(x))} \left\| H(x)^{-1/2} \left(\nabla f(x) + \int_0^1 H(x + sv_{\text{nsk}}) v_{\text{nsk}} ds \right) \right\|_2 \\ &\leq \frac{1}{(1 - (1 + \varepsilon')\lambda_f(x))} (M_1 + M_2), \end{aligned}$$

where

$$\begin{aligned} M_1 &= \left\| H(x)^{-1/2} \left(\nabla f(x) + \int_0^1 H(x + sv_{\text{nsk}}) v_{\text{ne}} ds \right) \right\|_2, \\ M_2 &= \left\| H(x)^{-1/2} \int_0^1 H(x + sv_{\text{nsk}}) (v_{\text{nsk}} - v_{\text{ne}}) ds \right\|_2. \end{aligned}$$

It remains to bound the terms M_1 and M_2 . Regarding M_1 , we have after re-arranging and using inequality (41) that

$$\begin{aligned} M_1 &= \left\| \int_0^1 \left(H(x)^{-1/2} H(x + sv_{\text{nsk}}) H(x)^{-1/2} - I_d \right) ds H(x)^{1/2} v_{\text{ne}} \right\|_2 \\ &\leq \left| \int_0^1 \frac{1}{(1 - s(1 + \varepsilon')\lambda_f(x))^2} ds - 1 \right| \left\| H(x)^{1/2} v_{\text{ne}} \right\|_2 \\ &= \frac{(1 + \varepsilon')\lambda_f^2(x)}{1 - (1 + \varepsilon')\lambda_f(x)}. \end{aligned}$$

Regarding M_2 , we have

$$\begin{aligned}
 M_2 &= \left\| \int_0^1 H(x)^{-1/2} H(x + sv_{\text{nsk}}) H(x)^{-1/2} ds H(x)^{1/2} (v_{\text{nsk}} - v_{\text{ne}}) \right\|_2 \\
 &\leq \left\| \int_0^1 \frac{1}{(1 - s(1 + \varepsilon')\lambda_f(x))^2} ds H(x)^{1/2} (v_{\text{nsk}} - v_{\text{ne}}) \right\|_2 \\
 &= \frac{1}{1 - (1 + \varepsilon')\lambda_f(x)} \left\| H(x)^{1/2} (v_{\text{nsk}} - v_{\text{ne}}) \right\|_2 \\
 &\leq \frac{\varepsilon' \lambda_f(x)}{1 - (1 + \varepsilon')\lambda_f(x)},
 \end{aligned}$$

where the last inequality follows from the assumption that the event $\mathcal{E}_{x,m,\varepsilon'}$ holds and as a consequence of Theorem 1. Plugging these bounds on M_1 and M_2 , we obtain that

$$\lambda_f(x_{\text{nsk}}) \leq \frac{(1 + \varepsilon')\lambda_f(x)^2 + \varepsilon' \lambda_f(x)}{(1 - (1 + \varepsilon')\lambda_f(x))^2}. \quad (42)$$

Recall that $\varepsilon' \leq \varepsilon \lambda_f(x)^\tau$. Combining this inequality with (42), we obtain

$$\lambda_f(x_{\text{nsk}}) \leq \frac{(1 + \varepsilon \lambda_f(x)^\tau) \lambda_f(x)^2 + \varepsilon \lambda_f(x)^{1+\tau}}{(1 - (1 + \varepsilon \lambda_f(x)^\tau) \lambda_f(x))^2} = \underbrace{\frac{\lambda_f(x)^{1-\tau} + \varepsilon \lambda_f(x) + \varepsilon}{(1 - (1 + \varepsilon \lambda_f(x)^\tau) \lambda_f(x))^2}}_{:=\alpha(\tau,x)} \lambda_f(x)^{1+\tau}.$$

On the event $\mathcal{E}_{x,m,\varepsilon'}$, we have according to Theorem 1 that $(1 + \varepsilon)\lambda_f(x) \leq \frac{(1+\varepsilon)\tilde{\lambda}_f(x)}{\sqrt{1-\varepsilon}} \leq \frac{(1+\varepsilon)\eta}{\sqrt{1-\varepsilon}} \leq \frac{1}{16}$, where the last inequality follows from Lemma 7. Hence, the denominator of $\alpha(\tau, x)$ satisfies

$$1 - (1 + \varepsilon \lambda_f(x)^\tau) \lambda_f(x) \geq 1 - (1 + \varepsilon) \lambda_f(x) \geq \frac{15}{16},$$

while the numerator of $\alpha(\tau, x)$ satisfies

$$\lambda_f(x)^{1-\tau} + \varepsilon \lambda_f(x) + \varepsilon \leq \frac{1}{16^{1-\tau}} + \frac{1}{32} + \frac{1}{2}$$

Combining these bounds together, we obtain that

$$\alpha(\tau, x) \leq \frac{8 + 1/2 + 16^\tau}{15} \leq 0.57 + \frac{16^\tau}{15} = \alpha(\tau).$$

It is easy to verify that $\alpha(\tau)^{1/\tau} \leq 2$ for any $\tau \in (0, 1]$. Furthermore, for $\tau = 0$, we obtain that $\alpha(0) \approx 0.63333 \leq 0.64 = \frac{16}{25}$, and this concludes the proof. Note that a similar linear convergence rate was obtained for the Newton sketch provided that $m \gtrsim d$ (see Lemma 3 in (Pilanci & Wainwright, 2017)). \square

B.5. Proof of Lemma 4

By induction, we obtain for any $t \geq 0$ that $\alpha^{\frac{1}{\tau}} \beta_t \leq (\alpha^{\frac{1}{\tau}} \eta)^{(1+\tau)^t}$. To have $\beta_t \leq \sqrt{\delta}$, it suffices that $(\alpha^{\frac{1}{\tau}} \eta)^{(1+\tau)^t} \leq \alpha^{\frac{1}{\tau}} \sqrt{\delta}$. Taking the logarithm on both sides, this yields $(1+\tau)^t \log(\alpha^{\frac{1}{\tau}} \eta) \leq \log(\alpha^{\frac{1}{\tau}} \sqrt{\delta})$, i.e., $(1+\tau)^t \log(1/\alpha^{\frac{1}{\tau}} \eta) \geq \log(1/\alpha^{\frac{1}{\tau}} \sqrt{\delta})$. By assumption, $\log(1/\alpha^{\frac{1}{\tau}} \eta) > 0$ and $\log(1/\alpha^{\frac{1}{\tau}} \sqrt{\delta}) > 0$. Therefore, after dividing both sides by $\log(1/\alpha^{\frac{1}{\tau}} \eta)$ and taking again the logarithm, we find that it is sufficient to have

$$\begin{aligned}
 t &\geq \left\lceil \frac{1}{\log(1+\tau)} \log \left(\frac{\log(1/\alpha^{\frac{1}{\tau}} \sqrt{\delta})}{\log(1/\alpha^{\frac{1}{\tau}} \eta)} \right) \right\rceil \\
 &= \left\lceil \frac{1}{\log(1+\tau)} \log \left(\frac{1 + \frac{\tau \log(1/\delta)}{2 \log(1/\alpha)}}{1 + \frac{\tau \log(1/\eta)}{\log(1/\alpha)}} \right) \right\rceil \\
 &= T_{\tau, \alpha, \delta}.
 \end{aligned}$$

\square

B.6. Proof of Theorem 2

We denote $N_1 := \frac{f(x_0) - f(x^*)}{\nu}$ and $\tilde{p} := \frac{p_0}{T+2}$, where $\bar{T} := N_1 + 1 + T_{\tau, \frac{3}{8}\delta}$. Recall that we pick $\varepsilon = 1/8$.

Our proof strategy proceeds as follows. In a first phase, we show that $f(x_{\text{nsk}}) - f(x) \leq -\nu$ until such a decrement cannot occur anymore, i.e., until $f(x_t) - f(x^*) < \nu$. Technical arguments for Phase 1 essentially follow from Lemma 2. Then, we enter a second phase where we observe a geometric decrease of the Newton decrement as described in Lemma 3.

We define

$$t := \inf \left\{ k \geq 0 \mid \tilde{\lambda}_f(x_k) \leq \eta \right\},$$

According to Lemma 8, we have $t \leq N_1$ with probability at least $1 - N_1 \tilde{p}$.

We turn to the analysis of Phase 2. We suppose that $T_f > t$ (i.e., the algorithm has not terminated during Phase 1), we define the additional number of iterations $J := \min\{T_{\tau, \frac{3}{8}\delta}, T_f - t - 1\}$, and we introduce the event

$$\mathcal{E}^{(2)} := \left\{ \mathcal{E}_{x_t, m_t, \varepsilon} \cap \bigcap_{j=0}^J \mathcal{E}_{x_{t+1+j}, m_{t+1+j}, \varepsilon \delta^{\frac{j}{2}}} \right\}.$$

Let us assume that $\mathcal{E}^{(2)}$ holds true, which happens with probability at least $1 - (2 + T_{\tau, \frac{3}{8}\delta})\tilde{p}$ according to Corollary 1. According to Lemma 9, we have for any $j = 0, \dots, J$ that $m_{t+1+j} = \bar{m}_2$ and,

$$\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1+j}) \leq (\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1}))^{(1+\tau)^j}.$$

Further, we have from Lemma 3 and Theorem 1 that $\lambda_f(x_{t+1}) \leq \frac{16}{25} \lambda_f(x_t) \leq \frac{\tilde{\lambda}_f(x_t)}{\sqrt{1-\varepsilon}} \leq \frac{\eta}{\sqrt{1-\varepsilon}} \leq \frac{1}{16}$. Hence, $\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1}) < 1/8$. As a consequence of Lemma 4, we must have that $\lambda_f(x_{t+1+j})^2 \leq \frac{3}{8}\delta$ for some $j \leq T_{\tau, \frac{3}{8}\delta}$, which further implies that

$$\tilde{\lambda}_f(x_{t+1+j})^2 \leq (1 + \varepsilon) \lambda_f(x_{t+1+j})^2 \leq \frac{3(1 + \varepsilon)}{8} \delta \leq \frac{3}{4} \delta.$$

The above inequality implies termination of the algorithm before the time $t + 1 + T_{\tau, \frac{3}{8}\delta}$. Using a union bound over $\{t \leq N_1\}$ and $\mathcal{E}^{(2)}$, we find that the algorithm terminates within $N_1 + 1 + T_{\tau, \frac{3}{8}\delta}$ iterations with probability at least $1 - (N_1 + 2 + T_{\tau, \frac{3}{8}\delta})\tilde{p}$.

It remains to guarantee that the algorithm returns a point \tilde{x} such that $f(\tilde{x}) - f(x^*) \leq \delta$. Note that the exit criterion guarantees that $\tilde{\lambda}_f(\tilde{x})^2 \leq \frac{3}{4}\delta$. Furthermore, the final sketch size \tilde{m} necessarily satisfies $\tilde{m} \geq \bar{m}_1$, so that, according to Theorem 1, we have with probability at least $1 - \tilde{p}$ that $\lambda_f(\tilde{x})^2 \leq \frac{1}{1-\varepsilon} \tilde{\lambda}_f(\tilde{x})^2 \leq \delta$. Self-concordance of f further implies that $f(\tilde{x}) - f(x^*) \leq \lambda_f(\tilde{x})^2 \leq \delta$.

In conclusion, we have shown that the algorithm returns a δ -accurate solution within $N_1 + 1 + T_{\tau, \frac{3}{8}\delta}$ iterations with probability at least $1 - (N_1 + 3 + T_{\tau, \frac{3}{8}\delta})\tilde{p} = 1 - p_0$. This concludes the proof. \square

B.6.1. COMPLEXITY GUARANTEES FOR THE SJLT

With the SJLT, consider the quadratic convergence case, i.e., $\tau = 1$. Let $p_0 > 0$ be a failure probability, and consider the sketch sizes

$$\bar{m}_1 \asymp \frac{\bar{d}_\mu^2 \log \log 1/\delta}{p_0}, \quad \bar{m}_2 \asymp \frac{1}{\delta} \frac{\bar{d}_\mu^2 \log \log 1/\delta}{p_0}.$$

We observe quadratic convergence with $T_f = \mathcal{O}(\log \log(\frac{1}{\delta}))$ iterations. Further, assuming that the sketching cost $\mathcal{O}(nd)$ dominates the cost $\mathcal{O}(\bar{m}^2 d)$ of solving the randomized Newton system, i.e., $n \gtrsim \frac{\bar{d}_\mu^4 \log(\log(1/\delta))^2}{\delta^2 p_0^2}$, then the total complexity results in

$$\mathcal{C} = \mathcal{O}(nd \log \log 1/\delta).$$

Similarly, we consider the linear convergence case, i.e., $\tau = 0$, and pick a failure probability $p_0 > 0$. Consider the sketch sizes

$$\bar{m}_1 \asymp \bar{m}_2 \asymp \frac{\bar{d}_\mu^2 \log 1/\delta}{p_0}.$$

We observe linear convergence with $T_f = \mathcal{O}(\log \frac{1}{\delta})$ iterations. Assuming again that the sketching cost dominates the cost of solving the randomized Newton system, i.e., $n \gtrsim \frac{\bar{d}_\mu^4 \log^2(1/\delta)}{p_0^2}$, we obtain the total time complexity

$$\mathcal{C} = \mathcal{O}(nd \log(1/\delta)).$$

□

B.7. Proof of Lemma 5

Let $S \in \mathbb{R}^{m \times n}$ be an embedding, and $C_S := H^{-1/2} H_S H^{-1/2}$. We use the notations $A := \nabla^2 f_0(x)^{1/2}$, and we let $A = U \Sigma V^\top$ be a thin SVD of A . Then, we have

$$\begin{aligned} C_S &= H^{-\frac{1}{2}} H_S H^{-\frac{1}{2}} = H^{-\frac{1}{2}} (H + (H_S - H)) H^{-\frac{1}{2}} \\ &= I_d + H^{-1/2} (H_S - H) H^{-1/2} \\ &= I_d + M^\top (U^\top S^\top S U - I_d) M, \end{aligned}$$

where $M := \Sigma V^\top H^{-1/2}$. According to (Cohen et al., 2015), it holds that $\|M^\top (U^\top S^\top S U - I_d) M\|_2 \leq \frac{\bar{d}_\mu}{2}$ (i.e., $\|C_S\|_2 \leq 1 + \frac{\bar{d}_\mu}{2}$) with probability at least $1 - p$, provided that $m \geq \Omega(\log^2(1/p))$ for a SRHT S , and, $m \geq \Omega(1/p)$ for a SJLT S .

Then, we use the fact that

$$\tilde{\lambda}_f(x)^2 = \langle H^{-1/2} \nabla f(x), H^{1/2} H_S^{-1} H^{1/2} H^{-1/2} \nabla f(x) \rangle \geq \frac{1}{\|C_S\|_2} \lambda_f(x)^2.$$

Conditional on $\|C_S\|_2 \leq 1 + \frac{\bar{d}_\mu}{2}$, it follows that

$$\lambda_f(x)^2 \leq \|C_S\|_2 \tilde{\lambda}_f(x)^2 \leq (1 + \frac{\bar{d}_\mu}{2}) \frac{\delta}{d} \leq \delta.$$

Using the self-concordance of f , we obtain that $f(x) - f(x^*) \leq \delta$. This concludes the proof. □

B.8. Proof of Theorem 3

We introduce the notations

$$\bar{T} = T_{\tau, \alpha(\tau, \varepsilon), \frac{\delta}{d}} + N_1, \quad \tilde{p} = \frac{p_0}{\bar{T}} \quad \text{and} \quad \varepsilon' = \varepsilon \left(\frac{\delta}{(1 + \varepsilon)d} \right)^{\tau/2}.$$

We consider \bar{m} a sketch size such that $\mathcal{E}_{x, \bar{m}, \varepsilon'}$ holds with probability at least $1 - \tilde{p}$, that is,

$$\begin{aligned} \bar{m} &= \Omega\left(\frac{d^\tau \bar{d}_\mu^2 \bar{T}}{p_0 \delta^\tau}\right) \quad \text{for the SJLT,} \\ \bar{m} &= \Omega\left(\frac{d^\tau}{\delta^\tau} \left(\bar{d}_\mu + \log\left(\frac{\bar{T} d^{\tau/2}}{p_0 \delta^{\tau/2}}\right) \log\left(\frac{\bar{d}_\mu \bar{T}}{p_0}\right)\right)\right) \quad \text{for the SRHT.} \end{aligned}$$

Phase 2. Let $t \geq 0$ be the first iteration such that $m_t \geq \bar{m}$, if any. Let $x \equiv x_{t+j}$ be an iterate after time t , for some $j \geq 0$. The sketch size is non-decreasing, whence $m \equiv m_{t+j} \geq \bar{m}$. We assume that $\mathcal{E}_{x, m, \varepsilon'}$ holds, and that the algorithm has not

yet terminated, i.e., $\tilde{\lambda}_f(x)^2 > \delta/d$. Note that $\varepsilon > \varepsilon'$, whence $\mathcal{E}_{x,m,\varepsilon}$ also holds. By Theorem 1, this implies in particular that $\tilde{\lambda}_f(x)^2 \leq (1+\varepsilon)\lambda_f(x)^2$, and we further obtain that $\lambda_f(x)^2 > \frac{\delta}{(1+\varepsilon)d}$, i.e.,

$$\varepsilon' < \varepsilon \lambda_f(x)^\tau.$$

There are two possible events.

- E_1 : Either $\tilde{\lambda}_f(x) > \eta$. Using the fact that $\mathcal{E}_{x,m,\varepsilon}$ holds, it follows from Lemma 2 that $f(x_{\text{nsk}}) - f(x) \leq -\nu$.
- E_2 : Or $\tilde{\lambda}_f(x) \leq \eta$. Using the facts that $\mathcal{E}_{x,m,\varepsilon'}$ holds and that $\varepsilon' < \varepsilon \lambda_f(x)^\tau$, it follows from Lemma 3 that $\lambda_f(x_{\text{nsk}}) \leq \alpha(\tau) (\lambda_f(x))^{1+\tau}$. Assuming further that the event $\mathcal{E}_{x_{\text{nsk}},m,\varepsilon'}$ holds, we have according to Lemma 6 that $\tilde{\lambda}_f(x_{\text{nsk}}) \leq \tilde{\lambda}_f(x) \leq \eta$ and then

$$\begin{aligned} \tilde{\lambda}_f(x_{\text{nsk}}) &\stackrel{(i)}{\leq} \sqrt{1+\varepsilon} \lambda_f(x_{\text{nsk}}) \leq \sqrt{1+\varepsilon} \alpha(\tau) (\lambda_f(x))^{1+\tau} \\ &\stackrel{(ii)}{\leq} \sqrt{1+\varepsilon} \alpha(\tau) (\tilde{\lambda}_f(x)/\sqrt{1-\varepsilon})^{1+\tau} \\ &= \alpha(\tau, \varepsilon) (\tilde{\lambda}_f(x))^{1+\tau}, \end{aligned}$$

where inequalities (i) and (ii) are immediate consequences of Theorem 1.

Hence, conditional on E_2 occurs once, then the event E_2 occurs K additional times in a row with probability at least $1 - K\tilde{p}$. According to Lemma 4, if $K \geq T_{\tau,\alpha(\tau,\varepsilon),\frac{\delta}{d}}$ then the algorithm terminates. On the other hand, the event E_1 can occur at most N_1 times.

In summary, conditional on $m_t \geq \bar{m}$, the algorithm must terminate within \bar{T} additional iterations with probability at least $1 - \bar{T}\tilde{p} = 1 - p_0$, and with final sketch size $m \leq 2\bar{m}$.

Phase 1. At each iteration, one of the following events must occur:

$$\begin{aligned} e_1 &:= \{\tilde{\lambda}_f(x) > \eta, f(x_{\text{nsk}}) - f(x) \leq -\nu\} \\ e_2 &:= \{\tilde{\lambda}_f(x) \leq \eta, \tilde{\lambda}_f(x_{\text{nsk}}) \leq \alpha(\tau, \varepsilon) (\tilde{\lambda}_f(x))^{1+\tau}\} \\ e_3 &:= \{m \leftarrow 2m\}. \end{aligned}$$

Fix any iteration $t \geq 0$, and suppose that the algorithm has not yet terminated. Consider the sequence of events $c_0, \dots, c_t \in \{e_1, e_2, e_3\}$ up to time t . According to Lemma 4, any subsequence of $\{c_j\}_{j=0}^t$ which contains only the event e_2 would result in termination of Algorithm 2 if its length is greater or equal to $T_{\tau,\alpha(\tau,\varepsilon),\delta/d} + 1$. Consequently, any such subsequence must have length smaller or equal to $T_{\tau,\alpha(\tau,\varepsilon),\delta/d}$. Between two consecutive longest subsequences containing only e_2 , either e_1 or e_3 occur. The event e_1 occurs at most N_1 times. By assumption on the choice of m_0 , once e_3 has occurred at least $\mathcal{O}(\log(\bar{d}_\mu))$ times then the sketch size is greater than \bar{m} . Consequently, there are at most $T_1 := \mathcal{O}((N_1 + \log(\bar{d}_\mu))T_{\tau,\alpha(\tau,\varepsilon),\delta/d})$ iterations before reaching a sketch size m such that $m \geq \bar{m}$ without termination. In the latter case, we enter Phase 2.

Combining Phase 1 and Phase 2. Combining the two above results, we obtain with probability at least $1 - p_0$ that Algorithm 2 terminates with a final sketch size m smaller than $2\bar{m}$ and within a number of iterations T scaling as

$$T = T_1 + T_2 = \mathcal{O}((N_1 + \log(\bar{d}_\mu))T_{\tau,\alpha(\tau,\varepsilon),\delta/d}) = \mathcal{O}(\log(\bar{d}_\mu) T_{\tau,\alpha(\tau,\varepsilon),\delta/d}),$$

where the last equality holds by treating N_1 as $\mathcal{O}(1)$.

Total complexity. The worst-case complexity per iteration is given as follows.

- (1) For a SJLT S , the sketching cost is at most $\mathcal{O}(nd)$ at each iteration, and forming and solving the linear system $H_S v_{\text{nsk}} = -\nabla f(x)$ with a direct method using the Woodbury identity takes time $\mathcal{O}(\bar{m}^2 d)$. Multiplying by the number of iterations, we obtain the total time complexity

$$\bar{C} = \mathcal{O} \left(\left(nd + \frac{\bar{d}_\mu^4 d^{2\tau+1} T_{\tau,\alpha(\tau,\varepsilon),\delta/d}^2}{\delta^{2\tau} p_0^2} \right) \log(\bar{d}_\mu) T_{\tau,\alpha(\tau,\varepsilon),\delta/d} \right).$$

For $\tau \approx 1$, we have that $T_{\tau, \alpha(\tau, \varepsilon), \delta/d} = \mathcal{O}(\log(\log(d/\delta)))$. For $n \gtrsim \frac{\bar{d}_\mu^4 d^2 \log(\log(d/\delta))^2}{\delta^2 p_0^2}$, the memory and time complexities simplify to

$$\bar{m} = \Omega\left(\frac{d \bar{d}_\mu^2 \log(\log(d/\delta))}{p_0 \delta}\right), \quad \bar{c} = \mathcal{O}(nd \log(\bar{d}_\mu) \log(\log(d/\delta))).$$

For $\tau = 0$, we have $T_{\tau, \alpha(\tau, \varepsilon), \delta/d} = \mathcal{O}(\log(d/\delta))$. For $n \gtrsim \frac{\bar{d}_\mu^4 \log(d/\delta)^2}{p_0^2}$, the memory and time complexities simplify to

$$\bar{m} = \Omega\left(\frac{\bar{d}_\mu^2 \log(d/\delta)}{p_0}\right), \quad \bar{c} = \mathcal{O}(nd \log(\bar{d}_\mu) \log(d/\delta)).$$

(2) We assume for simplicity that $\bar{d}_\mu \gtrsim \log^2(\log(d/\delta))$. For the SRHT, the sketching cost is $\mathcal{O}(nd \log \bar{m})$, whereas forming and solving the Newton linear system takes time $\mathcal{O}(\bar{m}^2 d)$. Thus, the total complexity is given by

$$\bar{c} = \mathcal{O}((nd \log \bar{m} + d \bar{m}^2) \log(\bar{d}_\mu) T_{\tau, \alpha(\tau, \varepsilon), \delta/d}).$$

For $\tau \approx 1$, we have $T_{\tau, \alpha(\tau, \varepsilon), \delta/d} = \mathcal{O}(\log(\log(d/\delta)))$. Picking $p_0 \asymp 1/\bar{d}_\mu$, we obtain the memory complexity

$$\bar{m} \asymp \frac{d}{\delta} (\bar{d}_\mu + \log(d/\delta) \log(\bar{d}_\mu)).$$

Consequently, $\log \bar{m} \lesssim \log(d/\delta)$ and $\bar{m}^2 \lesssim \frac{d^2}{\delta^2} (\bar{d}_\mu^2 + \log^2(d/\delta) \log^2(\bar{d}_\mu))$. Hence, provided that $n \gtrsim \frac{d^2 \bar{d}_\mu^2}{\delta^2}$, we obtain

$$\bar{c} = \mathcal{O}(nd \log(d/\delta) \log(\bar{d}_\mu) \log(\log(d/\delta))).$$

For $\tau = 0$, we have $T_{\tau, \alpha(\tau, \varepsilon), \delta/d} = \mathcal{O}(\log(d/\delta))$. Picking $p_0 \asymp 1/\bar{d}_\mu$, we obtain the memory complexity

$$\bar{m} \asymp \bar{d}_\mu.$$

Consequently, $\log \bar{m} \lesssim \log(\bar{d}_\mu)$ and $\bar{m}^2 \lesssim \bar{d}_\mu^2$. Assuming that $n \gtrsim \bar{d}_\mu^2 / \log(\bar{d}_\mu)$, the total time complexity is

$$\bar{c} = \mathcal{O}(nd \log(\bar{d}_\mu)^2 \log(d/\delta)).$$

This concludes the proof. \square

C. Auxiliary results

Lemma 6. *Let $x \in \text{dom } f$ and $\varepsilon \in (0, 1/4)$. Suppose that the event $\mathcal{E}_{x, m, \varepsilon} \cap \mathcal{E}_{x_{\text{nsk}}, m_{\text{nsk}}, \varepsilon}$ holds, and that $\tilde{\lambda}_f(x) \leq \eta$. Then, we have that*

$$\tilde{\lambda}_f(x_{\text{nsk}}) \leq \tilde{\lambda}_f(x) \leq \eta. \quad (43)$$

Proof. By assumption, the event $\mathcal{E}_{x_{\text{nsk}}, m_{\text{nsk}}, \varepsilon}$ holds. It follows from Theorem 1 that $\tilde{\lambda}_f(x_{\text{nsk}}) \leq \sqrt{1 + \varepsilon} \lambda_f(x_{\text{nsk}})$. We have by assumption that $\mathcal{E}_{x, m, \varepsilon}$ holds and that $\tilde{\lambda}_f(x) \leq \eta$. As a consequence of Lemma 3, we have $\tilde{\lambda}_f(x) \leq \frac{16}{25} \lambda_f(x)$. As a consequence of Theorem 1, we have $\lambda_f(x) \leq \frac{1}{\sqrt{1 - \varepsilon}} \tilde{\lambda}_f(x)$. Combining these bounds together, we obtain that

$$\tilde{\lambda}_f(x_{\text{nsk}}) \leq \sqrt{\frac{1 + \varepsilon}{1 - \varepsilon}} \frac{16}{25} \tilde{\lambda}_f(x).$$

Finally, using that $\varepsilon \in (0, 1/4)$, we get that $\sqrt{\frac{1 + \varepsilon}{1 - \varepsilon}} \frac{16}{25} \leq 1$, whence,

$$\tilde{\lambda}_f(x_{\text{nsk}}) \leq \tilde{\lambda}_f(x) \leq \eta.$$

\square

Lemma 7. For $\varepsilon \in (0, 1)$, it holds that

$$\eta \leq \frac{1-\varepsilon}{1+\varepsilon} \frac{1}{16} \leq \frac{1}{16}. \quad (44)$$

Proof. Set $\gamma = \left(\frac{1+\varepsilon}{1-\varepsilon}\right)^2$. We aim to show that $\eta \sqrt{\gamma} \leq 1/16$. Plugging-in the definition of η and using that $a \geq 0$, we have $\eta \sqrt{\gamma} = \frac{1}{8} \frac{1-\frac{\gamma}{2}-a}{\gamma} \leq \frac{1}{8} \frac{1-\frac{\gamma}{2}}{\gamma}$. Since $\varphi(\gamma) := \frac{1}{8} \frac{1-\frac{\gamma}{2}}{\gamma}$ is monotone decreasing and since $\gamma \geq 1$, we obtain that $\eta \sqrt{\gamma} \leq \varphi(1)$, i.e., $\eta \sqrt{\gamma} \leq \frac{1}{16}$. \square

C.1. Technical lemmas for the proof of Theorem 2

Lemma 8 (Phase 1). *It holds that*

$$t \leq N_1, \quad \text{with probability at least } 1 - N_1 \tilde{p}.$$

Proof. Let $j < t$ be any iteration before t_1 . Note by construction of Algorithm 1 that $m_j = \bar{m}_1$. Assuming that the event $\mathcal{E}_{x_j, m_j, \varepsilon}$ holds true, it follows from Lemma 2 that we observe the decrement $f(x_{\text{nsk}}) - f(x_j) \leq -\nu$. Consequently, under the event $\mathcal{E}^{(1)} := \bigcap_{j=0}^{t-1} \mathcal{E}_{x_j, m_j, \varepsilon}$, we obtain that

$$f(x^*) - f(x_0) \leq f(x_t) - f(x_0) = \sum_{j=0}^{t-1} f(x_{j+1}) - f(x_j) \leq -t\nu.$$

Hence, under $\mathcal{E}^{(1)}$, we must have $t \leq \frac{f(x_0) - f(x^*)}{\nu}$, i.e., $t \leq N_1$. According to Lemma 1 and the choice of \bar{m}_1 , each event $\mathcal{E}_{x_j, m_j, \varepsilon}$ holds with probability at least $1 - \tilde{p}$. Using a union bound, the event $\mathcal{E}^{(1)}$ holds with probability at least $1 - N_1 \tilde{p}$. \square

Lemma 9 (Phase 2). *Under the assumption that $\mathcal{E}^{(2)}$ holds, we have for any $j = 0, \dots, J$ that*

$$\begin{cases} m_{t+1+j} = \bar{m}_2, \\ \tilde{\lambda}_f(x_{t+1+j}) \leq \eta, \\ \left(\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1+j+1})\right) \leq \left(\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1+j})\right)^{1+\tau}. \end{cases}$$

Proof. We prove this claim by induction. We start with $j = 0$. By definition of the time t , we have $\tilde{\lambda}_f(x_t) \leq \eta$. Therefore, by construction of Algorithm 1, we have $m_{t+1} = \bar{m}_2$. From Lemma 6 and under $\mathcal{E}^{(2)}$, we get that $\tilde{\lambda}_f(x_{t+1}) \leq \tilde{\lambda}_f(x_t) \leq \eta$. Furthermore, before termination, we have that $\tilde{\lambda}_f(x_{t+1})^2 > \frac{3}{4}\delta$. It follows from Theorem 1 that

$$\lambda_f(x_{t+1})^2 \geq \frac{1}{1+\varepsilon} \tilde{\lambda}_f(x_{t+1})^2 > \frac{3}{4(1+\varepsilon)} \delta = \frac{2}{3} \delta,$$

and this implies in particular that $\varepsilon \delta^{\tau/2} \leq \varepsilon \left(\frac{3}{2}\right)^{\tau/2} \lambda_f(x_{t+1})^\tau \leq 2\varepsilon \lambda_f(x_{t+1})^\tau$. Consequently, the hypotheses of Lemma 3 are verified and we have $\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+2}) \leq \left(\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1})\right)^{1+\tau}$.

Now, we prove the induction hypothesis for any $j = 1, \dots, J$, assuming that it holds for $j - 1$. Since $\tilde{\lambda}_f(x_{t+1+j-1}) \leq \eta$, it follows by construction of Algorithm 1 that $m_{t+1+j} = \bar{m}_2$. From Lemma 6 and under $\mathcal{E}^{(2)}$, we get that $\tilde{\lambda}_f(x_{t+1+j}) \leq \tilde{\lambda}_f(x_{t+1+j-1}) \leq \eta$. Furthermore, before termination, we have $\tilde{\lambda}_f(x_{t+1+j})^2 > \frac{3}{4}\delta$. It follows from Theorem 1 that

$$\lambda_f(x_{t+1+j})^2 \geq \frac{1}{1+\varepsilon} \tilde{\lambda}_f(x_{t+1+j})^2 > \frac{3}{4(1+\varepsilon)} \delta = \frac{2}{3} \delta,$$

and this implies in particular that $\varepsilon \delta^{\tau/2} \leq \varepsilon \left(\frac{3}{2}\right)^{\tau/2} \lambda_f(x_{t+1+j})^\tau \leq 2\varepsilon \lambda_f(x_{t+1+j})^\tau$. Consequently, the hypotheses of Lemma 3 are verified and we have $\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1+j+1}) \leq \left(\alpha(\tau)^{\frac{1}{\tau}} \lambda_f(x_{t+1+j})\right)^{1+\tau}$. \square

Corollary 1. *The event $\mathcal{E}^{(2)}$ holds true with probability at least $1 - (2 + T_{\tau, \frac{3}{8}\delta})\tilde{p}$.*

Adaptive Newton Sketch

Proof. Recall that $m_t = \bar{m}_1$ by definition of the time t . According to Lemma 9, if $\mathcal{E}^{(2)}$ holds true, then $m_{t+1+j} = \bar{m}_2$ for $j = 0, \dots, J$. From Lemma 1, we have that $\mathbb{P}(\mathcal{E}_{x_t, \bar{m}_1, \varepsilon}) \geq 1 - \tilde{p}$ and $\mathbb{P}(\mathcal{E}_{x_{t+1+j}, \bar{m}_2, \varepsilon \delta^{\tau/2}}) \geq 1 - \tilde{p}$. We obtain by a union bound that $\mathbb{P}(\mathcal{E}^{(2)}) \geq 1 - (2 + T_{\tau, \frac{3}{8}\delta})\tilde{p}$. \square